

HiPEAC info 57

APRIL 2019

**Computing
Systems Week
Edinburgh**

High performance, HiPEAC style

Future chips: The European Processor Initiative

How HPC is transforming businesses



Welcome to Edinburgh



Toni Collis on the power of HPC



The European Processor Initiative

<p>3 Welcome <i>Koen De Bosschere</i></p> <p>4 News</p> <p>10 HiPEAC voices 'If all HPC programmers are male, we are limiting the scientific discovery of the human race' <i>Toni Collis</i></p> <p>14 High-performance computing special Homegrown high performance: The European Processor Initiative (EPI) <i>Mario Kovač</i></p> <p>18 High-performance computing special Performance, performance, performance: HPC, HiPEAC style <i>Estela Suarez, Eva Gellner, Konstantinos Nikas, Andreas Koch, Lukas Sommer, Jaco Hofmann, Martin Schulz, Tapasya Patki, Siddhartha Jana and Masaaki Kondo</i></p> <p>20 High-performance computing special Top of the FLOPS: Scaling up Europe's performance <i>Mike Ashworth, Georgios Goumas and Peter Hopton</i></p> <p>22 SME snapshot Campera: Simplifying FPGA complexity <i>Calliope-Louisa Sotiropoulou</i></p> <p>23 Industry focus To infinity and beyond: Preparing the Intel® Movidius™ Myriad 2 VPU for space <i>Aubrey Dunne</i></p> <p>24 Industry focus Breaking through the cloud I/O bottleneck with SUNLIGHT <i>Julian Chesterfield, Michail Flouris and Stelios Louloudakis</i></p> <p>26 HPC and innovation Taking HPC from the lab to the market with Eurolab4HPC <i>Per Stenström, Katrien Van Impe, Josep de la Puente, Jean-Thomas Acquaviva and Julian Kunkel</i></p>	<p>28 HPC and innovation Shaping up European companies with HPC <i>Chris Johnson</i></p> <p>29 Innovation Europe PRACE enters its sixth implementation phase <i>Stelios Erotokritou, Marjolein Oorsprong and Oriol Pineda</i></p> <p>31 Innovation Europe EVOLVE: Fusing HPC with cloud to extract maximum big-data value <i>Angelos Bilas, Dimitrios Soudris and Jean-Thomas Acquaviva</i></p> <p>32 Innovation Europe Energy, security and time awarded first place with TeamPlay <i>Emad Samuel Malki Ebeid</i></p> <p>33 Innovation Europe Exploiting each chip to its full potential: How UniServer overcomes energy scaling limits in commodity servers <i>Georgios Karakonstantis</i></p> <p>34 Technology opinion Is the peer-review system broken? <i>Ben Juurlink</i></p> <p>35 Technology opinion Standardized benchmarking in research papers <i>Vincent Hindriksen</i></p> <p>36 Peac performance ReFiRe: efficient deployment of Remote Fine-grained Reconfigurable accelerators <i>Dimitrios Pnevmatikatos</i></p> <p>38 HiPEAC futures The HiPEAC internship programme Career talk: Patience Masoso, Zimbabwe Centre for High-Performance Computing Mastering high-performance computer systems RISC-ing everything: Developing a Verilog core in a HiPEAC internship Three-minute thesis: Simplifying parallel prefix operations on heterogeneous platforms</p>
---	---



High-performance computing special



Space training for the Intel® Movidius™ Myriad 2 VPU



Career talk: Patience Masoso, ZCHPC



HiPEAC is the European network on high performance and embedded architecture and compilation.



hipecac.net



@hipecac



hipecac.net/linkedin



HiPEAC has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 779656.

Cover image: Visualization of the dynamical nature of Perovskite. Simulation carried out using data from ENEA (Italy), as part of the Energy Oriented Centre of Excellence (EoCoE) funded by Horizon 2020

Credit: Guillermo Marín / Fernando Cucchiatti
Back cover photo: Stefano Cherubin / Eneko Illarramendi

Design: www.magelaan.be

Editor: Madeleine Gray

Email: communication@hipecac.net

This spring, the HiPEAC network meets in Edinburgh for our Computing Systems Week on innovation and high-performance computing (HPC). The last time we organized an event in Edinburgh was nine years ago. That Computing Systems Week is for me still the most influential in the history of HiPEAC.

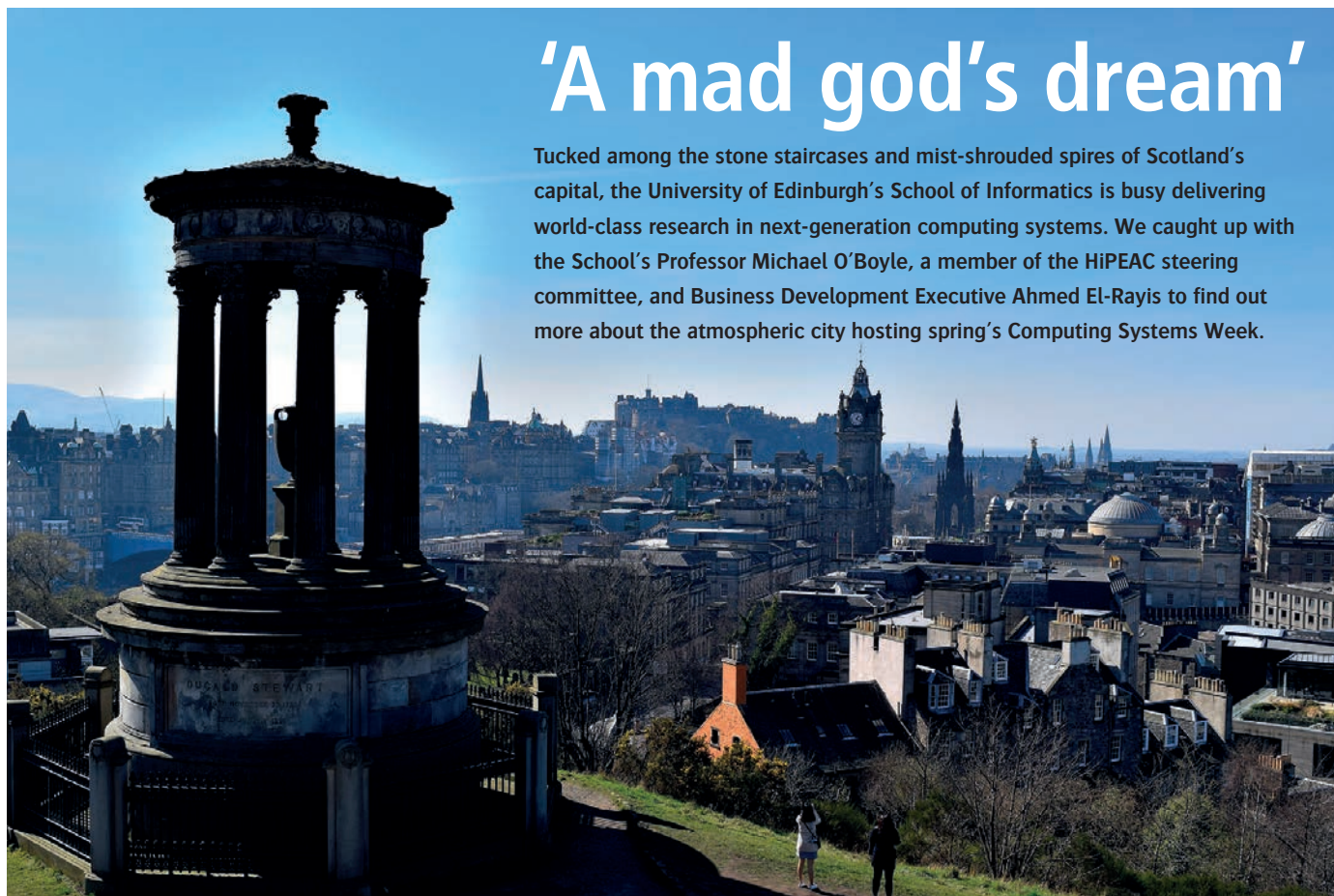
First of all, during that event we decided to move to a journal-first publication model for the HiPEAC conference, and to transform the conference from a publication event into a networking event. In 2012, we organized the first journal-first conference in Paris, and its attendance surpassed all our expectations. Instead of 200 attendees, we attracted more than 500. The hotel could barely host so many people. The journal we started collaborating with saw its number of submissions more than quadruple. It was clear that this was a transition the community had been waiting for. Since then, many conferences have been thinking about their publication model.

At the same Computing Systems Week, I was inspired by the keynote by Colin Adams about a project to make computer science students more entrepreneurial. In 2012, we officially started a student-entrepreneurship programme at Ghent University. After seven years, it has grown to be the largest such programme in Belgium, generating around 200 student enterprises every year in the city of Ghent. In 2018, we managed to convince the Belgian government to change the tax law to remove the fiscal discrimination between a student with a job and a self-employed student. Today, the higher education institutions in Ghent have a total of 10 business coaches working with students seeking advice for an entrepreneurial idea. The Flemish government decided this spring to make basic entrepreneurial education compulsory in all study programmes of higher education. I believe Flanders is the first region to pass such a law.

Sometimes, people ask us about the impact of HiPEAC. Well, it transformed the life of higher education students in Flanders, and it is making Flanders a more entrepreneurial region. I cannot say this was an impact we had foreseen in the project description. The Computing Systems Week in 2010 was part of HiPEAC2; today, we are halfway through HiPEAC5. It also shows that sometimes we should be patient to see the real impact of a project.

I hope this edition of the Computing Systems Week will be as inspiring as the one in 2010. In the meantime, this issue gives an update on high-performance innovation within the community, including the European Processor Initiative, European efforts to get to exascale and how high-performance computing is making businesses more competitive.

Koen De Bosschere, HiPEAC coordinator



'A mad god's dream'

Tucked among the stone staircases and mist-shrouded spires of Scotland's capital, the University of Edinburgh's School of Informatics is busy delivering world-class research in next-generation computing systems. We caught up with the School's Professor Michael O'Boyle, a member of the HiPEAC steering committee, and Business Development Executive Ahmed El-Rayis to find out more about the atmospheric city hosting spring's Computing Systems Week.

Photo credit: Harry J Burgess, Pixabay

Fàilte gu Dùn Èideann – Welcome to Edinburgh



What makes Edinburgh a special place for science, particularly computer science?

Founded in 1582, Edinburgh has had 23 Nobel Prize winners, including Richard Henderson for chemistry in 2017 and Peter Higgs, best known for the Higgs Boson, for physics in 2013. It has also been a leader in genetics with Dolly the Sheep, the first cloned mammal in 1996.

The School of Informatics has 450 academic and research staff and over 850 students, making it the largest in the UK and one of the largest in Europe. In the most recent research excellence evaluation, it produced more world-leading and internationally excellent research than any other computer science department in the UK.

Since 2018, we've been co-located with EPCC, who host ARCHER, the UK's primary academic research supercomputer. EPCC also develops and hosts the World-Class Data Infrastructure (WCDI), which underpins the Data-Driven Innovation programme of the Edinburgh and South East Scotland City Region Deal.

How does the innovation ecosystem support the transfer of computing technologies to the marketplace?

Edinburgh in general and Edinburgh University in particular support entrepreneurship on multiple levels. The ecosystem in Edinburgh includes multiple accelerators, incubators and various competitions, such as Converge Chal-

Mini phrase book

Scottish Gaelic was once spoken throughout Scotland. Here are some essential phrases to impress locals during your stay:

English

Scottish Gaelic

Hardware / software co-design is the way forward.

'S e an rathad air adhart co-dhealbhachadh bathar-cruaidh / bathar-bog

That is one powerful supercomputer.

'S e sin am for-choimpiutair cumhachdach

Did someone say 'whisky time'?

An tuir cuideigin 'uisge-beatha'?

Many thanks to Dr William E Lamb at the University of Edinburgh's School of Literatures, Languages and Cultures for the translation

lenge and Edge Scotland, that encourage spinouts and start-ups. Moreover, development agencies such as Scottish Enterprise support both the creation and the growth of small/medium enterprises (SMEs).

Edinburgh has been known for successful SMEs such as Calvatec, Skyscanner, Edesix (acquired by Motorola), VLSI Vision (acquired by STMicroelectronics), Wolfson (acquired by Cirrus Logic), FreeAgent (acquired by Royal Bank of Scotland), Veropath (acquired by Calero Software), SeeByte (acquired by Bluefin Robotics), Odos Imaging (acquired by Rockwell Automation), Camel Audio (acquired by Apple) and so on.

One special thing that differentiates the School of Informatics from other institutions is that it has its own commercialization team, which enables strong collaborations with industry but most importantly supports the creation of companies in relation to new innovations. The University of Edinburgh has been ranked top among UK Universities for 15 years in having the highest number of spinoffs and start-ups – an impressive 270.

Tell us a couple of things we probably don't know about Edinburgh.

Edinburgh is by the sea – there's a beach at Portobello where the water is delightful at this time of year. *[Editor's note: HiPEAC assumes no liability for medical insurance claims as a result of contracting hypothermia.]* There are mountains behind the city: the Pentlands Skyline walk is a 25 km circuit with 2,000 metres of ascent/descent – take the number

15 bus. We have a volcano in the city centre, visible from the Informatics Forum roof terrace. Known as Arthur's Seat, it's no longer smoking, and you can access it from the end of the Royal Mile. Oh, and we voted to stay in Europe! 74 per cent voted 'Remain', one of the highest numbers in the UK.

Any must-sees for visitors to the city?

The top selections have to be the Royal Mile (Castle to Palace), Calton Hill and Arthur's Seat. For shopping, head to Princes Street and George Street. For those who've been here before, the Water of Leith walk from Stockbridge to Dean Village is special, with stop-offs for good food in Stockbridge plus the excellent cafés at the Dean Gallery and Scottish Gallery of Modern Art. Finally, the Shore in Leith offers a glimpse into a completely different side to Edinburgh.

Where's best to go for a whisky after a hard day at Computing Systems Week?

The Doctor's and the Pear Tree are popular pubs within 100 metres of the Informatics Forum. In fact, there are hundreds of pubs within walking distance; Cowgate and Rose Street have a particularly high concentration, making for a convenient pub crawl. My favourites include the Café Royal on West Register Street and the Cask and Barrel on Broughton Street, along with the Jazz Bar on Chambers Street and the Devil's Advocate in the Old Town – if you can find it! Another one worth mentioning is the Canny Man in Morningside, which is unusual and has a wide whisky selection.

A few famous thinkers with links to Edinburgh



John Napier (1550 - 1617), above, mathematician, famous for the invention of logarithms

David Hume (1711 - 1776), philosopher
Joseph Black (1728 - 1799), physicist and chemist, first to isolate carbon dioxide

Adam Smith (1723 - 1790), economist, author of *The Wealth of Nations*

Mary Fairfax Somerville (1780 - 1872), science writer and polymath

Alexander Graham Bell (1847 - 1922), telephone pioneer

Charles Darwin (1809 - 1882), biologist, studied at the University of Edinburgh

Peter Higgs (born 1929), theoretical physicist

Lesley Jane Yellowlees (born 1953), inorganic chemist, first female president of the Royal Society of Chemistry

Bill Buchanan (born 1961), computer scientist

J. K. Rowling (born 1965), Harry Potter author



Arthur's Seat is visible from the Informatics Forum roof terrace



The Shore in Leith shows a different side of the city

Computing experts come together in Valencia for HiPEAC19



Over 540 delegates representing 263 institutions in 41 countries gathered in Valencia on 21-23 January 2019 for the fourteenth edition of the HiPEAC conference. With 26 workshops and 16 tutorials in addition to the ACM TACO paper track, the conference once again showcased the latest in European research on computing systems, with topics including deep learning, open-source hardware and quantum computing.

Valencia's famous hospitality made the Spanish city an elegant choice of location, as HiPEAC19 General Chair José Duato (Universitat Politècnica de València) explained: 'Valencia is the perfect blend of old and new. On the one hand, the historic buildings and excellent local restaurants; on the other, the unique and futuristic architecture, as well as the cutting-edge research being carried out at the city's two universities.'

The conference's three keynote speakers focused on three areas of major interest to the community. Monica Lam (Stanford University) explored the threat to the open web posed by virtual assistants and the paradigm shift associated with natural language processing; Alberto Sangiovanni Vincentelli (University of California, Berkeley) described his journey to founding multibillion dollar technology companies; while longstanding HiPEAC member Koen Bertels (Delft University of Technology) gave an eye-opening presentation on quantum computing.

'It's always a pleasure to come to the HiPEAC conference, because it's a place where you find a lot of knowledge,' commented Sandro D'Elia of the European Commission. 'It's not common to hear someone talking about quantum computing who really understands what is behind it, as I heard in the excellent keynote.'



In between attending events on the packed programme, delegates mingled in the conference's largest exhibition to date, where major multinational companies such as DeepMind, Arm, Atos and Xilinx mingled with European start-ups and scale-ups, as well as European Union-funded projects. In total, 60 companies were represented at the conference, with sponsoring organizations given the chance to present during the industrial session.

Once again, HiPEAC Jobs played a prominent role in the conference, with job and internship opportunities from around Europe featured on the HiPEAC Jobs wall. Following the success of last year's event, the science, technology, engineering and mathematics (STEM) student day introduced students to the HiPEAC community, allowing them to attend Tuesday's keynote talk, talk to companies attending and find out about the latest computing systems research.

