MASSIVIZING HIGH PERFORMANCE COMPUTING FOR AI AND ML: VU ON THE SCIENCE, DESIGN, AND ENGINEERING OF AI AND ML ECOSYSTEMS

@Large Research Massivizing Computer Systems



http://atlarge.science

bit.ly/MassivizingHPC22

Massivizing = Rich challenge of computer science \rightarrow high societal impact!



Sponsored by:



Contributions from the MCS/AtLarge teams. Many thanks! Many thanks to our collaborators, international working groups, authors of all images included here. Also to Gianfranco Bilardi for invitation, ScalPerf for discussion.







WE'RE MASSIVIZING COMPUTER SYSTEMS!

VU AMSTERDAM < SCHIPHOL < THE NETHERLANDS < EUROPE





http://atlarge.science





WE ARE HIRING A NEW ASST. PROF.!



Alumni



















WE ARE A FRIENDLY, DIVERSE GROUP, OF DIFFERENT RACES AND ETHNICITIES, GENDERS AND SEXUAL PREFERENCES, AND VIEWS OF CULTURE, POLITICS, AND RELIGION. YOU ARE WELCOME TO JOIN!

WHO AM I? PROF. DR. IR. ALEXANDRU IOSUP

- Education, my courses:
 - > Honours Programme, Computer Org. (BSc)
 - > Distributed Systems, Cloud Computing (MSc)
- Research, 15 years in DistribSys:
 - > Massivizing Computer Systems

• About me:

- > Worked in 7 countries, NL since 2004
- > I like to help... I train people in need
- > VU University Research Chair + Group Chair
- > NL ICT Researcher of the Year
- > NL Higher-Education Teacher of the Year
- > NL Young Royal Academy of Arts & Sciences
- > Knighted in 2020





WE ARE ALIGNED WITH COMMUNITY CONCERNS...

The Manifesto on

Computer Systems and Networking Research Clear vision for the field in the NL, 2021-2035

Signed 50+ PIs / Leads 7 universities 5 relevant societal stakeholders

Available

Full version (40+ pages) https://arxiv.org/pdf/2206.03259 Who's Who in CompSysNL? https://bit.ly/CompSysNLWhosWho

© 2022 Alexandru Iosup. All rights reserved.



ONE PROJECT TO MENTION...

Big Graph Processing: Used in AI/ML, FinTech, Sci/Pharma Industry 4.0, Energy Mgmt.*, etc.

Vision: Massivizing computer systems approaches are key to enable big graph ecosystems

contributed articles



Sakr, Bonifati, Voigt, Iosup, et al. (2021) <u>The Future Is</u> <u>Big Graphs!</u> CACM.

Prodan et al. (2022) The GraphMassivizer project: Towards Extreme and

Sustainable Graph Processing for Urgent Societal Challenges in Europe.



THIS IS THE GOLDEN AGE OF MASSIVE COMPUTER ECOSYSTEMS



THE ECONOMIC IMPACT OF MASSIVE COMPUTER ECOSYSTEMS



DIVERSE SERVICES FOR ALL

EVERY $eqref{eq:expansion} 1 \rightarrow eqref{eq:expansion} 15 \text{ added value}$

Impacting <u>>60%</u> of the NL GDP (1 trillion EUR/y)

Attracting <u>>20%</u> of all foreign direct investments in NL

Sources: Iosup et al., Massivizing Computer Systems, ICDCS 2018 [Online] / Dutch Data Center Association, 2020 [Online] / Growth: NL Gov't, Flexera, Binx 2020. Gartner 2019. IA 2017.



PHENOMENON: FAILURES IN CLOUD SERVICES

UNCOVERING THE PRESENCE OF FAILURES



PHENOMENON: PERFORMANCE IN CLOUD SERVICES

UNCOVERING THE PRESENCE OF PERFORMANCE ISSUES, EVEN LEADING TO CRASHES

Polygon

Source: http://bit.ly/EveOnline21Crash

NEWS

Players in Eve Online broke a world record — and then the game itself

Developers said they're not 'able to predict the server performance in these kinds of situations' By Charlie Hall | @Charlie_L_Hall | Jan 5, 2021, 2:54pm EST



PHENOMENON: CLOUD DATACENTER SUSTAINABILITY

UNCOVERING THE USE OF ENERGY AND WATER, THE IMPACT ON CLIMATE

Power consumption of datacenters: >1% of global electricity

Source: Nature, 2018 [Online]

Power consumption of datacenters in the Netherlands: $1 \rightarrow 3\%$ of national electricity

Source: NRC, 2019 [<u>Online]</u>

Water consumption of datacenters in the US: **>625Bn. l/y** (0,1%)

Source: Energy Technologies Area, 2016 [Online]

Other greenhouse emissions: Largely unknown

Source: Nature Climate Change, 2020 [Online]

Source: NASA Earth Observatory

THIS TALK: MASSIVIZING = LET'S THINK ECOSYSTEMS!

