# Developer life in academia & research

Tips and examples from PhD in databases

Iraklis Psaroudakis, 13/02/25, DevStaff.gr talk

# About

Iraklis Psaroudakis is working remotely from Greece as a Principal Software Engineer at Elastic, focusing on the distributed scalability of Elasticsearch.

Industry 3y

Previously, he was a Principal Member of Technical Staff at Oracle Labs in Zürich, Switzerland, concentrating on parallel and distributed programming, and graph-based analytical & machine learning workloads, especially in financial crime & compliance applications.

R&D 6y

Prior to Oracle, he completed his Ph.D. at the DIAS lab of EPFL in Lausanne, focusing on scaling up highly concurrent analytical database workloads on multi-socket multi-core servers through (a) sharing data and work across concurrent queries, and (b) adaptive NUMA-aware data placement and task scheduling. During his Ph.D., he cooperated with the SAP HANA database team.

Academia (& industry collab) 6y

Before his Ph.D., he completed his studies in Electrical & Computer Engineering at NTUA in Athens.

School 5y

# Grad life
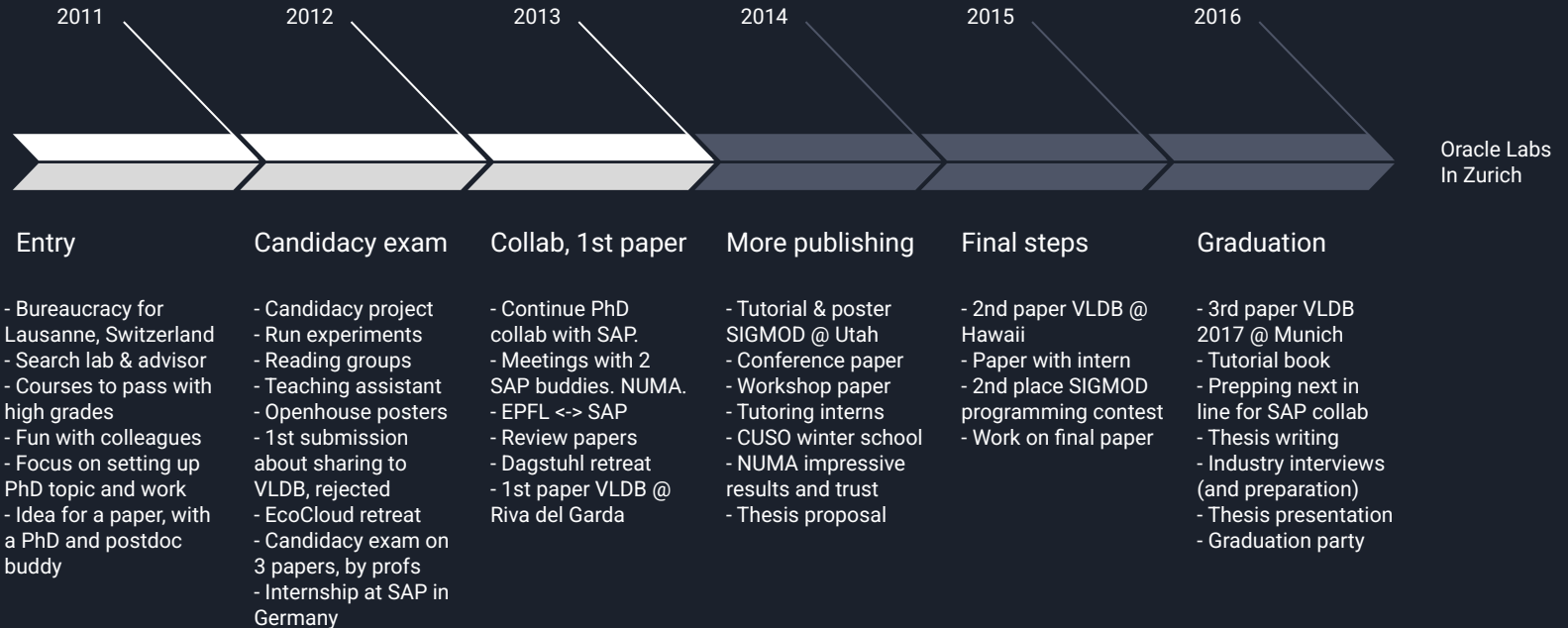
# Myths & truths

## Academia

- Goal: innovation, research
- Freedom
- Long-term research / impact
- Flexibility in organizing self
- Difficult work-life balance
- Progress through failures
- Struggle for money (grants, collabs)
- Uncertain job security
- Peer-review scrutiny
- Lonely
- Output: publish (or perish), presentations