The HiPEAC conference would not be the same without the generosity of its sponsors, and this year's event was no exception. A full list of sponsoring organizations is available on the HiPEAC19 website.

hipeac.net/2019

HiPEAC19 Google Photos album bit.ly/HiPEAC19_photos

HiPEAC19 YouTube playlist – including full keynote talks

bit.ly/HiPEAC19_videos

Article by *The Next Platform* on HiPEAC19

bit.ly/HiPEAC19_The_Next_Platform



Computing, architecture, innovation and more at the ACACES summer school



Photo credit: Magnus Sjölander

Taking place on 14-20 July in Fuggi, the fifteenth edition of HiPEAC's summer school, ACACES, offers an excellent, varied programme delivered by world-class experts in computing systems. Following the success of last year's event, we have once again teamed up with Eurolab4HPC to offer a high-performance computing (HPC) track, and with TETRAMAX to offer an entrepreneurial track.

This year's topics include:

- accelerated machine intelligence from edge to cloud
- fundamental limits in energy consumption for information technology
- reconfigurable multiprocessor systems-on-chip
- die stacking
- distributed memory
- intellectual rights protection

...and much more.

The deadline for applications is 31 May 2019.

Further information:

acaces.hipeac.net/2019



COEMS workshop at HiPEAC 2019

Volker Stolz and Svetlana Jakšić (Western Norway University of Applied Sciences)

The COEMS project workshop took place at HiPEAC conference in Valencia on 21 January 2019. In the COEMS (Continuous Observation of Embedded Multicore Systems) project, consortium members (University of Lübeck, Accemic, Airbus, Thales and Western Norway University of Applied Sciences) have joined forces to building a novel platform for online monitoring of multicore systems, which gives insight into the system's behaviour without affecting it.

Modern multicore processors provide highly compressed tracing information over a separate tracing port with little overhead. The COEMS system first reconstructs the sequence of instructions executed by the processor. This sequence is then analysed online by the compiled specification written in a specification language TeSSLa on a reconfigurable monitoring unit. To cope with the amount of tracing data generated by modern processors, we implemented the COEMS system in hardware using a field-programmable gate array (FPGA).

At the workshop, we presented the current state of project findings and developments. Topics covered included specification queries, data race detection, software-based tracing in the railway domain, runtime verification and real-time monitoring for correctness and robustness. Talks given by project partners and external domain experts Dr César Sánchez from the IMDEA Software Institute in Spain and Dr Stefan Jakšić from the Institute of Technology in Austria formed the basis for fruitful discussions. We gave the first demonstration of the COEMS hardware-based, non-intrusive, continuous observation and discussed performance and optimizations.

We are thankful that the workshop was organized as part of the HiPEAC conference, as this gave us a valuable opportunity to get feedback and communicate with experts and interested parties.

coems.eu

COEMS has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 732016

Long-term impact of European technology transfer in the ICT domain

Evidence from the TETRACOM and TETRAMAX projects

Rainer Leupers (RWTH Aachen) and Katrien Van Impe (Ghent University)



Academia-to-industry technology transfer is on the agenda of most European Union (EU)-funded projects, but it mostly becomes a reality only after projects have finished, by which time no infrastructure or resources for impact measurement are left. There is also little data available about the long-term impact. The full journey from a novel research idea to an innovative product may easily take more than five years – considerably longer than the duration of typical EU-funded projects.

The TETRACOM and TETRAMAX project series offered a unique opportunity to analyse long-term impact based on a comprehensive database of 50 Technology Transfer Experiments (TTX) in the domain of embedded information and communication technology (ICT). TETRACOM, an FP7 Coordination and Support Action, co-funded 50 TTX during 2013-2016, each focused on transferring a specific piece of hardware or software intellectual property (IP) from academic research to industry, prioritizing European small/medium enterprises (SMEs).

TETRACOM's total TTX impact was measured using indicators such as products improved, new jobs created or revenue increase thanks to the TTX, based on questionnaires completed by the client companies and research institutions. TETRAMAX, an ongoing H2020 Innovation Action, then repeated the same measurements in 2018, thus tracing the impact of the TTX around two to four years after their conclusion, or three to five years after their kick-off, respectively.

Long-term impact highlights include:

- contributions to standards more than tripled
- the number of new jobs created also more than tripled
- the number of open source-packages went up by two and a half times
- the number of new or improved products went up by 40%
- the TTX contributed to total cost savings of around €2 million and increased revenue of €4 million

A dramatic increase (over five times) in the number of TTX-related publications was also observed.

These figures are conservative, since confidentiality issues prevented some TTX clients from disclosing indicators. Besides enabling enhanced products for established companies, the TTX also helped 10 European start-ups, with more still in formation. Last but not least, in total €24 million in venture capital investment so far were at least partially attributed to the TTX.

Our conclusions are:

- EU-funded projects with a clear technology transfer focus can have significant, tangible economic impact, multiplying European taxpayers' investment (just €2 million in the case of TETRACOM).
- Project impact tends to grow over-proportionally. It is therefore more useful to capture and analyse impact metrics not primarily during the project's duration or immediately after its completion, but with more long term and systematic methods.

tetracom.eu

tetramax.eu

Book: lot for Smart Grids – Design challenges and paradigms

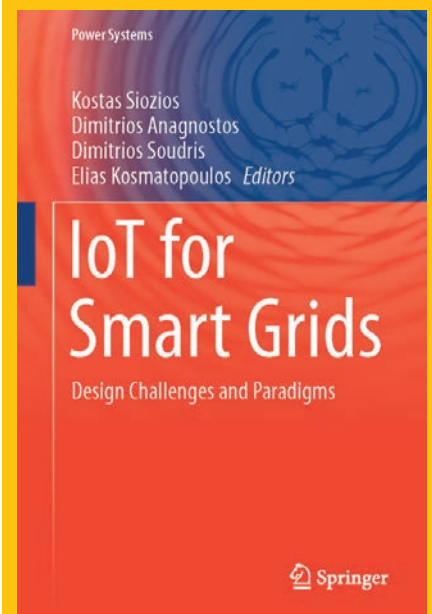
Editors: Kostas Siozios, Dimitrios Anagnostos, Dimitrios Soudris and Elias Kosmatopoulos

Recently, the convergence of emerging embedded computing, information technology, and distributed control has become a key enabler for applications and services with enormous societal impact and economic benefit. Energy systems, and more specifically smart grids, are an application domain where internet of things (IoT) technology is widely deployed.

This book studies topics related to the design and orchestration of IoT systems in order to provide higher flexibility to power generation and transmission systems. By analysing hardware and software solutions for decentralized energy systems, it is feasible to tackle the increasingly complex task of grid management.

Further information:

bit.ly/IoT_SG_Springer



HiPEAC Tech Transfer award winners 2018



In December, the winners of the HiPEAC Tech Transfer Awards 2018 were announced. These awards recognize successful examples of technology transfer, which covers technology licensing, providing dedicated services or creating a new company, for example.

For the purposes of the awards, technology transfer is defined as a contractually documented joint- or privately funded academia-industry project or technology licence agreement, with the goal of bringing a concrete research result into industrial practice. All applications are evaluated by an internal technology transfer committee.

For this edition, the nine projects selected demonstrate how the HiPEAC community is delivering real innovation for industry and thereby society. Spanning data fusion strategies, machine vision, cybersecurity, low-energy hardware, high-performance computing (HPC) simulation and much more, the awards highlight the breadth of expertise in the computing systems network.

The 2018 Tech Transfer Awards winners are as follows:

- **José M. Cecilia**, Universidad Católica de Murcia (UCAM): Multisensor data fusion strategies for the prevention of water-related disasters in El Salvador
- **Francesco Conti**, University of Bologna: Ultra-Low Energy Hardware Convolution Engine for GAP-8 IoT Application Processor
- **George Dan Mois**, Technical University of Cluj-Napoca: Thermal Printer, Bluetooth Low Energy and microSD Data Logger
- **Bjorn De Sutter**, Ghent University: Tightly-coupled self-debugging software protection
- **Oscar Deniz Suarez**, University of Castilla-La Mancha (UCLM): Eyes of Things
- **Magnus Jahre**, Norwegian University of Science and Technology: Non-Intrusive Power Monitoring for Embedded Systems
- **Diego R. Llanos**, University of Valladolid: RDNest, new start-up company in the field of IoT in Valladolid, Spain
- **Alessandro Pellegrini**, Sapienza University of Rome: Transparent HPC Simulation on Heterogeneous Distributed Architectures
- **Spela Stres**, Jožef Stefan Institute: A technology radiation dosage manipulation and surveillance

Congratulations to all the winners!

Further information: bit.ly/HiPEAC_TTAwards_2018_winners

Dates for your diary

EuroHPC Summit Week 2019

13-17 May 2019, Poznań, Poland

exdci.eu/events/eurohpc-summit-week-2019



Photo credit: Mahwish Arif

ACACES 2019: Fifteenth International Summer School on Advanced Computer Architecture and Compilation for High-Performance and Embedded Systems

14-20 July 2019, Fiuggi, Italy

Registration deadline: 31 May 2019

acaces.hipeac.net/2019

2019 IEEE Nordic Circuits and Systems Conference

29-30 October 2019, Helsinki, Finland

Submission deadline: 15 August 2019

norcas.org

Euro-Par 2019: 25th International European Conference on Parallel and Distributed Computing

26-30 August 2019, Göttingen, Germany

europar.org

Euromicro DSD/SSEA: Euromicro Conferences on Digital System Design and Software Engineering and Advanced Applications

28-30 August 2019, Kallithea, Chalkidiki, Greece

dsd-seaa2019.csd.auth.gr

EuroMPI 2019

10-13 September 2019, Zürich, Switzerland

eurompi19.inf.ethz.ch

FPL 2019: International Conference on Field-Programmable Logic and Applications

9-13 September 2019, Barcelona, Spain

fpl2019.bsc.es



A physicist by training, Toni Collis, Chief Business Development Officer at Appentra and Chair of Women in High Performance Computing (WHPC), took an unconventional route into high-performance computing (HPC). In this interview, she explains her passion for HPC and why it's so important to open up the field to a more diverse group of users.

'If all HPC programmers are male, we are limiting the scientific discovery of the human race'

What got you interested in HPC?

I was encouraged at the beginning of my physics PhD to learn parallel computing and HPC to help me use it in my studies, and I was very lucky to be able to study a master's in HPC at the same time. At the end of my doctoral studies I realized that if I worked in HPC I could actually enable more science than by staying in the traditional physics background, because so many physicists (and others) needed help in parallel programming.

Why is high-performance computing (HPC) important?

I truly believe that high-performance computing provides a solution to many of humankind's challenges – everything from climate change, to the current energy crisis, to how we feed the rapidly growing human population. Our imagination is the only limit to how many questions we can use HPC to help answer. Just as one example, there has been a lot of media attention recently on how many species mankind has wiped out in the last 50 years. That is something that high-performance computing, through modelling, could help us understand, address and fix more quickly than experiments.

There is so much that we can achieve with HPC, but we need to figure out how to enable more people to use it as a tool to answer their questions. At the moment, HPC and parallel programming is very much an elitist solution to problems. We need to make this a tool that everyone can make use of.

"HPC provides a solution to many of humankind's challenges"

Why is parallel computing important for science? What are the barriers to scientists using HPC?

The computational limit is often first reached by what can be computed in serial. The more we can move from serial to parallel algorithms, the more we can answer in a shorter amount of time. The main challenge is understanding which calculations in serial code could actually be run concurrently. Understanding the flow of your code is the most complicated but also the most important step in going from serial to parallel and reaping the benefits that that enables.

Parallel computing is now everywhere; even your mobile phone has multicores and has the ability to do hyperthreading. Despite that, the majority of applications that you run are still single threaded, and there are lots of cores on your laptop that are frequently sitting idle. Programmers in a scientific environment, where they're programming in addition to their day job – that is, being an expert in the science question they're answering – don't necessarily have the time to learn how to become an expert in parallel programming in addition to their speciality. Parallel programming is currently a complex skill in its own right, and we can't afford to retrain all potential users of HPC in parallel programming, as that then detracts from the questions that they need to answer.

What about the lack of diversity in HPC?

One of my passions is getting more women and other under-represented groups to use high-performance computing, hence why I started Women in High Performance Computing (WHPC). We know, for example, that although in many countries 50% or more of biologists are now female, it is not the case that 50% of the biologists using HPC are female. What is going on that means that



Toni with WHPC colleagues at the PRACE booth at SC17

even when we have a domain with many women, these women are not engaging in HPC to help with their science question?

By not utilizing all resources available to them, such as HPC, these scientists' discoveries, careers and more are intrinsically limited. This is not to say that they aren't having fantastic careers, but if not all avenues of scientific analysis are open to you, you are limiting your ability to investigate. However, this goes beyond individuals – we're also limiting our collective ability to advance science.

One of the key benefits of diversity – the diversity dividend – is the increased knowledge when you have a diverse group of individuals involved in scientific discovery. A diverse group, operating well, is more likely to have diversity of thought and diversity of ideas. We know that you increase your scientific output, you get better results, you answer more problems, and you get more patents if you have a diverse team. If all, or almost all, HPC programmers are male, we're not just impeding a woman's potential career path, we're also limiting the scientific discovery of the human race.

So how do we get more underrepresented groups to engage with HPC? It's a snowball effect: if we get more women involved, more women will be interested. Diversity starts with inclusion: if your workplace is more inclusive, diversity will follow.

Tools that can open up HPC to a wider audience, such as those developed by Appentra, where we are working on tools to improve the productivity of both parallel programmers and also by improving HPC education, can also help, but they are one step on a long path.

Making HPC more accessible, from being more inclusive to making it easier to use, is something I believe everyone in the HPC community, irrespective of their gender, needs to step up and help address. If we do that, HPC and the entire human race will benefit, as we will be able to answer more challenging questions every day.

What has WHPC achieved so far? What are the plans for the future?

WHPC is five years old this April. We now have nine chapter and affiliate organizations, which we hope to expand later this year, and members from across the globe. We've also just launched our new mentoring programme which will bring together women from around the globe all year round. The benefits to the community are slow, and hard to measure: we can't say for certain that we have had an impact on the percentage of women working in the global HPC workforce yet. But we are ensuring that everyone is now aware of the importance of diversity, the benefits of inclusion, and, crucially, the fact that we all need to measure.

When WHPC was launched, no one was counting how many women, or people from other underrepresented groups, were working in HPC. Today, a growing proportion of the sector is monitoring this, which is the first but most important step in bringing about positive change.

Video interviews with Toni Collis are available on the 'HiPEAC experts' webpage: hipeac.net/press/#/videos

appentra.com

womeninhpc.org



Can Europe really compete with the United States and Asia in high-performance processor development? The European Processor Initiative, which includes a number of HiPEAC members, is an investment in made-in-Europe technologies for the future. Here, EPI Chief Communications Officer Mario Kovač introduces the initiative.

Homegrown high performance: The European Processor Initiative (EPI)

The importance of high-performance computing (HPC) has been on the rise over the last few years, and it is expected that this trend will not only continue but rapidly grow. Industry reports show that annual global IP traffic will soon reach several zettabytes, vast amounts of new devices collect and store data, and scientists are exploring new computing approaches to solving global challenges.

At the same time, industry is changing the way products are designed, while we, as individuals, constantly expect more personalized services for many areas of our lives. For example, better drugs that get to market faster, have fewer side effects and cost less; faster diagnostic tools that help medical doctors perform better treatments; autonomous cars which are safer and available at a lower cost; and many others.

The need to collect and efficiently process these vast amounts of data comes at a price. The existing approach to HPC systems design is no longer sustainable for the exascale era, in which computers will execute a billion billion calculations per second. Energy efficiency is of enormous importance for the sustainability of future exascale HPC systems.

Europe has recognized this global HPC challenge and has developed a strategic plan to support the next generation of computing and data infrastructures. European Union (EU) efforts are synchronized in establishment of the EuroHPC Joint Undertaking, a legal funding entity which will enable pooling of national and EU-wide resources in high-performance computing to acquire, build and deploy the most powerful supercomputers in the world within Europe.

efficient computing performance. As recognized by high-level EU officials, EPI will benefit Europe's scientific leadership, industrial competitiveness, engineering skills and know-how – not to mention society as whole.

EPI will use a holistic approach to refine the system architecture and its component specifications. All aspects of the solution, and their interactions, will be considered and tackled simultaneously, taking a co-design approach:

- hardware platform architecture and components
- system and runtime software (operating system, middleware, developers kit, libraries, etc.)
- HPC end-user applications

To fulfill its objectives of working towards a hybrid exascale system, EPI will develop:

- a novel, exascale HPC-focused low-power processing system unit
- an accelerator to increase energy efficiency for computing intensive tasks in HPC and artificial intelligence (AI)
- an automotive demonstration platform to test the relevance of the above-mentioned components in this industry sector

Three streams

For this purpose, EPI will harmonize the heterogeneous computing environment by defining a common approach: the so-called **Common Platform (CP)**. The



European Processor Initiative partners at kick-off event

EPI: A cornerstone of EuroHPC

The European Processor Initiative (EPI) is one of the cornerstones of this EU HPC strategic plan. Drawing on the expertise of 23 partners from 10 European countries, EPI aims to bring a low-power microprocessor to market. It will ensure that the key competence of high-end chip design remains in Europe, a critical point for many application areas. Thanks to such new European technologies, European scientists and industry will be able to access exceptional levels of energy-

CP will be organized around a 2D-mesh network-on-chip connecting computing tiles that may be general purpose central processing unit (CPU) core with floating point unit (FPU) acceleration, RISC-V based accelerators, massively parallel processor array (MPPA) accelerators, embedded field-programmable gate array (eFPGA) or any other application-specific accelerators.