WE TAKE A HOLISTIC VIEW, BASED ON COMPUTER ECOSYSTEMS

Technology not ready, many issues*

A Why does this* happen? R What to do about it*?

* In modern computer systems, issues are often linked.

Source: Alexandru's personal library.



A new science, of complex, smart computer ecosystems

3 (operational simplicity for the <u>user</u>)

AN ANALOGY: MASSIVIZING CLIMATE SCIENCE

TAKE A HOLISTIC VIEW, BASED ON COUPLED NATURAL SYSTEMS

Can be understood only with coupled models

Climatologist Bjorn Stevens. . Source: HEOM/ire

* In climate science, issues are often linked. The same occurs in massive computer (eco)systems.

Structure: composites of smaller assemblies, but for ecosystems some assemblies are produced elsewhere, by teams with different practices

Operation: ecosystems exhibit many unknown, possibly emergent, phenomena, dynamics, and socio-technical issues

Lifecycle: some ecosystem constituents will perish, or be replaced with others that may not actually fulfill the needs

In plain English: in modern computer systems, several or all issues may be linked. Thus, looking at any single issue for an isolated system is no longer sufficient. We are now creating the science that explains this—Massivizing ...



ECOSYSTEM = SERVICES + COMPUTING + SMARTS + GOALS



Extreme Automation, Performance, Dependability, Sustainability

DISTRIBUTED ECOSYSTEMS, OUR DEFINITION

- 1. Set of 2+ constituents, often heterogeneous
- 2. Each constituent is a system or an ecosystem (recursively)
- 3. Constituents are autonomous, cooperative or in competition
- 4. Ecosystem structure and organization ensure responsibility
 - 1. Completing functions and providing services
 - 2. Providing desirable non-functional properties
 - 3. Fulfill agreements with both operators and clients, clients in the loop
- 5. Long and short-term dynamics occur in the ecosystem

Iosup et al., Lecture Notes in Distributed Systems, Section 1.1.1

Iosup et al., Massivizing Computer Systems, ICDCS 2018. [Online]

8 PRINCIPLES OF DISTRIBUTED SYSTEMS & ECOSYSTEMS

P1: The Golden Age

- P2: Design for massive scale, focus on scalability, elasticity
- P3: Phenomena to discover, seed innovation
- P4: Inherent functional requirements
- P5: Non-functional requirements
- P6: Resource management and scheduling
- P7: Interaction programming model system architecture
- P8: Super-distribution

Iosup et al. Distributed Sytems and Ecosystems,





(our theories are often frameworks, designs, etc.)

SERVERLESS AI/ML/DL OPERATIONS



 Λ

SERVERLESS AI/ML/DL OPERATIONS

ISSUES: COMPLEXITY, NON-TECHNICAL

Actual ML app is a very small part!

Rest is systems, HW+SW, including HPC

Adapted from:

Granhs | CACM

Sakr, Bonifati, Voigt, Iosup, et al. (2021) <u>The Future Is Big</u>



HOW TO DO RM&S ACROSS THE ECOSYSTEM?

IT'S OPERATIONS!

REFERENCE VIEW ON OPERATIONAL TECHNIQUES



SERVERLESS AI/ML/DL OPERATIONS



 Λ

HOW TO REPORT PERFORMANCE?



THE COMPLEXITY CHALLENGE

REFERENCE VIEW ON OPERATIONAL METRICS



VU

N. Herbst, E. Van Eyk, C. L. Abad, A. Iosup, et al. (2018) Quantifying Cloud Performance and Dependability: Taxonomy, Metric Design, and Emerging Challenges. TOMPECS 3(4): 19:1-19:36

SERVERLESS AI/ML/DL OPERATIONS





(our practical approaches are often instruments and datasets)

THE CHALLENGE OF REAL-WORLD AND ANALYTICAL APPROACHES

- Testing in the lab costs 6-9 months of effort → real-world experiments are costly
- Analytical approaches rely on limiting assumptions, and on the existence of a (valid, calibrated) model → analytical approaches rarely easy to do for new technology



OPENDC: SIMULATE DATACENTERS, TOGETHER



OpenDC

A Miniature Datacenter inside Your PC (and Beyond)



Georgios Andreadis

info@gandreadis.com LUMC & CWI

Fabian Mastenbroek <u>F.Mastenbroek@atlarge-research.com</u> Delft University of Technology



@Large Research Massivizing Computer Systems





© 2022 Alexandru Iosup. All rights reserved.