## Industry

- Goal: monetization
- Company needs
- Short-term/changing plan, direct impact
- More fine-grained, stricter plan
- 9-5 schedule
- Should rarely fail
- Direct money flow (e.g., sales)
- Stable pay flow
- Scrutiny by direct colleagues/managers
- Team effort
- Output: product, reputation, patents

(albeit, R&D is some middle ground)

Ref: "Academia vs industry", Setia Pramana, 2013, https://www.slideshare.net/slideshow/academia-vs-industry/16334490

# PhD timeline at EPFL

**2011**  **2012**  **2013**  **2014**  **2015**  **2016**

Oracle Labs
In Zurich

## Entry

- Bureaucracy for Lausanne, Switzerland
- Search lab & advisor
- Courses to pass with high grades
- Fun with colleagues
- Focus on setting up PhD topic and work
- Idea for a paper, with a PhD and postdoc buddy

## Candidacy exam

- Candidacy project
- Run experiments
- Reading groups
- Teaching assistant
- Openhouse posters
- 1st submission about sharing to VLDB, rejected
- EcoCloud retreat
- Candidacy exam on 3 papers, by profs
- Internship at SAP in Germany

## Collab, 1st paper

- Continue PhD collab with SAP.
- Meetings with 2 SAP buddies. NUMA.
- EPFL <-> SAP
- Review papers
- Dagstuhl retreat
- 1st paper VLDB @ Riva del Garda

## More publishing

- Tutorial & poster SIGMOD @ Utah
- Conference paper
- Workshop paper
- Tutoring interns
- CUSO winter school
- NUMA impressive results and trust
- Thesis proposal

## Final steps

- 2nd paper VLDB @ Hawaii
- Paper with intern
- 2nd place SIGMOD programming contest
- Work on final paper

## Graduation

- 3rd paper VLDB 2017 @ Munich
- Tutorial book
- Prepping next in line for SAP collab
- Thesis writing
- Industry interviews (and preparation)
- Thesis presentation
- Graduation party

# Software development I did in research

## Academia

- Open-source, C++
- Low-level (CPU instructions, caches, performance)
- Prototyping – things can be buggy and not full features (e.g., half-baked parser) – happy path
- Supervision by PhDs, postdoc, professor. Collab with other academic colleagues.
- Ultimate output: experimental results showcasing the potential

## Industry collaboration

- Still prototyping, but baking inside the production code
- PRs reviewed by colleagues
- Real-world measurable impact on product
- Big hardware and clusters to use
- Advice from other industry experts
- Final result may shape the product

# Paper anatomy: Intro, related work

Title

Authors

Abstract

Intro

Copyright

>= 3 contributions

Key figure

LaTeX template ~ feels like coding

Related work

# Visual cues

Figures help readability

UML diagrams

Tables

Algorithms

Bold, math, italics, acronyms, paragraphs

Workflows

Conceptual figures

# Experimental evaluation

Repetitions, STD

Setup

Takeaways

No red+green. Print in gray

Configuration

Detailed results

Other types of papers exist (e.g., vision)

# Conclusion and references

**References!**

**Final summary and takeaway**

**Strict page limit. We squeeze stuff!**

**Priorities and fairness.** Priorities and fairness are an orthogonal issue out of this paper's scope. We note that our task scheduler supports priorities and a degree of fairness (based on query submission time) [28]. In typical cases, if a workload has user-defined priorities, the prioritized tasks will highly contribute to the utilization of RUH which will be considered first by DP for moving or partitioning.

**Task classes.** Tasks in the same class should have similar memory throughput. We assume that classes are defined manually, as we do in Section 5 for our NUMA-aware operators. One can further specialize classes, e.g., by the involved predicates. As seen in Section 7, an aggregation's memory throughput can vary depending on the complexity of the predicate. Ideally, we need a way to classify complex predicates. This is left out of the paper's scope. For our implementation and experiments, we use rather typical predicates and the defined task classes can sufficiently capture the memory intensity of our NUMA-aware operators' different phases.

**Unit of data placement.** Our unit of data placement is a row-wise partition of a table. An alternative would be column-wise partitioning, i.e., placing a table's columns on different sockets. In such a case, a global dictionary (see Section 4) is not needed, but queries referencing columns on multiple sockets incur a lot of remote accesses. For this work, we assume that the organization of associated columns into tables is left to the administrator.