The HPC General Purpose Processor (GPP) will be the first implementation of the common platform targeting the HPC market. During the project, we will develop the first generation of the GPP with two revisions:

- a real chip (GPP generation 1-revision 1), tested and validated as suitable to build a pre-exascale HPC system
- specification and intellectual property (IP) ready for the second revision with the aim of meeting the performance targets (FLOPS/W, FLOPS/socket and b/FLOPS memory bandwidth)

The Accelerator stream will develop and demonstrate fully European processor IP based on the RISC-V instruction set architecture (ISA) and aims at providing a very low power and high computing throughput accelerator to the general-purpose cores. Targeting an accelerator design will allow us to focus on the

stated fundamental objective to achieve increased performance.

Using RISC-V gives us the opportunity to leverage many open-source resources at architecture and system software level and at the same time to innovate and develop aspects of our vision that will allow us to differentiate. The accelerator processor will be based on the RISC-V vector ISA and will include specialized blocks for stencil and deep learning acceleration. The vector and stencil capabilities will address workloads in HPC centres, while the deep learning block will target learning acceleration.

The design of a novel HPC processor family cannot be sustainable without thinking about possible additional markets that could support such long-term activities. Thus, EPI's automotive activities have been carefully selected to ensure the overall economic viability of the initiative.

Current main trends driving innovations in the automotive industry include the introduction of autonomous driving (class 4 and 5) and the 'connected car' infrastructure. New autonomous vehicle electric and electronic architectures require computing platforms able to execute complex vehicle perception algorithms that include modelling of the

surrounding environment, sensor/imaging processing, data fusion, low-latency deep machine learning for object classification and behaviour prediction with seamless, dependable and secure interaction between mobile high-performance embedded computing and stationary server-based high-performance computing.

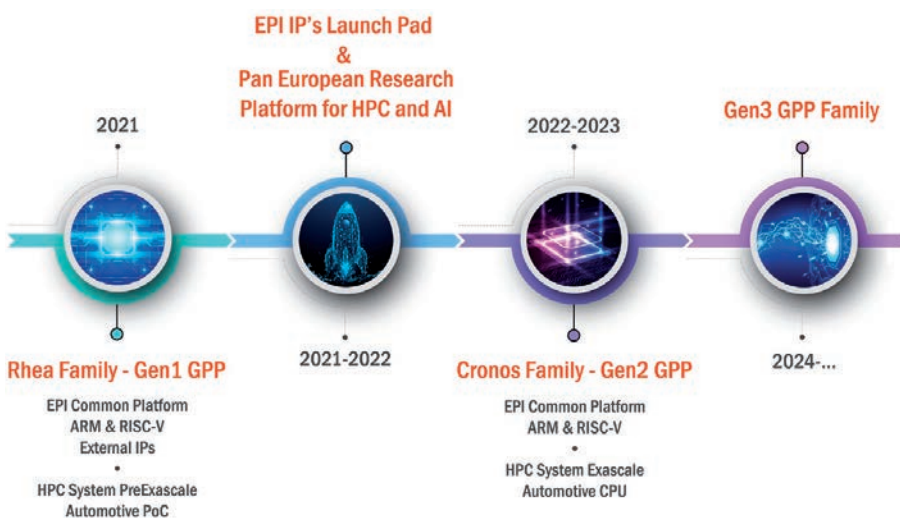
One specific objective for EPI in the automotive sector is to develop customized processors able to meet the performance needed for autonomous vehicles that would offer compute resources with the same characteristics (performance, performance per watt, bandwidth and so on) as their larger counterparts in exascale class supercomputers.

EPI brings together experts from the high-performance computing research community, the major supercomputing centres, and the computing, automotive and silicon industry as well as potential scientific and industrial users. Through a co-design approach, it will design and develop the first European HPC systems-on-chip and accelerators. Both elements will be implemented and validated in a prototype system that will become the basis for a full exascale machine based on European technology.

EPI will provide European industry and research with a competitive HPC platform and data processing solutions at world class level in the best interest of data security and ownership. We expect to achieve unprecedented levels of performance at very low power, and EPI's HPC and automotive industrial partners are already considering the EPI platform for their product roadmaps.

EPI has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 826647.

european-processor-initiative.eu



Squeezing the best performance out of computers has always been a priority, but with the complexity of modern computers and the demands of modern applications, it is arguably more challenging than ever. In this article, we find out about different solutions to the performance problem, how high-performance computing and high-performance data analytics are merging, how the international community is pushing towards exascale within a power budget, and how researchers are making it easier to use heterogenous hardware.

Performance, performance, performance: HPC, HiPEAC style

DEEP: A MODULAR SUPERCOMPUTER WHICH ADAPTS TO APPLICATION NEEDS

There's more than one way to build a supercomputer, and meeting the diverse demands of modern applications, which increasingly combine data analytics and artificial intelligence (AI) with simulation, requires a flexible system architecture. Since 2011, the DEEP series of projects has pioneered an innovative concept known as the modular supercomputer architecture, whereby 'multiple modules are coupled like building blocks', explains the project's coordinator Estela Suárez (Jülich Supercomputing Centre). 'Each module is tailored to the needs of a specific class of applications, and all modules together behave as a single machine,' she says.

Connected by a high-speed, federated network and programmed in a uniform system software and programming environment, the supercomputer 'allows an application to be distributed over several hardware modules, running each code component on the one which best suits its particular needs', according to Estela. Specifically, DEEP-EST, the latest project in the series, is building a prototype with three modules: 'a general-purpose cluster for low or medium scalable codes, a highly scalable booster comprising a cluster of accelerators, and a Data Analytics Module (DAM)', which will be tested with six applications combining high-performance computing (HPC) with high-performance data analytics (HPDA) and machine learning (ML).

The DEEP approach is part of the trend towards using accelerators to improve performance and overall energy efficiency – but with

a twist. 'Traditionally, heterogeneity is done within the node, combining a central processing unit (CPU) with one or more accelerators. In DEEP-EST we segregate the resources and pool them into compute modules, as this enables us to flexibly adapt the system to very diverse application requirements,' Estela says. In addition to usability and flexibility, the sustained performance made possible by following this approach aims to reach exascale levels, she adds.

One important aspect that makes the DEEP architecture stand out is the co-design approach, which is a key component of the project. 'Only by properly understanding what users and their applications really do can you build a computer that efficiently solves their problems,' says Estela. 'In DEEP-EST, we've selected six ambitious HPC/HPDA applications to drive the co-design process, which will be used to define and evaluate the hardware and software technologies developed.'

Careful analysis of the application codes allows a fuller understanding of their requirements, which informed the prototype's design and configuration, she notes. 'For instance, we used application benchmarks to choose the CPU version on each prototype module, while memory technologies and capacities have been selected based on the input/output needs of our co-design codes. Each application uses different module combinations, while dynamic scheduling and resource management software ensure the highest possible throughput.

In addition to traditional compute-intensive HPC applications, the DEEP-EST DAM includes leading-edge memory and storage technology tailored to the needs of data-intensive workloads, which occur in data analytics and ML. ‘Based on general-purpose processors with a huge amount of memory per core, this module is boosted by powerful general-purpose graphics processing units (GPGPU) and field-programmable gate array (FPGA) accelerators,’ says Estela. In the example of space



weather, the data analytics module is ideal for analysing high-resolution satellite images, Estela explains, while other parts of the application workflow – such as the interaction of particles emitted by the Sun with the Earth’s magnetic field, are distributed between the cluster module and the booster.

Through the DEEP projects, researchers have shown that combining resources in compute modules efficiently serves applications from multi-physics simulations to simulations integrating HPC with HPDA, to complex heterogeneous workflows such as those in artificial intelligence applications. ‘One of our most significant achievements was delivering a software environment which allows the “module farm” to work as a single machine,’ says Estela. ‘Now we have tangible results in the form of JURECA, the first modular computing system at production level.’

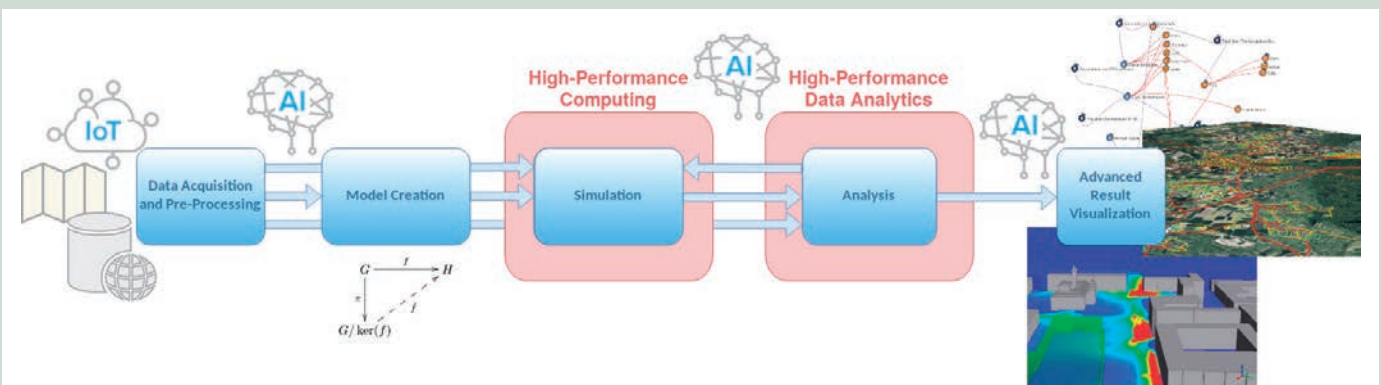
deep-projects.eu

The DEEP projects DEEP, DEEP-ER and DEEP-EST have received funding from the European Union’s Seventh Framework Programme for research, technological development and demonstration under grant agreement no ICT-610476 and no ICT-287530 as well as the Horizon2020 funding framework under grand agreement no. 754304.

SOLVING GLOBAL CHALLENGES WITH HIGH-PERFORMANCE COMPUTING AND DATA ANALYTICS

As noted in the HiPEAC Vision 2019, humanity is facing a number of unprecedented challenges which will require the full might of our collective knowledge and computational resources to solve. ‘From healthcare and climate change to stresses on the financial sector, these challenges involve highly complex systems and require an interdisciplinary, as well as international, response,’ says Konstantinos Nikas, a senior researcher at the Institute of Communication and Computer Systems, National Technical University of Athens. ‘Indeed, building evidence and understanding around global challenges is becoming vital to ensure the continuing (co)existence of humankind.’

To tame this complexity and channel new technologies towards solving global challenges, the European Union Centre of Excellence HiDALGO harnesses the power of high-performance computing (HPC) in combination with high-performance big data analytics (HPDA). ‘Bringing together stakeholders from different technology and application areas, HiDALGO focuses on coupled, multiscale simulations, highly dynamic workflows and, in particular, the integration of intelligent composition and adaptation methodologies,’ explains Konstantinos.



HiDALGO workflow

Image © 2019 Know-Center GmbH, Széchenyi István University, University of Stuttgart, and other members of the HiDALGO Consortium



The HazelHen supercomputer at HLRS

With high volumes of data being fundamental to the simulations, HiDALGO is exploring novel solutions that combine HPC with HPDA. In parallel, the project is investigating opportunities for integrating individual processing components such as data acquisition, data pre-processing and data analytics in an efficient, seamless workflow. ‘HiDALGO’s co-design process will also facilitate the adoption of our results in the exascale systems of the future – which will be data-centric by default,’ Konstantinos adds.

To this end, HiDALGO is creating a dedicated system architecture which focuses on the combination of multi-coupled simulations – HPC – with data analytics – HPDA – and interacts with a range of external elements including sensor data and visualization facilities. ‘HPC provides the computational muscle to execute highly complex simulations, while HPDA will analyse thousands of simulation results – up to petabytes in size – in situ, thus avoiding data transfer and input/output operation as far as possible,’ explains Konstantinos. ‘To complement this, artificial intelligence (AI) methods will tackle pain points such as data pre-and post-processing, and provide efficient mechanisms for parameter exploration. AI methods will also be explored in the analytics phase to automate trend or turbulence recognition, or to model behaviour,’ he adds.

Three compelling pilot applications have been chosen to test this technology, explains Konstantinos. ‘The first challenge we’re tackling relates to the movement of refugees. Understanding the dynamics of migration flows is crucial for the allocation of humanitarian resources,’ he says. ‘HiDALGO will extend FLEE, an agent-based simulation which leverages real-world data from UNHCR, the United Nations Refugee Agency, and the Armed Conflict Location & Event Data Project (ACLED). We will combine sophisticated HPC and HPDA methods to accurately predict massive refugee movements originating from different conflict regions.’

The second pilot application relates to air pollution, ‘one of the most urgent problems facing politicians and health organizations’, according to Konstantinos. ‘Currently, highly simplified models are used for pollution prediction in most cases, due to the complexity of accurate computational fluid dynamics (CFD) methods,’ he says. ‘HiDALGO will combine the 3DAirQualityPrediction application, developed as part of the MSO4SC project, with high-quality data from sensors in urban areas to develop HPDA-based prediction methods.’

Last but not least, HiDALGO aims to understand, model and simulate the processes behind information dissemination on social networks. ‘New clustering algorithms will be applied to derive synthetic graph models that accurately describe these networks. Next, stochastic models will be developed to build a reliable simulation framework to help us understand the dynamics behind the spread of messages,’ explains Konstantinos. ‘This understanding should enable easier identification of so-called “fake news” messages, giving decision makers the opportunity to adopt appropriate counter measures.’

HiDALGO – the European Centre of Excellence for High Performance Computing and Big Data Technologies in the domain of Global Challenges – has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement number 824115.

FURTHER INFORMATION:

hidalgo-project.eu

[@EU_HiDALGO](https://twitter.com/EU_HiDALGO)

[EU.Project.HiDALGO](https://www.facebook.com/EU.Project.HiDALGO)

FLEE on GitHub github.com/djgroen/flee-release

3DAirQualityPrediction mso4sc.eu/?page_id=110

HOT STUFF: SEAMLESS FPGA INTEGRATION WITH TAPASCO



Field-programmable gate arrays (FPGAs) were highlighted as the star of the show at the 2019 HiPEAC conference by *The Next Platform*, and their appeal for research is clear, according to Andreas Koch, leader of the Embedded Systems Group at Technische Universität Darmstadt. ‘In

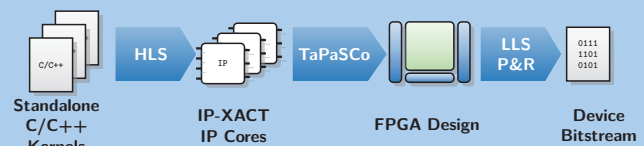
recent years, we’ve seen FPGAs being successfully employed for the implementation of many application-specific accelerators, for example in the field of machine learning inference,’ he notes.

The flexibility of FPGAs give them an edge over other processors in this area, says Andreas. ‘The generic design of central processing units (CPUs) means they cannot provide high efficiency for every task, while graphics processing units (GPUs) are tailored specifically to massively parallel execution. FPGAs, on the other hand, allow you to adapt the underlying architecture to the task in hand, meaning that they can achieve better energy efficiency.’

However, it is this very flexibility which makes developing an FPGA-based accelerator challenging, as Lukas Sommer, a research associate in the Embedded Systems Group, notes. ‘Compared to software programming for CPUs or GPUs, FPGA programming requires a completely different set of skills and techniques,’ he says. ‘Although this situation has improved with the advent of mature high-level synthesis (HLS) tools which allow programmers to use high-level languages such as C/C++ and OpenCL for FPGAs, the integration of FPGAs into a heterogeneous system remains a challenge for many practitioners.’



The TaPaSCo team at TU Darmstadt. Left to right: Carsten Heinz, Jaco Hofmann, Lukas Sommer and Andreas Koch



The TaPaSCo framework supports developers at every stage

It was this challenge that inspired the creation of the TaPaSCo framework as part of the European Union FP7 project REPARA, coordinated by the Universidad Carlos III de Madrid. ‘TaPaSCo was created to enable the speedy integration of FPGA-based accelerators into heterogeneous compute platforms or systems-on-chip (SoCs),’ explains Lukas. ‘Since REPARA ended, we’ve been enhancing TaPaSCo with numerous stability improvements, broader platform support, a unified kernel driver for all supported platforms, and many new features, such as support for local or scratchpad memories.’

As a result, TaPaSCo tools and application programming interfaces (APIs) provide a turnkey solution to heterogeneous system design on a variety of platforms, ‘from small FPGA SoCs such as Xilinx’s Pynq, to large accelerator cards attached by peripheral component interconnect express (PCIe), such as the Xilinx VCU-1525 with large Xilinx Virtex Ultrascale+ FPGAs’.

Given its potential usefulness to other researchers in the field, the group decided to make TaPaSCo open source. ‘The framework can support developers at all stages of the heterogeneous system development process,’ says Jaco Hofmann, also a research associate in the Embedded Systems Group. ‘A developer using TaPaSCo simply specifies the desired composition – the type and number – of cores, which we call “processing elements”. TaPaSCo will then automatically connect all processing elements to the memory and host interfaces before generating a complete bitstream which can be used to program the FPGA.’