... CAN WE AFFORD A? WHAT IF B? X vs. Y ... vs. Z?

TOO COSTLY TO CONDUCT REAL-WORLD EXPERIMENTS, SO USE A SIMULATOR



simulator



Learn more: opendc.org

- Short-term resource management
- Long-term capacity planning
- Sophisticated model

 many Qs, goals
- Supports many kinds of workloads
- Supports many kinds of resources
- Validated for various scenarios
- Work with major NL hoster
- Used in training, education, research



and more...

© 2022 Alexandru Iosup. All rights reserved.





@Large Research Massivizing Computer Systems



IMPACT OF UNDERSPECIFICATION ON PERFORMANCE RESULTS

- From the highly cited article on the scheduler of a Big Tech company
- Selected two underspecified stages
- Compared different credible policies that could have been used



TUDelft

Significant difference in performance!

VRIJE UNIVERSITEIT AMSTERDAM

@Large Research Massivizing Computer Systems

Follow-Up: Automatic Scheduler Design Exploration

- A novel approach for designing better schedulers
- We decompose schedulers into 33 stages
- Systematically construct new schedulers using building blocks
- Trying every combination takes over **10¹⁸ years!**
- A smarter exploration approach based on an automated process of natural selection
- Algorithm is able to adapt scheduler to workload
 ferns



Fabian Mastenbroek

@Large Research Massivizing Computer Systems



Experiments are also Expensive! An Environmental Perspective



Summary: Simulation-based experiments are very useful.

- 1. Simulators are very useful for many lines of inquiry
- Especially exploratory research can benefit □ simulation-based experiments are very efficient
- 3. We are building a datacenter simulator with exceptional capabilities
- 4. Already published results in top-tier conferences



CHALLENGE: MEANINGFUL REAL-WORLD EXPERIMENTS

THE CHALLENGE OF REAL-WORLD AND ANALYTICAL APPROACHES

Objective: analyse performance of (graph-based) AI/ML systems

Analytical modeling Profiling Tracing





© 2022 Alexandru Iosup. All rights reserved.

GradeML: PERFORMANCE MEASUREMENT AND ANALYSIS FOR COMPLEX ML WORKFLOWS

+

HIGHLY EFFICIENT MEASUREMENT, ADVANCED AND VERSATILE ANALYSIS



Also easy to use, low overhead in terms of modeling

© 2022 Alexandru Iosup. All rights reserved.

GradeML: AUTOMATED BOTTLENECK DETECTION AND PERFORMANCE ISSUE IDENTIFICATION FOR AI/ML



GradeML: AUTOMATED BOTTLENECK DETECTION AND PERFORMANCE ISSUE IDENTIFICATION FOR AI/ML







Hegeman et al., Grade10: A Framework for Performance Characterization of Distributed Graph Processing, IEEE Cluster







Hegeman et al., Grade10: A Framework for Performance Characterization of Distributed Graph Processing, IEEE Cluster







CHALLENGE: FALSE SENSE OF KNOWING HOW THINGS WORK

THE CHALLENGE OF TOO MUCH EXPERTISE

An anecdote

- For grid computing, we built them from scratch, so we though we knew all about how they work
- The community leads thought the future is "Big, parallel jobs"
- Whether we actually have parallel jobs or not really matters
- Unfortunately, the belief was not true

Iosup and Epema: Grid Computing Workloads. IEEE Internet Computing 15(2): 19-26 (2011)

48



DISCOVERY = LARGE-SCALE, LONG-TERM STUDY

UNCOVERING THE MYSTERIES OF OUR PHYSICAL UNIVERSE





James Cordes, The Square Kilometer Array, Project Description, 2009 [Online]

The Square Kilometer Array Factsheet, How much will it cost?, 2012 [Online]

Phil Diamond and Rosie Bolton, Life, the Universe & Computing: The story of the SKA Telescope, SC17 Keynote. [Online]



DISCOVERY = LARGE-SCALE, LONG-TERM STUDY

UNCOVERING THE MYSTERIES OF OUR UNIVERSE, PHYSICAL AND DIGITAL







© 2021 Alexandru Iosup. All rights reserved.