**Balancing memory throughput.** We balance primarily the CPU utilization under the constraint of not creating memory bandwidth bottlenecks. This is to allow newly placed data to potentially increase their utilization. Since we balance local-only CPU utilization, this can indirectly balance memory TP as well as shown in many of our experiments. This is not guaranteed, however. One may wish DP to continue balancing memory TP after CPU utilization is balanced, under the constraint of not increasing the CPU imbalance. DP's possible actions can be extended to consider exchanging TBP or TGP between sockets. We have found only a few cases where balancing memory TP is required to slightly improve IPC and TP.

## 10. CONCLUSIONS

In this paper, we show that main-memory column-stores should not employ a static strategy of always partitioning data across all sockets, and always allowing inter-socket task stealing. We show that unnecessary partitioning involves an overhead of up to 2x in comparison to not partitioning. For this reason, we develop an adaptive data placement algorithm that can track a utilization imbalance across sockets, and can move or repartition tables at run-time to fix the imbalance. Also, we show that inter-socket stealing of memory-intensive tasks can hurt throughput by up to 4x in comparison to not stealing. For this reason, we develop an adaptive technique that disallows stealing at run-time for tasks whose memory intensity exceeds a fixed threshold for a NUMA server.

## 11. REFERENCES

[1] SAP HANA Platform SPS 11 Administration Guide, Dec. 2015. http://help.sap.com/hana_platform.
[2] SAP HANA Data Distribution Optimizer Administration Guide, Mar. 2016. http://help.sap.com/hana_options_dwf.
[3] TPC Benchmark H Rev. 2.17.1, 2016. http://www.tpc.org/.
[4] S. Agrawal et al. Database Tuning Advisor for Microsoft SQL Server 2005. In VLDB, pp. 1110–1121, 2004.
[5] M. Albutiu et al. Massively Parallel Sort-Merge Joins in Main Memory Multi-Core Database Systems. PVLDB, 5(10):1064–1075, 2012.
[6] C. Balkesen et al. Multi-Core, Main-Memory Joins: Sort vs. Hash Revisited. PVLDB, 7(1):85–96, 2013.
[7] S. Blagodurov et al. A Case for NUMA-aware Contention Management on Multicore Systems. In USENIX, 2011.
[8] M. Dashti et al. Traffic management: A holistic approach to memory placement on NUMA systems. In ASPLOS, pp. 381–394, 2013.
[9] R. Dementiev et al. Intel Performance Counter Monitor, Mar. 2016. https://software.intel.com/articles/intel-performance-counter-monitor.
[10] K. P. Eswaran. Placement of records in a file and file allocation in a computer network. In Information Processing, pp. 304–307, 1974.
[11] F. Färber et al. The SAP HANA Database – An Architecture Overview. IEEE Data Eng. Bull., 35(1):28–33, 2012.
[12] J. Giceva et al. Deployment of Query Plans on Multicores. PVLDB, 8(3):233–244, 2014.
[13] L. Golab et al. Distributed data placement to minimize communication costs via graph partitioning. In SSDBM, pp. 20:1–20:12, 2014.
[14] T. Gubner. Achieving many-core scalability in Vectorwise. 2014. Master's thesis, TU Ilmenau.
[15] G. Hill and A. Ross. Reducing outer joins. The VLDB Journal, 18(3):599–610, Aug. 2008.
[16] T. Kissinger et al. ERIS: A NUMA-Aware In-Memory Storage Engine for Analytical Workloads. ADMS, pp. 74–85, 2014.
[17] C. Lameter et al. NUMA (Non-Uniform Memory Access): An Overview. ACM Queue, 11(7):40:40–40:51, 2013.
[18] H. Lang et al. Massively Parallel NUMA-aware Hash Joins. In IMDM, pp. 1–12, 2013.
[19] P.-A. Larson et al. Enhancements to SQL server column stores. In SIGMOD, pp. 1159–1168, 2013.
[20] V. Leis et al. Morsel-Driven Parallelism: A NUMA-Aware Query Evaluation Framework for the Many-Core Age. In SIGMOD, pp. 743–754, 2014.
[21] C. Lemke et al. Speeding up queries in column stores: a case for compression. In DaWaK, pp. 117–129, 2010.
[22] Y. Li et al. NUMA-aware algorithms: the case of data shuffling. In CIDR, 2013.
[23] N. Mukherjee et al. Distributed Architecture of Oracle Database In-memory. In PVLDB, volume 8, pp. 1630–1641, 2015.
[24] I. Müller et al. Cache-Efficient Aggregation: Hashing Is Sorting. In SIGMOD, pp. 1123–1136, 2015.
[25] M. T. Özsu et al. Principles of Distributed Database Systems, Third Edition. Springer, 2011.
[26] D. Porobic et al. ATraPos: Adaptive transaction processing on hardware Islands. In ICDE, pp. 688–699, 2014.
[27] I. Psaroudakis et al. Task Scheduling for Highly Concurrent Analytical and Transactional Main-Memory Workloads. In ADMS, pp. 36–45, 2013.
[28] I. Psaroudakis et al. Scaling Up Concurrent Main-memory Column-store Scans: Towards Adaptive NUMA-aware Data and Task Placement. PVLDB, 8(12):1442–1453, 2015.
[29] I. Psaroudakis et al. Scaling Up Mixed Workloads: A Battle of Data Freshness, Flexibility, and Scheduling. In TPCTC, pp. 97–112, 2015.
[30] V. Raman et al. DB2 with BLU Acceleration: So much more than just a column store. PVLDB, 6(11):1080–1091, 2013.
[31] J. Rao et al. Automating Physical Database Design in a Parallel Database. In SIGMOD, pp. 558–569, 2002.
[32] O. Steinau et al. Method for calculating distributed joins in main memory with minimal communicaton overhead. US Patent App. 11/018,697.
[33] F. Transier et al. Aggregation in parallel computation environments with shared memory, 2012. US Patent App. 12/978,194.
[34] V. Viswanathan et al. Intel Memory Latency Checker v3.0, Mar. 2016. https://software.intel.com/articles/intelr-memory-latency-checker.
[35] M. Wagle et al. NUMA-Aware Memory Management with In-Memory Databases. In TPCTC, pp. 45–60, 2015.
[36] T. Willhalm et al. Vectorizing database column scans with complex predicates. In ADMS, pp. 1–12, 2013.
[37] Y. Ye et al. Scalable aggregation on multicore processors. DaMoN, 2011.
[38] E. Zamanian et al. Locality-aware Partitioning in Parallel Database Systems. In SIGMOD, pp. 17–30, 2015.
[39] M. Zukowski et al. Vectorwise: Beyond Column Stores. IEEE Data Eng. Bull., 35(1):21–27, 2012.