Another interesting feature, according to Jaco, is the automatic design-space exploration, which can be used to determine the best possible composition and operating frequency for the heterogeneous system design. ‘This is particularly useful for researchers interested in investigating trade-offs between area and frequency,’ he points out. Thanks to its flexible API for interacting with the accelerators, TaPaSCo also enables rapid prototyping once the FPGA has been programmed with the bitstream generated.



High-performance computing special

TaPaSCo is still being developed and improved, says Andreas. ‘We are always working on adding more platforms, for example support for well-known cloud computing providers,’ he says. Features such as support from RISC-V softcores are under development, while others – such as Ethernet support – have been added by user request. ‘We invite all researchers in the field to download TaPaSCo and give us feedback – or even contribute to its development,’ says Andreas.

REPARA received funding from the European Union’s Seventh Framework Programme for research, technological development and demonstration under grant agreement no 609666.

Download TaPasCo from GitHub:

github.com/esa-tu-darmstadt/tapasco

REPARA project

repara-project.eu

POWERSTACK: A GLOBAL RESPONSE TO THE POWER MANAGEMENT PROBLEM FOR EXASCALE

As explored in the HiPEAC Vision 2019, energy is an issue affecting the entire compute continuum, from tiny devices at the edge to enormous data centres. In the push towards exascale computing, improving energy efficiency is a driving factor for a number of reasons, says Professor Martin Schulz, chair for Computer Architecture and Parallel Systems at Technische Universität München (TUM) and a member of the Board of Directors at Leibniz Supercomputing Centre (LRZ).

‘One of the primary challenges for exascale is the total cost and variability associated with the energy and power consumed, not only by the high-performance computing (HPC) system but by the infrastructure supporting it,’ adds Tapasya Patki, a computer scientist at the Center for Applied Science Computing and principal investigator of the ECP Power Steering project at Lawrence Livermore National Laboratory (LLNL). ‘Many potential exascale sites are bound by power constraints of around 20-30MW. There may also be external factors such as a shortage of electricity, natural disasters and/or government-

issued mandates that limit the supply of power even further. For others, reducing electricity costs in order to improve purchasing power is a key motivation.’

In the race to greater performance using less energy, the focus is often on hardware; however, as Tapasya argues, system software also plays a crucial role in achieving higher throughput and better utilization in constrained scenarios. To achieve this, the HPC community needs to better understand the underlying technical aspects of power and energy management, she stresses. ‘For example, there is a misguided assumption that giving more power to an application will always improve its performance, and that enforcing a power cap will always slow it down. Although true for processor-bound applications, such as high-performance LINPACK, this does not apply to most scientific applications, which tend to be bound by memory, input/output or network usage.’

Another less understood aspect, according to Tapasya, is processor manufacturing variability, where processors with the exact same



The PowerStack core committee members (left to right, top to bottom): Aniruddha Marathe (LLNL), Barry Rountree (LLNL), Carsten Trinitis (TUM), Christopher Cantalupo (Intel), Jonathan Eastep (Intel), Josef Weidendorfer (TUM), Martin Schulz (LRZ, TUM), Masaaki Kondo (RIKEN, University of Tokyo), Matthias Maiterth (LMU, Intel), Ryuichi Sakamoto (University of Tokyo), Siddhartha Jana (EEHPC-WG, Intel), Tapasya Patki (LLNL, ECP)

microarchitecture exhibit different power and performance characteristics. ‘This is attributed to the chip fabrication process, and several vendors, including Intel and IBM, have confirmed that such variability is expected to worsen in the future and at larger scales.’

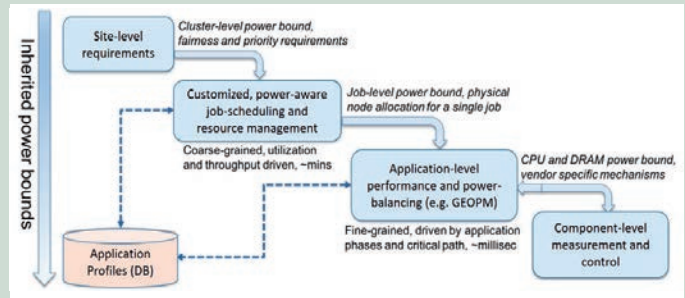
The answer? A software stack that can steer power based on site requirements, application characteristics and dynamic behaviour – all in a vendor-neutral way, says Siddhartha Jana, an HPC research scientist at Intel and subteam co-lead in the global Energy Efficient HPC Working Group (EE-HPC WG). ‘HPC PowerStack is a community-wide consortium that started in 2016 with the aim of bringing together experts from academia, research laboratories and industry to design a holistic, extensible power management framework,’ explains Siddhartha.

PowerStack explores hierarchical interfaces between components at three specific levels: batch-job schedulers, job-level runtime systems and node-level managers, according to Siddhartha. ‘Site-specific requirements such as cluster-level power bounds, user fairness or job priorities will be translated as inputs to the job scheduler. The job scheduler will choose power-aware scheduling plugins, managing allocations across multiple users and diverse workloads,’ he says. ‘Such allocations will serve as inputs to a fine-grained, job-level runtime system that manages application ranks, in turn relying on a vendor-agnostic, node-level measurement and control mechanisms.’

PowerStack forms part of the wider United States Exascale Computing Project (ECP), which aims to develop an HPC ecosystem – system software, applications, platforms and computational science, along with workforce development – using a co-design approach. ECP comprises three different projects targeting different levels of the stack:

- Runtime system for application-level power steering, focusing on the safe execution and performance optimization of applications running in a power-constrained environment.
- Operating system and resource management for exascale, focusing on improving and augmenting the operating system (Argo) and associated resource management frameworks (Flux).
- Exascale performance application programming interface (API), focusing on designing a ‘consistent interface and methodology’ for monitoring hardware and software-based performance events.

‘PowerStack is perfectly aligned with the first two projects. In the runtime project, the aim is to extend the use of the Global Extensible Open Power Manager (GEOPM), an open-source, community-driven power management application,’ says Tapasya.



An overview of the envisioned PowerStack

‘PowerStack is also actively collaborating with the ECP Argo and Flux teams to develop a more holistic power management stack. The vision is to make resource managers and job schedulers interoperable with power / performance management frameworks such as GEOPM,’ she adds.

As for the collaboration’s results so far, PowerStack has demonstrated that ‘HPC sites can improve system efficiency by overprovisioning their resources and incorporating a scalable power-aware resource management framework, demonstrated using the widely used Slurm workload manager’, according to Tapasya. ‘At application-node level, contributors have demonstrated anywhere from 5% to 30% of performance improvement depending on application design and architecture of power-constrained systems using GEOPM’, adds Siddhartha.

Critically, PowerStack solutions are tested across multiple HPC facilities to ensure that they cater to the needs of a range of global sites, says Siddhartha. ‘We’ve been collaborating with the Energy Efficient HPC Working Group, which led the first global survey analysing the current solutions for HPC sites in France, Italy, Japan, Saudi Arabia, Germany, the United Kingdom, and the United States.’ Since 2016, the size of the PowerStack consortium has grown considerably, and participants include national labs, system integrators, chip vendors, job scheduler and resource management vendors, and academic institutions.

‘We would like to invite more collaborators, and are actively planning forums during the ISC19 and SC19 timeframes,’ adds Professor Masaaki Kondo, a research lead at the Advanced Institute for Computational Science in RIKEN and an associate professor of the Graduate School of Information Science and Technology at the University of Tokyo.

An extended version of this article, including full list of reference papers, is available at bit.ly/HiPEAC_PowerStack

Contact points for PowerStack include Martin Schulz (Technical University of Munich), Masaaki Kondo (University of Tokyo/RIKEN), Tapasya Patki (LLNL/ECP), Siddhartha Jana (EEHPC-WG), and Jonathan Eastep (Intel/GEOPM PI).

powerstack.lrr.in.tum.de

With Europe currently lagging behind the United States and China in terms of the highest performing computers, the European Union is investing in research to help us reach exascale. In this article, HPC application researcher Mike Ashworth (University of Manchester), EuroEXA coordinator Georgios Goumas (Institute of Communication and Computer Systems, National Technical University of Athens) and EuroEXA dissemination lead Peter Hopton (ICETOPE) explain how EuroEXA's groundbreaking, co-designed architecture is providing the template for exascale.

Top of the FLOPS: Scaling up Europe's performance

High-performance computing (HPC) is about using the biggest, fastest computers around for simulation and data processing in support of scientific research in academia, industry and government. The performance of these HPC systems has been growing exponentially over at least five decades, with the raw number crunching power increasing by about 1,000-fold every ten years. Exascale refers to the next target, the next 1,000-fold leap in compute power.

This awesome power is measured in the number of calculations the system can compute every second, measured as floating point operations per second or 'FLOPS'. The current world-leading systems operate at around 100 petaFLOPS or 10^{17} FLOPS – the current number one system, 'Summit' at Oak Ridge National Laboratory in the USA, is rated at 143.5 petaFLOPS. Another factor of ten is needed to get us to exascale, or 10^{18} FLOPS.

HPC has become critically important across a wide range of human endeavours with immeasurable impacts upon our daily lives, our economy and our quality of life. A report by the PRACE Scientific Steering Committee, *The Scientific Case for Computing in Europe 2018-2026*, lists the following key areas:

- fundamental sciences
- climate, weather and earth sciences
- life sciences and medicine
- infrastructure and manufacturing
- chemistry and materials sciences
- complexity and data
- next generation computing

The report concludes that 'enhanced synergies between scientists working on hardware, algorithms, and applications is required for advancing the frontiers of science and industry in Europe for the benefit of its citizens'.

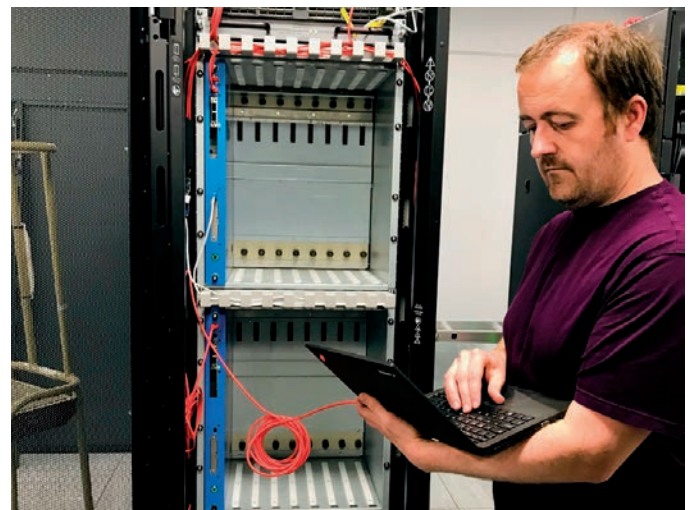
A number of critical research areas and problem classes, including weather prediction with fine granularity, climate change, large eddy simulation for turbulence modelling in aeronautics and the challenges in fusion energy research, are beyond current computing capabilities; they need exascale performance and even more.

The impact of exascale computing in many of these areas has been analysed in *The Opportunities and Challenges of Exascale Computing*, a report from the United States Department of Energy Office of Science. This report concludes that '[e]xascale computing will uniquely provide knowledge leading to transformative advances for our economy, security and society in general. A failure to proceed with appropriate speed risks losing competitiveness in information technology, in our industrial base writ large, and in leading-edge science'.

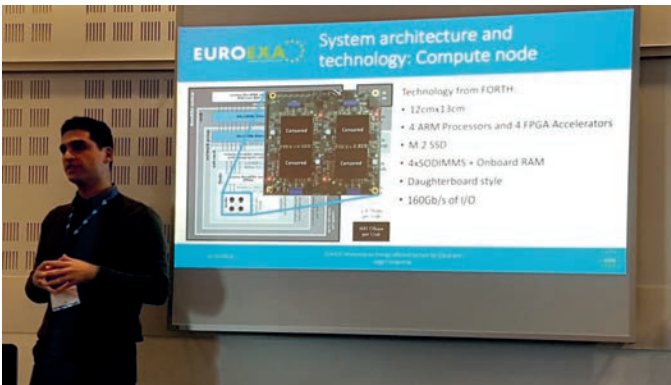
Building on European research

Launched in September 2017, EuroEXA draws upon previous European research to deliver the recipe for an exascale computer by the project's end in 2021. Specifically, this major (€20 million) project builds on the work of three projects, ExaNoDe, ExaNeSt and ECOSCALE.

Using co-design principles, we will physically demonstrate a testbed deployment providing an expected 2-4 petaFLOPS peak performance in an operational environment. The groundbreaking system architecture will be developed and optimized through a series of three testbed systems with increasing levels of performance and increasing sophistication, which are being developed during the course of the project.



Andrew Attwood (University of Manchester) beginning the installation of Testbed 1



Project coordinator Georgios Goumas (left) and Paul Carpenter (Barcelona Supercomputing Center) presenting EuroEXA at the 2018 HiPEAC conference

EuroEXA systems will be equipped with an optimized system software stack building on the work on the operating and runtime systems from the ExaNoDe project. They will be assessed through a wide set of HPC applications representing a diverse range of scientific subject areas. The testbed architecture will be shown to be capable of scaling to world-class peak performance in excess of 400 petaFLOPS with an estimated system power of around 30 MW peak.

Groundbreaking architecture responding to application needs

EuroEXA combines state-of-the-art computing components using a groundbreaking system architecture. This applies the design flexibility of UNIMEM, a scalable memory scheme first developed during the FP7 EUROSERVICES projects and used in ExaNeSt and ECOSCALE. The architecture delivers high levels of performance to the selected applications and balances compute and reconfigurable acceleration resources with the demands of applications. Through co-design between the enabling technologies, the system software and the applications, EuroEXA is delivering an innovative solution that achieves both extreme data processing and extreme computing.

Work on applications started with an assessment of application requirements and has progressed to porting and optimization work, using testbed systems similar to the EuroEXA architecture, to begin offloading computational kernels onto field-programmable gate arrays (FPGAs). The project is now working to implement a rich system software and runtime stack that will ensure applications can fully exploit the novel characteristics of the underlying architecture.

EuroEXA is evolving both traditional HPC programming paradigms (such as MPI and OpenMP) and novel ones including programming support for FPGA acceleration (such as OmpSs@FPGA), task-based, multi-node programming (such as OmpSs@clusters, OpenStream, UNIMEM-based) and streaming dataflow programming (Maxeler dataflow).

The first of three testbeds, Testbed 1, has been designed, built and installed. Testbed 1 consists of Quad-FPGA Daughter Boards (QFDB) connected by high-speed links in an innovative and state-of-the-art packaging design. A key part of the design process is co-design, with application requirements feeding into the design parameters. We have updated the specification for Testbed 2 resulting in a new design, the Co-Design Recommended Daughter Board (CRDB), which is forecast to result in four times greater performance.

Initial designs have also been produced for a novel, hybrid, hierarchical, low-latency, high-performance network, as well as for a multi-central processing unit (CPU) custom application specific integrated circuit (ASIC) for the core of the Testbed 3 compute node. The ASIC features custom hardware for implementing UNIMEM global addressing and memory compression.

At an early stage, the co-design process revealed a requirement for a redesign of the daughter board for the Testbed 2 system, meaning that we had to rethink some of the tasks in the project. However, thanks to flexible and collaborative working, we have managed to adapt swiftly and turn this significant challenge into a major success.

Towards the end of 2019, we'll be reaching a major milestone with the deployment of the Testbed 2 system. Application porting and optimization, together with porting and optimization efforts of parallel language runtimes, will start to reveal the benefits of the EuroEXA architecture, especially the memory and communications aspects, which are key to whole system performance across the cluster. Moreover, ongoing work on performance modelling will be able to provide initial evidence on how the EuroEXA testbeds can lead to effective exascale machines.

EuroEXA has received funding from the European Union's Horizon 2020 research and innovation programme under Grant Agreement no 754337

Simplifying FPGA complexity



COMPANY: Campera Electronic Systems Srl

MAIN BUSINESS: high-performance FPGA design service, FPGA IP core design, embedded systems hardware prototyping

LOCATION: Pisa (Navacchio), Italy

campera-es.com

CONTACT: Dr. Calliope-Louisa Sotiropoulou, Research and Development Manager

c.sotiropoulou@campera-es.com

Campera Electronic Systems is an innovative Italian small/medium enterprise (SME) founded to respond to the needs of the electronics market for high-performance, high-quality field-programmable gate array (FPGA) systems and custom designed intellectual property (IP). The company was founded in 2014 by Andrea Campera and Gabrielle Dalle Mura, two electronics engineers with extensive experience in FPGA systems development. The idea was to build a company with a robust workflow to provide custom designed IPs and FPGA design services as fast as possible and to the highest standard, all while offering excellent customer service.

Although relatively young, the company has been constantly growing and is now well established in the Italian market. Its major client share comes from the military and radar market, who benefit from the company's expertise in video and radar processing applications. Campera Electronic Systems also has strong collaborations with major research and academic institutes in Italy (University of Pisa, the National Institute of Nuclear Physics (INFN), National Institute of Astrophysics (INAF) and others), and continues to expand its collaboration network through a variety of research projects and applications.