TAKE-HOME MESSAGE



Massivizing \rightarrow computer ecosystems with good functional and non-functional properties, for all

The ecosystem is vast, challenging: many apps, many platforms, may goals, many approaches

Many modern, open challenges: resource management and scheduling, telemetry, analysis, simulation, experimentation, etc.

https://atlarge-research.com/publications.html



- 1. Iosup et al. Massivizing Computer Systems. ICDCS 2018 ← start here
- 2. Andreadis et al. A Reference Architecture for Datacenter Scheduling, SC18
- Van Eyk et al. Serverless is More: From PaaS to Present Cloud Computing, IEEE IC Sep/Oct 2018
- 4. Uta et al. Exploring HPC and Big Data Convergence: A Graph Processing Study on Intel Knights Landing, IEEE Cluster 2018
- 5. Talluri et al. Big Data Storage Workload in the Cloud. ACM/SPEC ICPE 2019.
- 6. Toader et al. Graphless. IEEE ISPDC'19.
- 7. Jiang et al. Mirror. CCPE 2018.

FURTHER READING

- 8. Ilyushkin et al. Autoscalers. TOMPECS 2018.
- 9. Versluis et al. Autoscaling Workflows. CCGRID'18.
- 10. Uta et al. Elasticity in Graph Analytics? IEEE Cluster 2018.

- 11. Herbst et al. Ready for rain? TOMPECS 2018.
- 12. Guo et al. Streaming Graph-partitioning. JPDC'18.
- 13. Iosup et al. The OpenDC Vision. ISPDC'17.
- 14. Iosup et al. Self-Aware Computing Systems book.
- 15. Iosup et al. LDBC Graphalytics. PVLDB 2016.

Etc.

FURTHER READING

https://atlarge-research.com/publications.html



- 1. Crusoe, Iosup, et al. (2022) Methods Included: CWL. CACM
- 2. Sakr, Bonifati, Voigt, Iosup, et al. (2021) The Future Is Big Graphs! CACM
- 3. Andreadis et al. (2021) Capelin: Data-Driven Capacity Procurement for Cloud Datacenters using Portfolios of Scenarios. TPDS, under review.
- 4. Versluis et al. The Workflow Trace Archive: Open-Access Data From Public and Private Computing Infrastructures. TPDS 2020.
- 5. Uta et al. (2020) Beneath the SURFace: An MRI-like View into the Life of a 21st-Century Datacenter. login USENIX
- 6. losup, Hegeman, et al. (2020) The LDBC Graphalytics Benchmark. CoRR. <u>https://arxiv.org/abs/2011.15028</u>
- 7. Hegeman et al. (2021) GradeML. HotCloudPerf.

 Abad, Iosup, et al. An Analysis of Distributed Systems Syllabi With a Focus on Performance-Related Topics. WEPPE 2021. <u>https://arxiv.org/abs/2103.01858</u>
 Etc.

FURTHER READING

https://atlarge-research.com/publications.html



- Iosup et al. The AtLarge Vision on the Design of Distributed Systems and Ecosystems. ICDCS 2019 ← Start here
- 2. Uta et al. Is big data performance reproducible in modern cloud networks? NSDI 2020
- 3. Van Eyk et al. The SPEC-RG Reference Architecture for FaaS: From Microservices and Containers to Serverless Platforms, IEEE IC 2019
- 4. Papadopoulos et al. Methodological Principles for Reproducible Performance Evaluation in Cloud Computing. TSE 2019 and (journal-first) ICSE 2020
- van Beek et al. Portfolio Scheduling for Managing Operational and Disaster-Recovery Risks in Virtualized Datacenters Hosting Business-Critical Workloads. ISPDC 2019
- van Beek et al. A CPU Contention Predictor for Business-Critical Workloads in Cloud Datacenters. HotCloudPerf19

 Iyushkin et al. Performance-Feedback Autoscaling with Budget Constraints for Cloud-based Workloads of Workflows. Under submission

Etc.





A LARGER VISION OF HOW COMPUTING WILL HELP OUR SOCIETY



ONE PROJECT TO MENTION...

GraphMassiziver

Big Graph Processing: Used in Al/ML, FinTech, Sci/Pharma, Industry 4.0, Energy Mgmt.*, etc.

Prodan et al. (2022) The GraphMassivizer project: Towards Extreme and Sustainable Graph Processing for Urgent Societal Challenges in Europe. IEEE Cloud Summit.



Data Center and HPC Resources

© 2022 Alexandru Iosup. All rights reserved.

Computational Continuum $DC \rightarrow Endpoint-edge-cloud$

Vision: Massivizing computer systems approaches enable holistic understanding and management in the computational continuum



Caring of Data at the Edge, HotEdge

APPROACH: WITH BROAD PARTICIPATION, UNIFY CURRENT COMPUTING MODELS, FORM A COMPUTATIONAL CONTINUUM



for the Edge Continuum, CoRR abs/2207.04159