# Submission and reviews

Typical timeline:

- Prepare paper for months ahead of deadline
- Submit by deadline (e.g., Microsoft CMT)
- Submit conflicts of interest of authors
- Get (blind) reviews a couple of months later
- Possible rebuttal phase
- Acceptance/rejection notification
- If accepted, camera-ready deadline
- Copyright form
- At least 1 presented live in the conference
- Prepare presentation & do dry-runs
- Travel typically paid by academic institution
- Present & possible poster session

## VLDB 2025: Important Dates

All deadlines below are 5 p.m. PT, unless otherwise noted.

| | Chairs | Submissions Open | Submissions Deadline | Notifications | Camera-Ready Submissions |
|---|---|---|---|---|---|
| **Research Track** | Themis Palpanas<br>University Paris Cite<br><br>Nesime Tatbul<br>Intel Labs and MIT | 20th of every month starting March 2024 | 1st of every month until March 2025 | 15th of the next month following the deadline | Proceedings chairs will contact the authors with CR instructions by the 5th of the next month following the acceptance notification |
| **Industrial Track** | Surajit Chaudhuri<br>Microsoft Research, USA<br><br>Nikos Ntarmos<br>Huawei, United Kingdom<br><br>Jingren Zhou<br>Alibaba Group, China | TBD | March 17, 2025 | May 19, 2025 | July 14, 2025 |
| **Tutorials** | Hakan Ferhatosmanoglu,<br>University of Warwick and AWS<br><br>Madelon Hulsebos,<br>CWI | TBD | April 15, 2025 | May 30, 2025 | Camera-ready: July 1, 2025<br>Slides availability: August 20, 2025 |
| **Demonstrations** | Sourav S Bhowmick,<br>NTU, Singapore<br><br>Philippe Bonnet,<br>University of Copenhagen, Denmark | TBD | March 30, 2025 | May 27, 2025 | TBD |
| **Panels** | Jana Giceva,<br>Imperial College London<br><br>Alexandra Meliou,<br>University of Massachusetts, Amherst | TBD | May 15, 2025 | May 30, 2025 | June 31, 2025 |
| **Workshop Proposals** | John Paparrizos,<br>Ohio State University<br><br>Norman Paton,<br>University of Manchester | December 1, 2024 | January 17, 2025 | January 31, 2025 | N/A |
| **PhD Workshop Track** | Sonia Bergamaschi,<br>University of Modena and Reggio Emilia<br><br>Raul Castro Fernandez,<br>The University of Chicago | TBD | May 22, 2025 | June 30, 2025 | July 22, 2025 |