The Campera Electronic Systems team

The company has developed a set of IP core libraries architected, developed, verified, released and maintained through a rigorous and efficient process. The IPs are vendor independent, 'off the shelf' VHDL cores for FPGAs (Xilinx, Altera, Lattice and Microsemi), optimized in terms of speed, power and resource usage. The mathematics and digital signal processing (DSP) IP cores can be supplied both in fixed point and in IEEE-Standard-754 compliant floating point. The current portfolio includes: the utility library (available for free under GPL, the General Public License), mathematics library, input/output library (such as MAC, Ethernet, SATA), video library (such as digital stabilization and tracking), DSP (for example, super sample rate fast Fourier transform (FFT) cores) and radar library (for example, constant false alarm rate (CFAR) and pulse compression filter).



Thanks to its expertise, the company is participating in some major projects, such as the Square Kilometre Array (SKA) project, an international effort to build the world's largest radio telescope, with a square kilometre (one million square metres) of collecting area. Campera Electronic Systems has been chosen as the FPGA and algorithms designer with the Arcetri Astrophysical Observatory, part of the National Institute for Astrophysics (INAF). The company designed a super sample rate polyphase filter bank (PFB) channelizer custom IP cores in VHDL (up to 8 GSaPS).

Campera Electronic Systems is now entering the biomedical market as the project leader of ultraVISTA, a project founded by Tuscan Region research funds. The ultraVISTA project proposes a unique combination that is missing from the biomedical market: a headset that offers advanced image magnification with real-time digital stabilization and video processing, combined with the advantages of an augmented reality environment, all within an attractive price range. The final prototype – planned to be commercialized as soon as possible – benefits from the company's extensive video processing IP library and video stabilization experience.

FURTHER READING:

Square Kilometre Array skatelescope.org

ultraVISTA ultra-vista.it

Hardware which is sent into space is subjected to some extreme conditions, which it has to withstand for several years in a specific mission. In this article, Aubrey Dunne (Ubotica Technologies) explains how the Intel® Movidius™ Myriad 2 Vision Processing Unit was subjected to radiation testing at CERN to prepare it for its long journey.

Preparing the Intel® Movidius™ Myriad 2 VPU for space

To infinity and beyond

In November 2018, at a particle accelerator in CERN, the Myriad 2 Vision Processing Unit (VPU) was subjected to an extensive series of radiation tests in order to characterize it for use in space applications. The Intel® Movidius™ Myriad™ 2 is a low-power, high-performance system-on-chip (SoC) designed specifically for computer vision and artificial intelligence (AI) inference. Space characterization of the Myriad 2 is being performed by HiPEAC member company Ubotica Technologies, a Dublin-based computer vision and AI company, under contract with the European Space Agency (ESA).

Myriad 2 was irradiated over the course of three days by the H8 beamline of the Super Proton Synchrotron (SPS), a particle accelerator in the CERN complex in France that is most famous for being the particle injector to the Large Hadron Collider (LHC). The SPS accelerated a beam of lead ions to a speed just below the speed of light and emitted this ion beam into a device test chamber encased by one-metre-thick concrete blocks. The Myriad 2 test board was mounted in the path of this beam, with tests conducted and monitored remotely from a control room.

Ambient radiation in space is sufficient to cause undesirable single event effects (SEEs) within the silicon structure of semiconductor devices. Such SEEs can be internal latch-ups that cause excessive



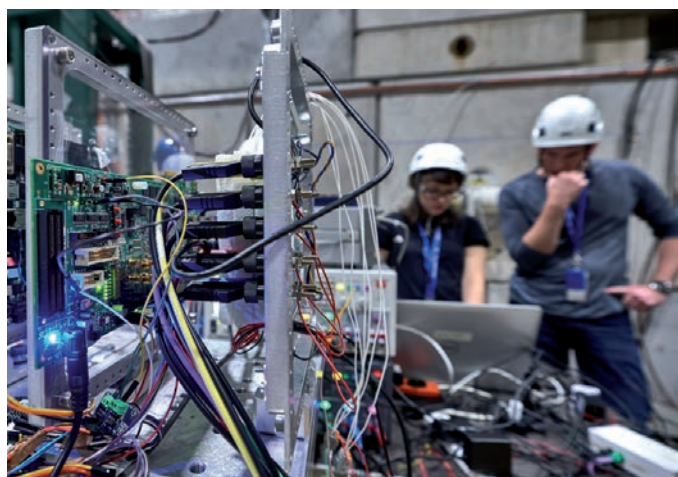
Myriad 2 test board

Photo credit: Maximilien Brice/CERN

current draw, bit upsets in memories and registers, or functional interrupts that cause processor operation divergence. While the first mode of failure is destructive if not handled by latch-up protection circuitry, the latter two modes can often be handled functionally by suitable software mitigation techniques.

The purpose of the space characterization of Myriad 2 was to determine its performance when irradiated in a test campaign that replicates but accelerates the likely radiation exposure in space. By carefully testing each of the functional blocks within Myriad 2 during ion bombardment, valuable information for software mitigation and latch-up protection is acquired that can be used to predict the radiation impacts on Myriad 2 over the course of multi-year space missions.

Following on from the successful test campaign at CERN, Myriad 2 is being integrated into an in-orbit demonstration CubeSat for launch in the second half of 2019. This exciting development is the first time that an AI chip will fly on a CubeSat mission, with the aim of demonstrating the suitability of Myriad 2 as a space-borne neural-network inference engine. Such a low-power AI engine can apply AI to a broad class of autonomous identification and classification tasks, such as cloud detection, ship detection and fire detection, which are currently performed on the ground due to size and power constraints. Together Ubotica Technologies, ESA and the Intel® Movidius™ Myriad 2 are bringing the AI revolution to space.



Myriad 2 test setup in the SPS H8 beam room at CERN

Photo credit: Maximilien Brice/CERN

Breaking through the cloud I/O

Based in Cambridge, the OnApp spinout Sunlight.io develops technology offering both flexibility and performance for data centres. In this article, Julian Chesterfield, Michail Flouris and Stelios Louloudakis introduce the company's hyperconverged platform and explain what gives it a competitive edge.



Around five years ago we witnessed the emergence of 'hyper-convergence' as an information technology (IT) framework that combines storage, computing and networking into a single system in an effort to reduce data centre complexity and increase scalability. Hyperconverged platforms include a hypervisor for virtualized computing, software-defined storage and virtualized networking.

Hyperconverged infrastructure is a truly disruptive technology that has significantly changed the compute landscape. Initially focused around the storage layer, it enabled early adopters to free themselves from the storage array vendor technology lock-in that was becoming so prevalent. One of the major benefits is that it allows scalable infrastructure using commodity off-the-shelf systems without the same sort of inherent scaling issues that are common in traditional virtualized infrastructure with disk arrays, in order to deliver simplicity and flexibility in comparison to legacy solutions.

However, one of the limitations with pure hyperconvergence is that it is not possible to scale storage capacity and performance independently from the compute capacity, which for certain workloads that have asymmetric compute and storage demands, such as big data applications, can be inefficient.

The problem

Flexibility is therefore essential, but there is one very important factor that is missing from all the above: performance. For a long time, the industry believed that you could have flexibility, or you could have performance, but you couldn't have both. This was due to a number of reasons, which have since changed.

First, access latencies of non-volatile memory express (NVMe) flash storage have significantly decreased. On the current generation of Intel Optane 3D X-point technology you can read a 4K block in under 10 microseconds, and on many storage drives on the market today from a variety of vendors you can achieve between 800,000 and one million I/O operations per second (IOPs). In comparison, five years ago we were talking about deploying disk arrays of 32

or 64 SAS drives in order to achieve this sort of performance. In addition, the access latency at that time, and for that level of performance, was one or even two orders of magnitude larger than NVMe flash drives today, with all the unpredictability of mechanical head access latency depending on where you were reading data physically from the drive.

Second, in the past, high-performance storage and networking was far less of a concern than virtualized central processing unit (CPU) core efficiency. In fact, storage and network devices were always a massive bottleneck in the whole system. Fast forward 15 years and we see a very different picture, where driving even a single NVMe or single 100Gbit Ethernet network interface controller (NIC) at maximum line rate suddenly requires multiple CPU cores and ultra-efficient, non-uniform memory access (NUMA)-aware scheduling. This is where SUNLIGHT comes in.

Sunlight.io is a spinout of OnApp, the company powering over one in three of all the public clouds for managed services providers (MSPs), telecommunications and large hosting providers around the world. The Sunlight.io technology was developed in the OnApp emerging technologies research and development (R+D) labs over a period of years.

The Sunlight.io solution

Sunlight.io has developed and released the fastest fully converged technology platform on the market today, based on cutting-edge technology components which have been in production systems for several years, helping to power the public cloud. Sunlight.io is the premier platform for 5G and network function virtualization (NFV) applications and is constantly updated in order to meet new market demands.

Traditional platforms are not well suited to the high-performance network requirements of NFV. However, as a fully converged platform that supports both hyperconverged and disaggregated deployments, all at the fastest bare-metal virtualization speed, we believe that the Sunlight.io platform is the best available solution

The Sunlight.io platform: vital statistics

- over one million storage I/O operations per second (IOPs) per virtual machine (VM)
- 200Gbit/s Ethernet connectivity and I/O virtualization acceleration for network and storage traffic (for environments in need of the convenience and security of full I/O virtualization)
- OpenStack-compliant application programming interfaces (APIs)
- bare-metal performance at the speed of today's fastest Ethernet and NVMe flash storage technology

0 bottleneck with Sunlight.io

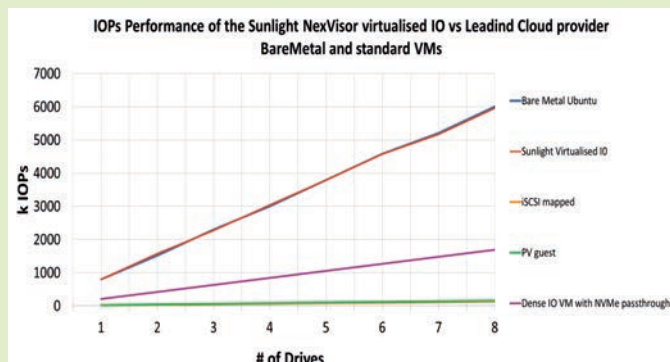
Intuitive and practical interface

The Sunlight.io Enterprise Software Platform is a fully converged infrastructure system optimized to provide high-performance I/O for both networking and storage traffic to and from tenant virtual machines. The Sunlight.io Enterprise Software Platform is comprised of a lightweight virtualization hypervisor, called the NexVisor, a distributed block storage volume management layer based on the mature Sunlight.io Integrated Storage technology; and an OpenStack-based controller that runs as a tenant Virtual Service on top of one of the NexVisors in the cluster.

The Sunlight.io Enterprise Software Platform is optimized on a specific set of hardware platforms, such as clusters of high-end Intel BobcatPeak processor systems, Xeon-D Broadwell systems and Kaleao KMAX 12 Blade/192 server systems. As new hardware becomes available, additional optimizations are developed constantly, in order to support new systems.

The SUNLIGHT Company

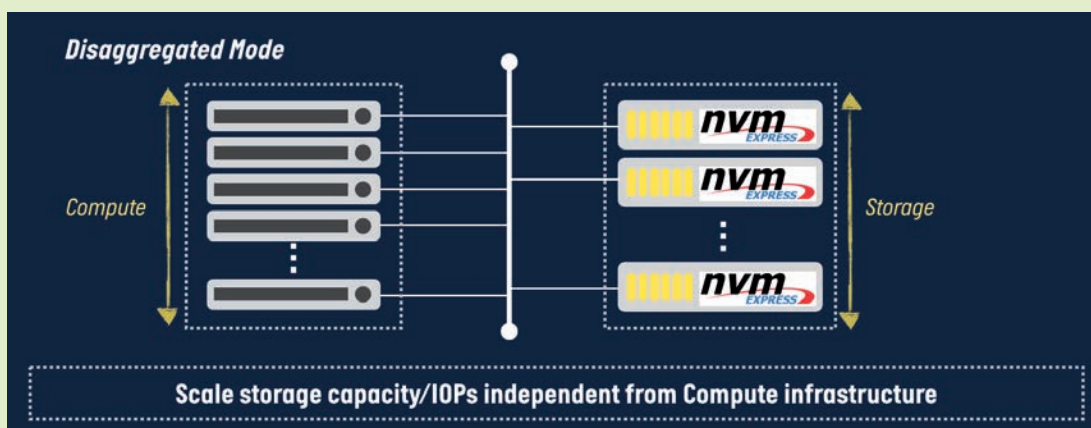
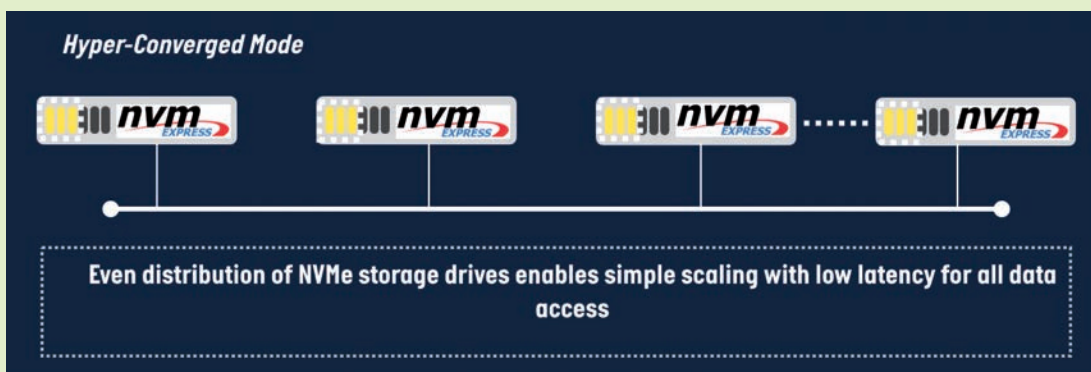
Using cutting-edge technologies and fully exploiting its knowhow, the Sunlight.io team, through its incubation within OnApp, has successfully participated in several European research projects, in which it has delivered innovative products and services. With a track record of exploiting the research and development



Virtualized IOPs scaling as measured on the Sunlight platform on a leading cloud provider baremetal server platform. Sunlight virtualized performance scales equally to baremetal at six million IOPs.

work, the team has successfully aligned commercial needs to the research goals of the projects.

Sunlight.io headquarters are located in Cambridge, UK with teams based around the world, including the US. The company is backed by a world-class team of engineers, with expertise in low-level hardware design, firmware/drivers, hypervisors, storage and network I/O, cloud management/orchestration and web UI technologies.





While advanced computing can be immensely beneficial for society, it can be hard to bridge the gap from research to industry. Eurolab4HPC has been meeting this challenge head-on with practical and financial support to help technology reach the market. We caught up with Eurolab4HPC Coordinator Per Stenström (Chalmers University) to find out more.

Taking HPC from the lab to t

Why is it important for advanced computing technologies to reach society?

The reasons are manifold. First of all, industry is undergoing a revolution in digitization in which more and more advanced computing techniques are integrated into new products and in which the manufacturing techniques themselves rely on advanced computing techniques.

“Industry is going through a revolution in digitization”

Second, the internet of things (IoT) has started a brand new industry in which computerized devices close to the sensors must be capable of analysing large volumes of data in real time. This calls for innovations in advanced computing technologies.

Finally, we have all witnessed the revolution of artificial intelligence (AI). This came as advanced computing technologies passed a performance threshold, making it feasible to run deep learning algorithms sufficiently fast to enable a broad spectrum of use cases for deep learning. Indeed, we are living at an exciting time when it comes to the need of advanced computing technologies.

High-resolution maps for speedy earthquake response

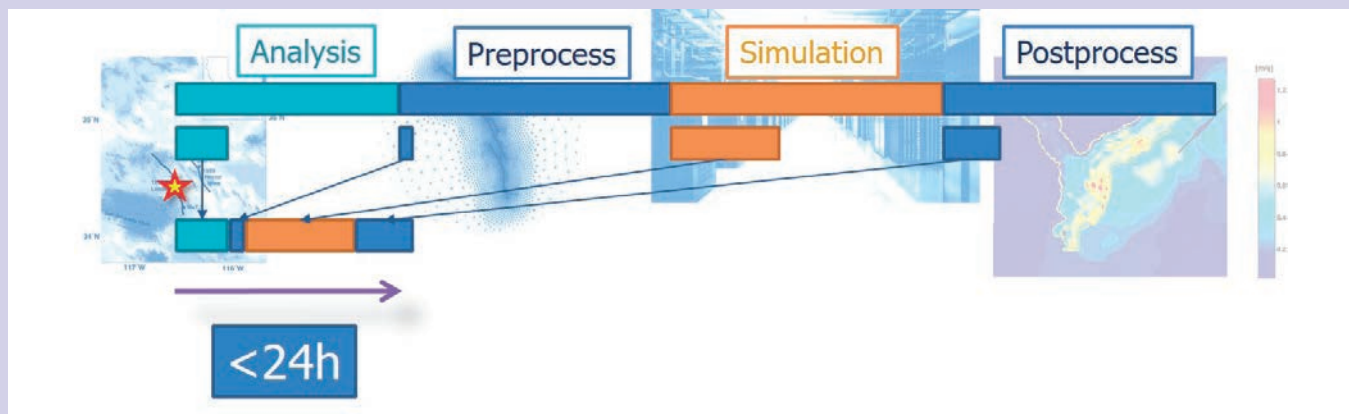


Josep de la Puente (Barcelona Supercomputing Center)

By means of the project USER, and with help from the CHEESE Centre of Excellence, we are trying to enable the use of supercomputers for civil protection purposes. Our main goal is using advanced simulations to generate high-resolution ‘shake maps’ after an earthquake. The maps contain detailed information of the most affected

zones and thus can help the disaster response logistics. The challenge is being able to generate such maps shortly after the earthquake. Only the most powerful supercomputers are able to achieve this task.