# Presentation tips

# Scholar profile

**Important for research jobs**

**Most cited on top**

**Citation indices**

**Consistent citations**

**Popular co-authors**



Iraklis Psaroudakis

Elastic
Verified email at elastic.co - Homepage

databases    graph processing    distributed processing

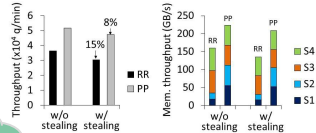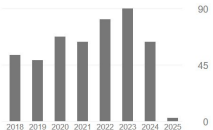| TITLE | CITED BY | YEAR |
|---|---|---|
| Parallel execution of parsed query based on a concurrency level corresponding to an average number of available worker threads<br>A Ailamaki, T Scheuer, I Psaroudakis, N May<br>US Patent 9,329,899 | 125 | 2016 |
| Scaling Up Concurrent Main-Memory Column-Store Scans: Towards Adaptive NUMA-aware Data and Task Placement<br>I Psaroudakis, T Scheuer, N May, A Sellami, A Ailamaki<br>Proceedings of the VLDB Endowment 8 (12), 1442-1453 | 83 | 2015 |
| Adaptive NUMA-aware data placement and task scheduling for analytical workloads in main-memory column-stores<br>I Psaroudakis, T Scheuer, N May, A Sellami, A Ailamaki<br>Proceedings of the VLDB Endowment 10 (2), 37-48 | 63 | 2016 |
| Scaling Up Mixed Workloads: A Battle of Data Freshness, Flexibility, and Scheduling<br>I Psaroudakis, F Wolf, N May, T Neumann, A Böhm, A Ailamaki, KU Sattler<br>Proceedings of the 6th TPC Technology Conference (TPCTC) | 61 | 2014 |
| Sharing data and work across concurrent analytical queries<br>I Psaroudakis, M Athanassoulis, A Ailamaki<br>Proceedings of the VLDB Endowment 6 (9), 637-648 | 51 | 2013 |
| Task Scheduling for Highly Concurrent Analytical and Transactional Main-Memory Workloads<br>I Psaroudakis, T Scheuer, N May, A Ailamaki<br>4th International Workshop on Accelerating Data Management Systems Using … | 45 | 2013 |
| Dynamic Fine-Grained Scheduling for Energy-Efficient Main-Memory Queries<br>I Psaroudakis, T Kissinger, D Porobic, T Ilsche, E Liarou, P Tözün, …<br>Proceedings of the 10th International Workshop on Data Management on New … | 28 | 2014 |
| {aDFS}: An Almost {Depth-First-Search} Distributed {Graph-Querying} System<br>V Trigonakis, JP Lozi, T Faltin, NP Roth, I Psaroudakis, A Delamare, …<br>2021 USENIX Annual Technical Conference (USENIX ATC 21), 209-224 | 22 | 2021 |
| Extending database task schedulers for multi-threaded application code<br>F Wolf, I Psaroudakis, N May, A Ailamaki, KU Sattler<br>Proceedings of the 27th International Conference on Scientific and … | 19 | 2015 |
| Parallel execution of parsed query based on a concurrency level<br>A Ailamaki, T Scheuer, I Psaroudakis, N May<br>US Patent 9,983,903 | 17 | 2018 |
| Csr++: A fast, scalable, update-friendly graph data structure<br>S Firmli, V Trigonakis, JP Lozi, I Psaroudakis, A Weld, D Chiadmi, S Hong, …<br>24th International Conference on Principles of Distributed Systems (OPODIS'20) | 15 | 2020 |

**Cited by** — VIEW ALL

| | All | Since 2020 |
|---|---|---|
| Citations | 604 | 375 |
| h-index | 12 | 10 |
| i10-index | 14 | 11 |



2018  2019  2020  2021  2022  2023  2024  2025

**Public access** — VIEW ALL

0 articles          1 article

not available          available

Based on funding mandates

**Co-authors** — EDIT

Anastasia Ailamaki
Professor of Computer and Com…

Norman May
SAP SE

Danica Porobic
Oracle

Erietta Liarou
Postdoctoral Researcher at EPFL

Pınar Tözün
IT University of Copenhagen

Jean-Pierre Lozi
Researcher at Inria

Sungpack Hong
Oracle Labs

Kai-Uwe Sattler
Professor of Computer Science, …

# Conclusion

- Software developer life in academia & research involves
  - 35% coding
  - 35% experimenting, writing
  - 15% reading related work
  - 15% talking, presenting, getting feedback
- More flexibility and challenging topics than industry
- Less life-work balance than industry
- Publish or perish
- Experimental evaluation is key
- Presenting in a succinct manner is key
- PhD necessary for academia, wished for in R&D, not necessary for industry

# Thank you!

Questions? www.kingherc.com