EurolabHPC is proving an invaluable instrument to push academic development into business prototypes. It is a very helpful first step towards testing out our technology in the real world. In our case we are making sure that we can match our value proposition with potential customers’ needs.



he market with Eurolab4HPC

How is Eurolab4HPC approaching technology transfer?

EuroLab4HPC is approaching technology transfer in what we refer to as the innovation pipeline. As a first step, we arrange ‘idea-to-business’ meetings where participants get trained in how to build a business case around a unique technical solution. Topics include stating hypotheses on what the value is, who are the customers and how can the value be transformed into a revenue stream. Then we also arrange open-source innovation camps that inspire the participants to come up with innovative solutions to potential pain problems in business.

We also have open calls for what we refer to as business prototyping projects (BPPs) and technology transfers, where the idea is to provide grants for the specific purpose of validating a business hypothesis and transferring a technology to a company, respectively. A successful BPP project will have reached a point where investors or venture capitalists (VCs) find the project investable. To help at that stage, we also arrange venture capital events where BPPs are presented to VCs.

In the stories cited here, you can see a couple of examples of how Eurolab4HPC has been supporting participants to transform their technology into usable products and services.

Implementing advanced storage monitoring with help from Eurolab4HPC



DDN offers a variety of enterprise storage products like the Infinite Memory Engine (IME). IME is designed as a burst-buffer which is optimized for random input/output (I/O). ‘However, understanding the reasons that lead to observed performance behaviour is non-trivial,’ explains **Jean-Thomas Acquaviva** (DDN France).



Performance analysis of the I/O path in computer systems is a crucial task, according to **Julian Kunkel** (Department of Computer Science, University of Reading), as it paves

the way for subsequent optimization not only of applications but the storage stack itself. ‘In particular, online analysis enables the detection of “misbehaving” situations and offers means to mitigate them. From the user-perspective, non-intrusive monitoring, i.e., without code changes, is key for widespread adoption.’

The goal of this project is to enhance the monitor capabilities in DDN’s IME FUSE client, first by incorporating relevant performance counters into the FUSE module and then by enabling fine-grained online monitoring in the module supporting parallel applications. The extended module will be evaluated on relevant I/O benchmarks.

‘Ultimately, this will lead to a prototype and proof-of-concept system with the extended monitoring capabilities that exhibit substantially better analytic capabilities,’ says Julian. ‘Without the funding from Eurolab4HPC, we wouldn’t be able to support productizing advanced performance analysis for I/O. The close relationship with DDN and the proximity to the market will also stimulate and accelerate our long-term research activities.’

Jean-Thomas adds: ‘Eurolab4HPC is a real opportunity for DDN to share expertise and have access to leading class research in a collaborative framework. We strongly believe this programme is a way to foster a community, to increase cross-fertilization between public and private research and overall to improve European competitiveness.’



High-performance computing (HPC) is a formidable resource, but one which doesn't often reach Europe's small and medium enterprises (SMEs). The SHAPE programme, run by the Partnership for Advanced Computing in Europe (PRACE), aims to fix this, as SHAPE Manager Chris Johnson (EPCC) explains.

Shaping up European companies with HPC

Can small and medium companies really benefit from HPC resources?

Yes – definitely! We have helped over 40 SMEs across Europe to get a foot on the HPC ladder. We recently asked SMEs from completed SHAPE projects for feedback on their experiences and many reported that due to the work carried out in their SHAPE project they were now able to demonstrate a faster time to market for their products, along with an associated reduction in costs and improved sales.

What's your favourite application of HPC in industry?

One of the great things about SHAPE is that as well as supporting projects within traditional HPC areas (engineering, climate science, etc.) we have also had projects from areas less well represented within HPC, such as finance and medicine.

For example, we had one project with an SME that creates cranial orthoses (head supports) via 3D-printing. These designs need testing. Of course, this can be done in a laboratory, but simulation is more efficient. During the SHAPE project, the SME was able to try out a standard open-source HPC package, ESPRESO, which

could then be compared with the well-established commercial package they were already using (ANSYS). This was a success as they were able to show that the open-source package did not compromise results in any way.

Why did PRACE decide to launch the SHAPE programme?

The original aims of the programme were to raise awareness of HPC and to provide European SMEs with the expertise necessary to take advantage of the innovation possibilities created by HPC, thus increasing their competitiveness. These aims appear to have been realized: Most SMEs report that they are continuing to use HPC after the project, even when they had not been using HPC before, and also report an improvement in their business process.

Can you give some examples of the kinds of projects carried out within the SHAPE programme?

The range of projects is quite diverse. For example, a project could involve the porting or parallelization of the company's code to allow it to run on an HPC system, getting an SME up and running with a code already installed on an HPC system, or optimizing an already running code for the company's specific use case.

How can companies get involved in SHAPE?

Calls are issued twice a year; the current one runs until Friday 31 May 2019. We can put SMEs in touch with a PRACE centre to help develop a proposal – the actual application process is very lightweight. During the project the SME works closely with the PRACE partner to complete the required work while allowing the company's staff to gain HPC skills they did not previously have. In addition, SMEs get access to the HPC hardware resources they require. This may be on a conventional HPC system, but in some cases more novel hardware such as GPU nodes or visualization suites may be appropriate.

“SMEs reported a faster time to market for their products, reduced costs and improved sales”



Visualization of the printed cranial orthosis

Image credit: IT4 Innovations / Invent Medical Group

Innovation Europe

PRACE ENTERS ITS SIXTH IMPLEMENTATION PHASE



PRACE, the Partnership for Advanced Computing in Europe, is Europe's key to world-class computing infrastructure. HiPEAC caught up with Stelios Erotokritou (The Cyprus Institute), Marjolein Oorsprong (PRACE aisbl) and Oriol Pineda (Barcelona Supercomputing Center) to find out what the next phase of the project holds.

Why is PRACE so crucial to the European computing systems community, and society in general?

The overarching goal of PRACE is to provide a federated European supercomputing infrastructure that is science driven and globally competitive. It aims to strengthen European science by providing access to high-end computing and data analysis resources, which will help drive discoveries and new developments in all areas of computational science. The goal of these actions is to help create a fertile basis for research, technology development, and industrial competitiveness in Europe.

PRACE is thus crucial, as high-performance computing (HPC) systems are tools that unlock the potential for new scientific breakthroughs and innovations, and PRACE is the European



HPC computing infrastructure which provides access to these resources and related services.

What have the initiative's greatest achievements been so far?

PRACE's greatest achievement is that it is now positioned as the European HPC infrastructure providing a range of HPC resources and services to academia and industry in Europe. These include:

- **Peer review of HPC projects**
PRACE has been running a recognized, competitive and world-class peer review of applications for HPC resources solely based on scientific excellence.
- **Training the next generation of HPC users**
PRACE has established an extensive training programme, with more than 100 training events each year taking place throughout Europe. During the past 10 years, since 2008, PRACE has hosted more than 500 events, training more than 11,500 participants.
- **Maintaining the operation of the European-wide HPC infrastructure**
PRACE has integrated the PRACE European Tier-0 (Europe-wide) and Tier-1 (national) systems. These systems have a common working environment and are connected by a high-speed network so users can use any system and switch between systems without much disruption to their working process, especially with regard to the transfer of their scientific data.
- **Application enabling and support activities**
Over the PRACE-IP projects, PRACE has carried out application enabling activities and has supported users and industry through Preparatory Access, porting and optimization of code and PRACE's SME HPC Adoption Programme in Europe (SHAPE).
- **Dissemination activities**
PRACE has highlighted the importance and achievements of HPC in Europe to the general public – hoping to inspire the next generation of HPC users, academia and industry – encouraging new research fields and companies to use HPC in their research.



What are some of the main trends in HPC that you've seen over the past few years? How do you see these evolving?

A number of trends have been materializing in HPC. These include:

- **urgent computing**, which aims to provide computing resource in the event of emergencies, such as floods or tsunamis
- **virtualization services**, which is developing use cases of containers in the HPC infrastructures
- **high-performance data analytics**: the use of HPC resources for tasks involving sufficient data volumes and algorithmic complexity
- leveraging **artificial intelligence** and **machine learning tools** on HPC

These are expected to develop to help save lives and reduce damage costs, help users in using HPC systems, extract insights from raw information in a more efficient manner and help automate and unlock potential new scientific discoveries, respectively.

Additional HPC trends including evolving hardware and system software/programming environments, as well as a move towards more energy-efficient HPC systems.

How will PRACE develop during PRACE-6IP, the sixth implementation phase?

PRACE-6IP will continue activities from previous projects and build upon these, launching calls for access to Tier-0 HPC systems, organizing training events, application enabling, providing user support and presenting at or organizing conferences – such as the European HPC Summit Week, science festivals and other outreach activities.

New activities within PRACE-6IP will include Calls for Distributed European Computing Initiative (DECI) Tier-1 HPC resources and greater engagement with industry, as PRACE will have a dedicated industry liaison officer. In addition, PRACE

will have various projects which will deliver open-source, forward-looking software solutions addressing the challenges of diversity of hardware and software complexity in the pre-exascale landscape.

The PRACE aisbl association will also be assuming a greater leadership role in PRACE-6IP, acting as a single point of contact for European Union (EU) HPC activities to reduce overlap and redundant actions while strengthening PRACE's role as the main actor in the field.

How do you see PRACE working with EuroHPC?

PRACE and EuroHPC can work together in a number of ways as these two initiatives focus on various aspects of HPC. EuroHPC can build on the established PRACE European HPC infrastructure by incorporating the supercomputing systems to be acquired in future. PRACE can continue and develop its existing activities – such as peer review of HPC projects - including those requesting access to EuroHPC computational resources, training, operation of the European HPC infrastructure, application enabling, and support and dissemination.

Will Europe be able to compete with the United States and Asian countries thanks to initiatives like this?

Contrary to the United States, Japan or China, no single European country can purchase and sustain a top-of-the-range exascale system by its own. The EuroHPC Joint Undertaking will pool European resources to procure two generations of world-class HPC systems and increase the computational capacity of PRACE. This extraordinary investment, together with the rest of PRACE services, will place Europe at the forefront of HPC.

prace-ri.eu

PRACE-6IP has received funding from the European Union's Horizon2020 research and innovation programme under grant agreement number 823767.

EVOLVE: FUSING HPC WITH CLOUD TO EXTRACT MAXIMUM BIG-DATA VALUE



Leading the Big Data Revolution

Launched in December 2018, EVOLVE is a €14 million innovation action comprising 19 key partners from 11 European countries. The project is to introduce important elements of high-performance computing (HPC) and cloud into big data platforms, taking advantage of recent technological advancements to enable cost-effective applications in seven different pilot cases, in order to keep up with the unprecedented data growth we are experiencing.

EVOLVE aims to build a large-scale testbed by integrating technology from:

- The **HPC** world: an advanced computing platform with HPC features and systems software.
- The **big data** world: a versatile big-data processing stack for end-to-end workflows.
- The **cloud** world: ease of deployment, access and use in a shared manner, while addressing data protection.

EVOLVE aims to take concrete and decisive steps in bringing together the big data, HPC, and cloud worlds in a unique testbed that will increase our ability to extract value from massive and demanding datasets. EVOLVE aims to bring the following benefits for processing large and demanding datasets:

- **Performance:** Reduced turnaround time for domain experts, industry (both large companies and small/medium enterprises (SMEs)), and end users.
- **Experts:** Increased productivity when designing new products and services, thanks to the ability to process large datasets.
- **Businesses:** Reduced capital and operational costs for acquiring and maintaining computing infrastructure.
- **Society:** Accelerated innovation via faster design and deployment of innovative services that unleash creativity.

EVOLVE intends to build and demonstrate the proposed testbed with real-life, massive datasets from demanding application areas. To realize this vision, EVOLVE follows a lean approach. EVOLVE will deliver three fully operational prototypes with an expanding set of features to support both general-purpose and accelerated capabilities for computing, data processing and hierarchical storage.



EVOLVE brings together technology and pilot partners from European Union (EU) industry with demonstrated experience, established markets, and vested interest. Furthermore, EVOLVE will conduct a set of proof-of-concepts with stakeholders from the big data value chain to build up digital ecosystems and achieve broader market penetration.

NAME: EVOLVE: HPC and Cloud Enhanced Testbed for Extracting Value from Large-scale Diverse Data

START/END DATE: 01/12/2018-30/11/2021

KEY THEMES: HPC, cloud computing, big data applications, data analytics, innovation

COORDINATOR: Jean-Thomas Aquaviva, DataDirect Networks (DDN), France

[✉ jtacquaviva@ddn.com](mailto:jtacquaviva@ddn.com)

PARTNERS: France: DDN (coordinator), Bull, Thales, Cybeletech; Ireland: IBM; Greece: Foundation for Research and Technology – Hellas (FORTH), Institute of Communication and Computer Systems (ICCS), Space Hellas; UK (Gibraltar): OnApp; Norway: MemoScale; Austria: webLyzard, Virtual Vehicle, AVL; Portugal: Globaz; Luxembourg: Neurocom; Italy: MemEx, Tiemme; Germany: Bayerische Motoren Werke (BMW); Bosnia and Herzegovina: Koola

BUDGET: €14 million

WEBSITE: evolve-h2020.eu

EVOLVE has received funding from the European Union's Horizon H2020 research and innovation programme under grant agreement no. 825061.

ENERGY, SECURITY AND TIME-CRITICALITY GIVEN FIRST PLACE WITH TEAMPLAY



TEAMPLAY

Imagine if you could get hours' more use out of your smartphone before needing to charge it again, or if it differentiated between games and banking apps so that games ran fast while banking apps were run securely to defend them against hackers. These are some of the ambitious aims of the European Union Horizon 2020 project TeamPlay, which was launched in January 2018.

As outlined in the HiPEAC Vision 2019, energy efficiency in computing systems is becoming increasingly important, as mobile applications, the internet of things and cyber-physical systems become more and more prevalent in everyday life. Other properties of our software should not be compromised, such as the degree of security or the ability to react in a timely fashion.

However, such concerns – referred to as 'non-functional properties' – are often treated as secondary compared with performance. Software designers need tools that help them optimize performance while meeting application energy constraints, reaction times and security requirements.

TeamPlay is taking a radically new approach that will enable energy, time and security transparency at the source-code level, turning energy usage, time, security and other non-functional program properties into first-class citizens during the design phase. Application programmers will be able to manipulate and analyse energy, time and security as normal program requirements thanks to analytical and optimization frameworks.

Real-life industry use cases

TeamPlay is applying its tools to real-life use cases provided by our industrial partners, working in domains where energy usage, time and security are critical factors such as computer vision, satellites, drones, medical applications and cybersecurity. As an example, TeamPlay technologies are being tested on drones that perform search and rescue, which need to lower energy usage in order to increase flight time but without compromising security or the reaction time of the flight controller.



TeamPlay partners being shown round the University of Southern Denmark's Unmanned Aerial Systems (UAS) Center



A more sustainable computing paradigm

Information and communication technology (ICT) is estimated to represent an alarming 10% of the world's energy consumption – 50% more energy than the aviation sector. By making energy, time and security first-class system design goals, TeamPlay will pave the way for low-energy, high-performance and high-security products and services – including smartphone apps, web browsers or GPS systems in cars – to become a reality.

At the same time, TeamPlay will make it possible for developers to explicitly prioritize security or reaction time over other properties, ensuring that these products and services provide new levels of security and remain responsive in a wide range of scenarios.

Today we purchase light bulbs, kitchen white goods, and cars based on energy ratings. In future, consumers should also be able to demand low-power, highly secure, and high-performance ICT products. Ultimately, TeamPlay aims to lay the foundation for Europe to lead the low-power, high-performance and high-security computing revolution.

NAME: TeamPlay: Time, Energy and Security Analysis for Multi-/Many-core heterogeneous PLATforms

START/END DATE: 01/01/2018-31/12/2021

KEY THEMES: time-criticality, energy, security, non-functional properties, parallel software

COORDINATOR: Olivier Zendra, INRIA, France

✉ Olivier.Zendra@inria.fr

PARTNERS: France: INRIA (coordinator), Secure-IC; Spain: Thales Alenia Space; Germany: AbsInt, Technische Universität Hamburg; UK: University of Bristol, University of St Andrews; Denmark: Sky Watch, Syddansk Universitet; Netherlands: Universiteit van Amsterdam

BUDGET: €5.42 million

WEBSITE: teamplay-h2020.eu

TeamPlay has received funding from the European Union's Horizon H2020 research and innovation programme under grant agreement no. 779882

EXPLOITING EACH CHIP TO ITS FULL POTENTIAL: HOW UNISERVER OVERCOMES ENERGY SCALING LIMITS IN COMMODITY SERVERS



Imperfections in today's electronic chip manufacturing processes lead to substantial variations of critical device parameters, which are expected to further worsen as transistor dimensions reach the atomic scale. Such increased variability causes otherwise identical nanoscale circuits to exhibit different performance or power-consumption behaviours, even though they are produced using the same processes.

Currently, manufacturers try to deal with the huge performance and power variability in manufactured chips – and hide it from the software layers – by adopting pessimistic safety timing and voltage margins in accordance with rare worst-case scenarios that are assumed during the design phase. The consequence is that all chips are artificially constrained to operate at the low speed and high power consumption of the outliers – i.e. the relatively few worst-case chips – and not according to their true capabilities.

Saving energy one chip at a time

If, instead, we were able to harness the differences in variability between circuits and use each to its maximum potential, we could improve performance and energy efficiency. This is the fundamental but radical idea behind UniServer, a Horizon-2020 project that is developing a complete system stack for servers with mechanisms to reveal and exploit the true capabilities of every product without compromising their dependability.

Over the past three years, the consortium has developed unique automated online and offline characterization processes and diagnostic 'daemons'. These collect and analyse various parameters about the hardware components as they operate, and have revealed extensive margins within multicore servers.



UniServer XGene2

Our results show that some cores could use 10% less than the nominal supply voltage advised by the manufacturer, leading to up-to 38% power savings. Similarly, for dynamic access random memory chips (DRAM) the refresh rate and supply voltage could be decreased by 98% and 5% from nominal levels, leading to average power savings of more than 22% across a range of benchmarks and temperatures.

The revealed capabilities of each hardware core and DRAM chip are being exploited by an enhanced error-resilient software stack that includes new task monitoring and allocation mechanisms in the virtualization layer (QEMU-KVM hypervisor) and OpenStack resource manager for improving the energy efficiency, while maintaining high levels of system availability.

The developed system stack is ported to the world's first 64-bit Arm based server-on-chip family built by AppliedMicro/Ampere in its current and forthcoming X-Genie platforms. A number of prototypes were evaluated using classical cloud and new emerging edge applications, such as detection of malicious network jamming devices and smart advertisement. Results show that server energy savings can range from 12% to 21%. The project results led to more than 35 papers in top-tier venues, a patent and several distinctions, and were disseminated through numerous tutorials and workshops.

UniServer's technologies can support the upcoming transformation of the internet by providing the energy efficiency and performance boosts needed to make fog / edge computing a reality, while supporting classical cloud and HPC ecosystems.

NAME: UniServer- Universal Micro-Server EcoSystem by Exceeding the Energy and Performance Scaling Boundaries

START/END DATE: 01/02/2016-31/07/2019

KEY THEMES: energy efficiency, variability, edge and cloud computing

COORDINATOR: Georgios Karakonstantis, Queen's University Belfast, United Kingdom

✉ G.Karakonstantis@qub.ac.uk

PARTNERS: **UK:** Queen's University Belfast; **Cyprus:** University of Cyprus, Meritorius Audit Ltd; **Greece:** University of Athens, University of Thessaly; **Germany:** Applied Micro Circuits Corporation Deutschland GmbH; **UK:** ARM Holdings; **Ireland:** IBM Ireland Ltd; **Spain:** Worldsensing, Sparsity

BUDGET: €4.8M

WEBSITE: uniserver2020.eu

UniServer has received funding from the EU Horizon2020 programme under grant agreement no.688540



With so many competing claims on our time, it can be hard to find time to review papers. Here, Ben Juurlink (Technische Universität Berlin) makes the case for carefully considering peer-review requests, in order to help keep a proven system going.

Is the peer review system broken?

We live in very busy times. As an example, I am a university professor, and modern professors aren't just expected to simply perform research and teach. They are also expected to acquire funding from all kinds of sources, to manage the projects acquired, to present their research to a general audience, to develop professional videos of the courses they teach, to open up their research data, to co-develop the strategy of the faculty, to have performance appraisal interviews with their team members, to reply to at least 100 emails per day, and so on and so forth.

One of the things I do in addition to all of the things above (with the help of my team members) is that I'm editor of a journal: the Elsevier Journal on Microprocessors and Microsystems – Embedded Hardware Design (MICPRO for short). My role is to find suitable reviewers for the papers assigned to me, to evaluate their reviews, and to make a decision (accept, reject, revise) based on these reviews.

I started my PhD in 1992 and have been in academia ever since. Over this time I have got to know a relatively large number of people in the field of embedded computer architecture and beyond. Given this experience, I consider it my duty to help young researchers find expert reviewers for their manuscripts.

That's where the problem lies. Since we are all very busy, no one can find the time to do these reviews. In recent years it has been my experience that, roughly estimated:

- about 50% of the review requests are neither accepted nor declined
- about 30% of the review requests are declined without giving a reason or alternative reviewers
- approximately 10% of the requests are declined due to work overload but without suggesting alternative reviewers.
- only 10% or so are accepted

“Sometimes I feel that the entire peer-review system is about to collapse. But peer review is a very good system”



Photo credit: © Lamai Prasitsuwani | Dreamstime.com

Occasionally, reviewers accept but then simply don't complete the reviews, while in some cases reviewers review the original manuscript but then decline to review the revision.

As a result, my life as an editor has become very hard. Sometimes I feel that the entire peer-review system is about to collapse. But peer review is a very good system: our papers are evaluated by our peers, by experts like us. By the pigeon-hole-principle, for the peer-review system to work, everybody who tries to publish papers should do as many reviews as the number of manuscripts (s)he submits. Here I assume for simplicity that the average number of co-authors equals the average number of reviews for a paper.

To prevent the peer-review system from collapsing, I appeal to you, dear member of the HiPEAC community. If you receive a review request, please do not ignore it but take it seriously. It is perfectly understandable that you are too busy at certain times and that you have to decline, but then please suggest alternative reviewers, such as team members. The community can only benefit as a result.

Find out why Ben thinks about working in embedded systems and how you can get high performance at low power in our interviews with him, available on the HiPEAC experts webpage:

🔗 hipeac.net/press/#/videos



While processor benchmarking is necessary, standardization is important so that we can ensure that the results are fair and reproducible, argues Vincent Hindriksen (Stream HPC).

Standardized benchmarking in research papers

In high-performance computing (HPC)-related research, a good proportion of papers are about new and faster algorithms for specific processors, showing benchmark results that should convince the reader there has indeed been an improvement. Unfortunately, there is no standard of what is expected from this type of paper, giving too many papers too low a score on reproducibility. When this happens, it gives an incomplete view of the problems researched, making the paper look like a white paper from the processor vendor.

So what should you take into account when writing a benchmark-focused paper? Below, I've outlined a few things which I think should be essential for this task.

How to write a benchmark-focused research paper

- **Benchmark on multiple processors.**

There are various ways to benchmark on recent processors. For graphics processing units (GPUs) there are a few 'HPC zoos' around, including the one produced by my company Stream HPC, and these can be used for non-commercial purposes. Furthermore, there are research clouds, public clouds (especially for field-programmable gate arrays (FPGAs) and high-end GPUs) and of course the hardware of other research groups. Bear in mind that you will need to reserve access in advance.

- **Use multiple metrics** – floating point operations per second (FLOPS), FLOPS per watt, FLOPS/costs, gigabytes per second throughput, latency, etc. On technical websites we often get a better overview when hardware is benchmarked, as multiple metrics are used, showing a winner per metric. This is in contrast

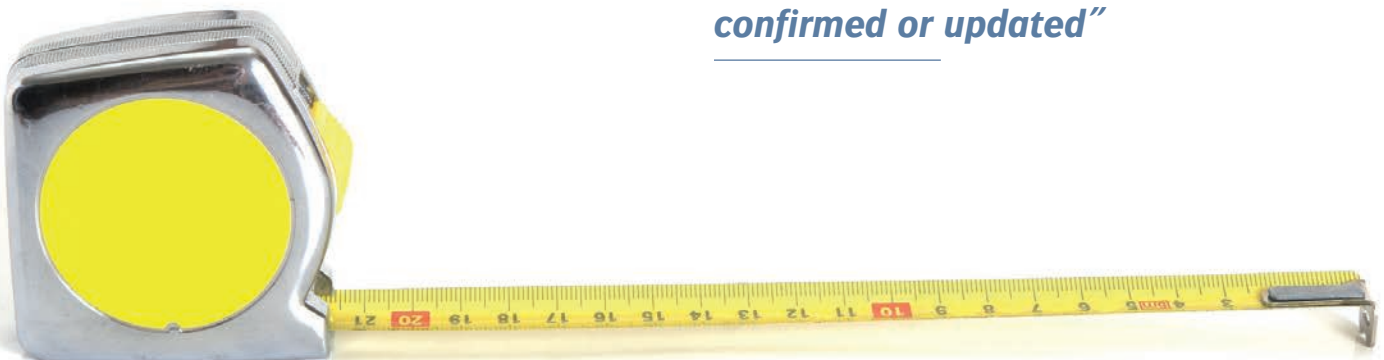
with research papers where too often only a single metric is used, and a single hardware winner is shown.

- **Set a clear context for the experiments.** First, explain why specific processors and hardware have been chosen, and why other hardware has not been included. Second, explain which data and experiment sizes the focus is on. Third, discuss which metrics are only shown but not used, and why. The reader should understand these decisions in full.
- **Allow reproducibility.** Everybody can make a mistake, and by allowing others to reproduce your results, you make sure the research results can be confirmed or updated. This takes away doubts about code not being fully optimized for all processors. Using a DOI code, you can add links to code and full benchmark results – go to doi.org for more information and to see if your university is already a member. Also, consider publishing benchmark results on a Wiki or in GitHub, so others can add their own results.

What would this change?

It's not the responsibility of researchers to benchmark all processors, but neither it is a researcher's responsibility to promote the processor of a particular sponsor. The reasons given (lack of access to recent hardware, lack of time and lack of interest) can often be solved by sharing code.

“By allowing others to reproduce your results, you make sure the research results can be confirmed or updated”





This article presents ReFiRe (Remote Fine-grained Reconfigurable acceleration), a framework that enables the efficient offloading of tasks to remove hardware reconfigurable accelerators. Dionisios Pnevmatikatos of the Technical University of Crete and Foundation for Research and Technology- Hellas (FORTH) explains how ReFiRe outperforms software defined system-on-chip-generated systems by increasing the computation-synchronization ratio per application, rather per accelerator, and by encapsulating reconfiguration requests and dynamic accelerator chaining to construct larger pipelines.

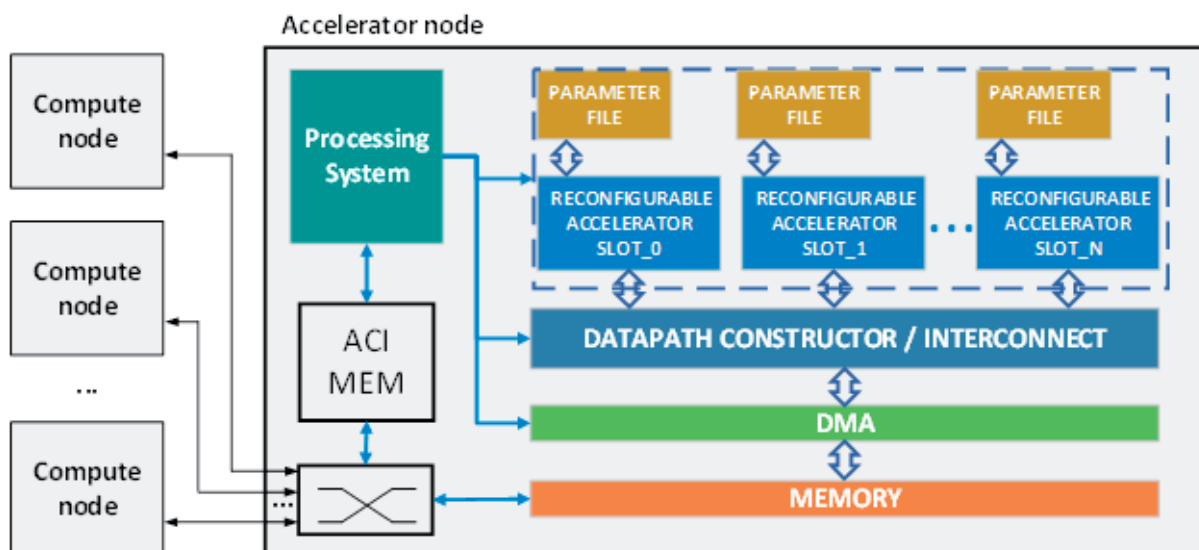
ReFiRe: efficient deployment of Remote Fine-grained Reconfigurable accelerators

The need for specialized hardware acceleration in today's computing platforms is intensified by an insatiable demand for compute power. In the upcoming disaggregated computing environments, where all data transfers between remote nodes are realized via packet exchanges over a rack-scale network, reducing communication and synchronization is a prerequisite to the effective employment of remote acceleration.

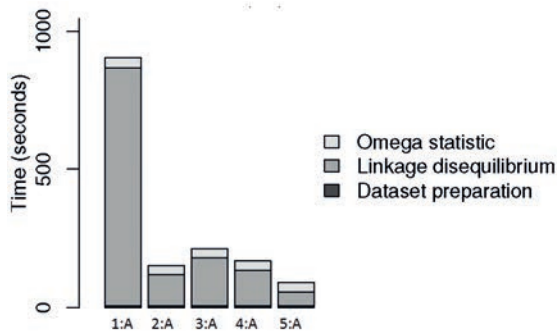
reduced. This is achieved by off-loading control flow and partial reconfiguration decisions to the remote side through arbitrarily long instructions that encapsulate complex sequences of operations and their respective synchronization requirements. One or more compute nodes can issue ReFiRe instructions, which we dub advanced co-processor instructions (ACIs), to a single ReFiRe-enabled hardware platform.

ReFiRe is a generic deployment framework with support for partial reconfiguration that allows communication needs between a processor and remote accelerators to be considerably

reduced. This is achieved by off-loading control flow and partial reconfiguration decisions to the remote side through arbitrarily long instructions that encapsulate complex sequences of operations and their respective synchronization requirements. One or more compute nodes can issue ReFiRe instructions, which we dub advanced co-processor instructions (ACIs), to a single ReFiRe-enabled hardware platform.



The ReFiRe architecture



Experimental results

Communication between accelerators is facilitated by the so-called datapath constructor, which synchronizes the exchange of data between cores. Dedicated on-chip memory, the parameter file, allows per-task configuration parameters to be passed on to the accelerators. External memory is only accessed by the first and last accelerator cores in the datapath, via dedicated direct memory access (DMA) engines.

To evaluate ReFiRe, we tested different accelerator execution scenarios for a population genetics application that detects positive selection. We employed the open-source software OmegaPlus as reference, and introduced the required ReFiRe application programming interface (API) routines in the source code to accelerate linkage disequilibrium (LD) computations, which quantify the non-random association between mutations in the genomes. One accelerator core performs pairwise population counter operations on arbitrarily long binary vectors that represent mutations, while another calculates the squared Pearson's correlation coefficient r^2 (a commonly used measure of LD) between two vectors based on the output values of the previous accelerator.

We measured the execution times per application stage for five different accelerator-deployment scenarios:

1. **baseline reference** (software)
2. **software-defined system-on-chip (SDSoC) system** (accelerators are local to the host)
3. **ReFiRe system** (one ACI per call to a remote accelerator)
4. **ReFiRe system** (one ACI per 32x32 group of LD scores, performing accelerator chaining)
5. **ReFiRe system** (one ACI per LD task, performing accelerator chaining and iteratively invoking the constructed pipeline)

All runs analyse a total of 20,000 binary vectors of 8,000 bits each. The total numbers of LD and ω scores are 1,910, 589 and 167,655, respectively.

The graph to the left shows execution times for the three main stages of OmegaPlus for several execution configurations. The peak throughput performance that the employed accelerator cores can achieve for LD (population counter and r^2 calculator) is $2,809 \times 10^3$ scores/second. Increasing the computation-to-synchronization ratio for these accelerators by creating deeper pipelines (configuration 4) and offloading considerably more computations to the remote node through loops (configuration 5) outperforms an SDSoC-generated implementation that employs the same accelerator cores locally by up to 2.2 times, achieving $2,430 \times 10^3$ scores/second, thus enabling near-peak accelerator performance at the application level.

To conclude, an ACI introduces several degrees of flexibility in employing fine-grained accelerators, exploiting five basic specialization principles as outlined in the second paper cited below – that is, concurrency, computation, communication, caching, and coordination. ReFiRe shifts all control flow and partial reconfiguration decisions to the remote side, that is, near the accelerators, thus considerably reducing the frequency of synchronization events between the specialized hardware and a host processor. This enables near-peak accelerator performance at the application level, despite performing computations on remote nodes.

This work was supported in part by the dReDBox project, which has received funding from the European Union's Horizon2020 research and innovation programme under grant agreement number 687632.

FURTHER READING:

Alachiotis, Nikolaos, Alexandros Stamatakis, and Pavlos Pavlidis. 'OmegaPlus: a scalable tool for rapid detection of selective sweeps in whole-genome datasets.' *Bioinformatics* 28.17 (2012): 2274-2275

Nowatzki, Tony, et al. 'Domain specialization is generally unnecessary for accelerators.' *IEEE Micro* 37.3 (2017): 40-50

“Reducing communication and synchronization is a prerequisite to the effective employment of remote acceleration”

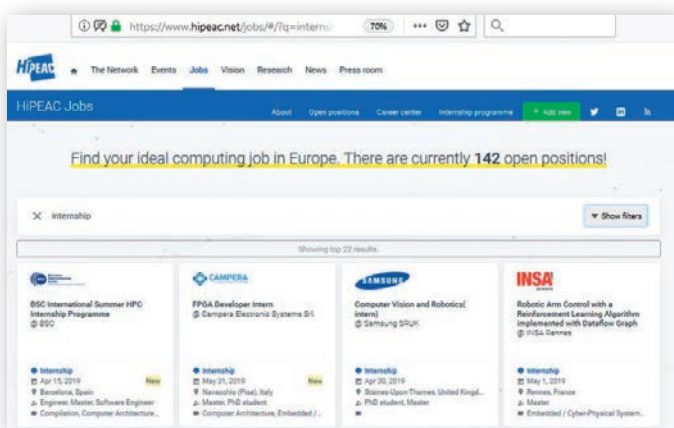
HiPEAC Internship Programme

HiPEAC mobility has been one of most successful instruments for structuring and connecting our community over the years. Since 2006, more than 300 PhD students have taken part in HiPEAC collaboration grants and internships. Many of the students who received a mobility grant while they were doing a PhD are now excellent senior researchers, successful entrepreneurs, or hold leading engineering positions at the most well-known companies in the field, such as Arm, Intel, Samsung, etc.

Over the last year, HiPEAC has been working to make the internship process more user-friendly and agile. Internships are now integrated into the HiPEAC Jobs portal, which has a number of advantages:

- Internships are available throughout the year. Companies can upload opportunities any time, instead of waiting for a specific call to open.
- The process is more flexible and tailored to the different needs of each company. Companies just need to indicate how they expect students to contact them and can follow their own recruitment processes.
- HiPEAC will promote the internships among the students within our network throughout year, without deadlines. We will also present them at different events and via different activities, such as the HiPEAC Jobs wall, our 'Inspiring Futures' careers sessions, presentations, roadshow booths, science, technology, engineering and mathematics (STEM) student days, etc.
- Students can apply at any point during the year and agree a time that suits them and their supervisor to do the internship with the company.

HiPEAC will continue to provide funding for internships at small/medium enterprises, while large companies will get other benefits, such as a free pass for the HiPEAC conference or ACACES summer school.



How does it work for companies?

1. Companies upload the internship to the HiPEAC Jobs portal. Candidates apply for the position as indicated in the description (by email, via the company's website, etc.). Indicate clearly the necessary information they will need to provide (CV's, recommendation letters, etc.) and if there are any restrictions regarding visa requirements, timing of the internship, etc.
2. Select the candidate following the company's own procedures and best practices (such as a phone interview or face-to-face interview)
3. Confirm the intern by indicating the name and contact on the evaluation form in the manage jobs section. HiPEAC will get in touch to make sure you receive funding or other benefits.

All HiPEAC companies are eligible to submit any number of internship proposals, provided they are based in Europe and that they take place during 2019 (and end by 28 February 2020 at the latest).

How does it work for students?

1. Find your ideal placement at the HiPEAC Jobs portal. Candidates apply for the position as indicated in the internship description (email, via the company's website, etc.).
2. Once the company has notified HiPEAC, we will contact the selected student to finalize arrangements for the internship, if required. Students will need to have an account on the HiPEAC website (hipeac.net/accounts/signup) to complete the process.
3. Selected students will need to submit a short summary of the internship after completing it, to be published on the student's HiPEAC profile

Eligibility for funding:

To be eligible for the HiPEAC Internship programme:

- Students must be currently enrolled at a European institution (whether at postgraduate or undergraduate level).
- The internship placement must involve movement to a different location.

FURTHER INFORMATION:

HiPEAC Jobs portal

hipeac.net/jobs

Internship information

hipeac.net/jobs/#/internships

recruitment@hipeac.net

Follow us: @HiPEACjobs



Career talk: Patience Masoso, Zimbabwe Centre for High Performance Computing (ZCHPC)



Can you tell me about your background?

How did you get into high-performance computing (HPC)?

I'm an HPC applications engineer at the ZCHPC, with a BSc in information systems from Zimbabwe's Midlands State University. I'm passionate about familiarizing myself with computer systems and adding value to information technology through research and system

improvement. During my undergraduate programme I did two projects on system development, and ended up specializing in HPC after I realized that it promises a wide range of opportunities, as it's constantly changing. The field of HPC is not limited to information technology but incorporates various domains.

Currently, I'm head of the life sciences department at our institution, and I'm participating in life science projects being done in our country, which has given me the chance to learn things besides information technology. But I've got experience from different areas: during my work-related learning, I worked at a government ministry where I interacted with different users and worked in a support role. I've also worked as an information technology teacher at a high school before coming to ZCHPC, and another thing I'm passionate about is educating undergraduate students in high-performance computing.

"Inherently, technology focuses on breaking boundaries and improving people's lives"

All this has offered me the opportunity to work with people at different levels and give them advice on how to improve their work using information technology. Inherently, technology focuses on breaking boundaries and improving people's lives. If I get the opportunity, I would like to gain a master's in HPC and study for a PhD in areas related to HPC and data-intensive computing.

When was ZCHPC established?

ZCHPC was conceived in 2011; the implementation began in May 2014 and took less than a year, with the centre opening in 2015. The system has a theoretical peak performance of 36 teraFLOPS.

Zimbabwe is now officially the fourth country in Africa to have HPC. HPC was identified as one of the key solutions to assisting the nation and the Southern African Development Community (SADC) region to solve issues such as climate change, food security, poverty, disease, energy and human capital development problems through advanced research.

What are some of the main projects ongoing at the moment?

ZCHPC has a portfolio of projects in different domains, including for prominent sub-Saharan diseases and to ensure food security. It is mandated to ensure that organizations use HPC as an essential part of their business strategy, helping them to design new products and solve production problems in order to become more innovative and competitive in support of industry.

The main projects in which HPC is being used in Zimbabwe are as follows:

- Weather forecasting: the system is running collective calculations, predicting seismic occurrences, building complicated climate models and forecasting weather accurately. The Meteorological Service Department (MSD) and higher learning institutions are making use of HPC in this regard. MSD uses HPC to carry out the Zimbabwe weather forecast on daily basis. They have managed to improve the resolution of weather forecasts as well as reduce the time needed to generate daily forecasts. Researchers in the domain of agriculture use the system to make predictions.

- Computer-aided identification, modification and validation of anti-schistosomiasis compounds and their targets. This project seeks to come up with new drug compounds and novel drug targets that can be used in the discovery and development of drugs for schistosomiasis.
- Alien invasive species DNA barcoding: ZCHPC is collaborating with various organizations in a DNA barcoding project. The project seeks to come up with a virtual biobank with genetic information of alien species for the whole of the SADC region. ZCHPC will provide storage for the database.
- Integrating cancer bioinformatics, systems biology and cheminformatics to design novel drugs.
- Developing potent new plasmodium falciparum PI3K (PfPI3K) inhibitors as anti-malarial drugs.

What is the field of HPC like in Zimbabwe?

The field of HPC in Zimbabwe is growing, as more research scientists have incorporated HPC to improve their tasks. The centre has managed to introduce training programmes for users in different sectors, including the private sector and academic institutions.

HPC has been used for civil protection and geoscience through the Scientific and Industrial Research and Development Centre (SIRDC). The centre had a collaboration with SIRDC and managed to develop a flood map for the Muzarabani area for civil protection purposes in the event of national disasters using the HPC system. This saved lives and mitigated flood risks through the use of HPC to predict natural disasters.

However, there is still a need for the centre to raise awareness on how HPC can be used by undergraduate students in doing their projects.



The official opening of ZCHPC in 2015



Harare, where ZCHPC is located

HPC is being utilized in the following areas:

- **Zimbabwean universities:** HPC is being used to carry out advanced scientific research within tertiary institutions. Academia has the potential to develop different products that will have benefits in manufacturing, automobile design, patient specific medicine and infrastructural design. To date, tertiary institutions have received training in agriculture, health, industry and commerce, the latter with the aim of generating revenue through HPC usage.
- **Human capacity development:** the HPC system is currently being used as a teaching aid to equip researchers in Zimbabwe with technological skills in artificial intelligence and machine learning, weather forecasting and climate modelling, bioinformatics, computational chemistry and finance.

What plans do you have for the future?

ZCHPC thrives on being a service provider for computation in business intelligence, both in industry and commercial companies. It aims at increasing the capacity utilization of the system by 50%. Its mandate is to be part of flagship projects in Zimbabwe, therefore it is working on increasing the number of industrial and academic researchers by at least six per year. The centre also aims to introduce modules related to HPC to university curricula in Zimbabwe.

“Zimbabwe is now officially the fourth country in Africa to have HPC”



So you're in the final year of your computer science or engineering degree, and you're not sure whether you want to start working or embark upon a PhD? Pedro Trancoso, who is directing a new master's in high-performance computer systems at Chalmers University, told us why he thinks a master's degree is a great way to develop your skills.

Mastering high-performance computer systems

Why should students specialize in high-performance computing (HPC)?

HPC is no longer a niche topic for a few people – it's now widely used by industry and business for the development and deployment of novel products. It responds to needs such as processing large amounts of data (data science), doing complex running simulations (weather, finance, and so on) or taking decisions quickly in autonomous environments (such as self-driving cars or drones). Some of the most exciting emerging technologies, including machine learning and quantum computing, also rely on efficient high-performance systems.

In the area of high-performance computer systems, the traditional approach was to work on separate areas and improve them independently. However, now we're reaching certain limits – such as the end of Moore's Law – we need new approaches that are more holistic and domain specific, which requires different skills.

What are the benefits of a master's degree?

A master's degree gives students exposure to many of the latest topics in research developments; by the end, they should have a much clearer view of what's going on in the area, be able to use the different technologies available and understand the trade-offs between existing and future systems.

At Chalmers, students have the opportunity of doing a master's thesis within industry, meaning that they get a chance to apply the material presented in class to real case studies.

What does the master's in high-performance computer systems offer?

While we've known for some time how to develop and exploit large supercomputers, we now want to offer the same computational capacity to smaller and autonomous systems, which is a considerable challenge. With this programme, we aim to prepare engineers to develop and exploit new high-performance systems under tight constraints, such as energy consumption and real time. What marks this course out is that it combines different skills in the areas of computer systems, computer graphics, and real-time systems, allowing students to look at a system as a whole and exploit the interactions and optimizations across the different layers of hardware and software. Students also get to develop their value creation skills through the entrepreneurial track.

Classes are taught by world-class teachers and researchers, many of whom have received prestigious grants and participate in major European projects such as the European Processor Initiative featured in this magazine (see p.12).

You've participated in the HiPEAC Student Challenge. Why is it important for students to take part in extra-curricular activities like this?

These activities are so important they should almost be mandatory! We're increasingly seeing that successful professionals require not only a solid technical background, but also 'soft skills' like communication, team work, entrepreneurship and so on – skills which these challenges help develop. Student challenges are also usually excellent opportunities to develop original problem solving strategies, which are also an extremely important aspect of an engineer's profile.

FURTHER INFORMATION:

bit.ly/Chalmers_Masters_HPCS



We've revamped the HiPEAC internship programme: as set out on p. 38, now companies can post internships and students can apply all year round via the HiPEAC Jobs portal. For further information, visit the HiPEAC website: hipeac.net/jobs/#/internships. During his internship at Belgian company Embedded Computing Specialists, Spanish student Daniel Báscones worked on a Verilog core able to run the RISC-V 'I' instruction set.

HiPEAC internships: your career starts here



NAME: Daniel Báscones
RESEARCH CENTER: Universidad Complutense de Madrid
HOST COMPANY: Embedded Computing Specialists
DATE OF INTERNSHIP: 22/09/2018 – 22/12/2019

RISC-ing everything: Developing a Verilog core in a HiPEAC internship



Embedded
Computing
Specialists

My internship at Embedded Computing Specialists focused on the development of a Verilog core capable of running the RISC-V 'I' instruction set. Along with this base, the extensions 'C' and 'E' were also developed. The core was first synthesized and simulated using the Vivado program from Xilinx. After its functionality had been tested, it was then implemented on a MicroZed board containing a ZYNQ 7010 field-programmable gate array (FPGA) and further tested with a suite of benchmarks.

Compilation was done using the RISC-V-GCC toolchain. Custom initialization routes were developed in C and assembly for the specific linker scripts. Bash scripts were developed to compile generic C programs for the core, by linking against the custom routines. In addition, a loader was developed to initialize the core's memory dynamically from the computer, as well as controlling the core's functionality. A command-like interface was also created to facilitate core control and memory loading.

Results

Benchmarks were run on the core and the results compared against RISC-V's SPIKE (a golden reference simulator), as well as against the MicroBlaze processor, a well-tested soft-core.

Three versions of MicroBlaze were tested with different pipeline depths: three-, five- and eight-stage. Our core was tested with a five- and six-stage pipeline. Two factors were taken into account when comparing the results: maximum theoretical core frequency and number of cycles to finish each benchmark, which were used to obtain the runtime.

With regard to the MicroBlaze cores, while the frequency was similar for all three versions – as we used the simplest version of the core – the five-stage MicroBlaze had the lowest number of cycles. As for the RISC-V core, while the six-stage core was able to run at a higher frequency, it also had a higher number of cycles than the five-stage version. The runtime was thus similar for both versions. RISC-V achieved an overall lower runtime than MicroBlaze across different compiler optimization levels, proving that the RISC-V instruction set architecture can be more efficient than that of MicroBlaze.



Embedded Computing Specialists co-founder Philippe Manet commented: 'This internship came about thanks to my participation in a HiPEAC "Inspiring Futures" careers session at the ACACES summer school in 2018. I would recommend this to anyone looking for highly qualified interns and recruits – it's a great way to meet students and chat about what you can offer them as an employer.'



Daniel's internship took him from Madrid (left) to Brussels (right)

The HiPEAC network includes almost 1,000 PhD students who are researching the key topics of tomorrow's computing systems. This issue's featured thesis is by Adrián Pérez Diéguez of the University of A Coruña, and focuses on parallel operations on heterogeneous platforms.

Three-minute thesis



NAME: Adrián Pérez Diéguez
RESEARCH CENTER: University of A Coruña
SUPERVISORS: Margarita Amor López and Ramón Doalla Biempica
DATE DEFENDED: 17/01/2019
THESIS TITLE: Parallel prefix operations on heterogeneous platforms

Graphics processing units (GPUs) have shown remarkable advantages with regard to computing performance and energy efficiency, and represent one of the most promising hardware trends for the future of high-performance computing (HPC). However, these devices can be difficult to program, and considerable effort is required to ensure portability between different generations.

Parallel prefix algorithms (fast Fourier transform – FFT, scan primitive, sorting or tridiagonal system solvers) are a set of widely used parallel algorithms. Their efficiency is crucial in many applications, including signal processing, image filtering, facial recognition, video games, database and map/reduce operations. GPUs can accelerate the computation of such algorithms, but if the algorithms don't map correctly to the GPU architecture or fully exploit the GPU's parallelism, they can also pose limitations.

A number of proposals exist to facilitate the programming of these architectures: autotuning, directives – such as OpenACC or hiCuda – automatic compilers, or accelerated libraries. However, all of these have drawbacks. My thesis provides a library that fully exploits the power of current, improving usability and allowing portability.



Specifically, this thesis focuses on two approaches. On the one hand, completely new parallel prefix algorithms were algorithmically designed; although focusing on the GPU computing model, they can be implemented in any other parallel programming paradigm. On the other hand, a GPU tuning methodology is provided for a set of parallel prefix algorithms. Our methodology is able to perform an efficient execution on any CUDA GPU architecture, and its code can be easily updated for future architectures.

To accomplish this, the methodology identifies the GPU parameters which influence performance and, following a set of performance premises, obtains the convenient values of these parameters. To do this, the different algorithms, problem sizes and GPU architectures were taken into consideration.

Three different approaches are proposed, varying in accordance with the size of the dataset. The first two approaches solve small and medium-large datasets, respectively, on a single GPU. Meanwhile, the third approach deals with extremely large datasets on a multiple-GPU environment, where communication latencies have to be taken into consideration. Our methodology has also been tested on CUDA embedded systems, obtaining highly satisfactory results.

As a result, a GPU-tuned library is built considering previous methodology and new algorithms developed. The library is based on a modular design, which allows us to solve different parallel prefix algorithms. The library outperforms the state of the art for most of the parallel prefix operations and problem sizes. To the best of our knowledge, there is not any other GPU library with the same features.

Our library was tested for the multiplication of large integers, which is a highly important operation in cryptography (web security, blockchain etc.), using both a FFT-based approach and a tiling-based parallel-multiplication, which surpassed other GPU approximations.

Many trending real-time applications (such as autonomous driving, blockchain) use parallel prefix operations (for example in convolution, in body recognition or in cryptography), which needs efficient implementations that can be easily extended for future architectures. In these cases, our research can be of great use.



HIPEAC

ACACES 2019
Fiuggi, 14-21 July 2019



Fifteenth international summer school on Advanced Computer Architecture and Compilation for High-Performance and Embedded Systems

Featuring: TETRAMAX innovation track • Eurolab4HPC high-performance computing track • Accelerating machine intelligence along the compute continuum • Fundamental limits on energy consumption • Die stacking • Processor architecture security • ...and much more!

acaces.hipeac.net/2019