

# Bruno Magalhaes

## Machine Learning, High Performance Computing and Big Data

✉ [brunomaga@gmail.com](mailto:brunomaga@gmail.com) 📞 +41 77 487 89 92 🛠 C, C++, CUDA, HPX, MPI, Python, PyTorch  
🖱 <https://brunomaga.github.io> 🌐 [brunomaga](#) 🗺 Portuguese and Swiss  
🏠 Lausanne, Switzerland 🗣 fluent in English, French, Portuguese, Spanish 🏊 waterpolo, skiing



## 📁 Work Experience

present Sep 2019	<b>AI Resident » Researcher » Senior Researcher, Microsoft Research, Cambridge, UK</b> <ul style="list-style-type: none"><li>as Sr Researcher, 2022-: porting Transformer-based models to confidential optical hardware (ongoing). Large ML models scaling via model parallelism, sharding, pipelining, gradient accumulation, checkpointing, IO offloading, shared memory, mixed precision, model compression, and distillation. Likelihood estimators, Multi-armed bandit and Gaussian Processes for error quantification and fine-tuning of optical systems. Information encoding (Gray, Huffman), error correction (LDPC) and channel capacity (Blahut-Arimoto) for non-binary systems. Mentoring of junior members and PhD interns.</li><li>as Researcher, 2021: computer vision models for thousand-object classification on 3D glass at Project Silica. Distributed data parallelism. Presenter of talks on the topics of <i>CPU/GPU optimization</i>, <i>distributed algorithms</i> and <i>AI SuperComputing</i>.</li><li>as AI Resident, 2019-20: RNNs, GRUs, Encoder-Decoders, and Bayesian Optimization for regression on time series, to improve load balancing of Exchange email servers on distributed exabyte-scale COSMOS databases. Graph Neural Nets for a recommendation system on a trillion-edge graph of meetings, documents, emails and users, stored on a distributed spark databases.</li><li>always: full-stack MLOps and CI/CD for cluster and cloud (AzureML) environments. Fine-tuning ML for hardware (network, memory) and business specifications (cost vs accuracy vs runtime trade-off). Performance modelling and analysis at scale.</li></ul>
Aug 2019 Mar 2015	<b>PhD candidate » postdoctoral researcher, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland</b> <ul style="list-style-type: none"><li>Research, development (C, C++) and publication on the field of distributed-parallel asynchronous variable-order variable-step simulation and optimization of detailed spiking neural networks, on Cray and SGI supercomputers with 10K+ compute nodes.</li><li>Technologies: C, C++, distributed asynchronous runtime systems (HPX-5) for communication (InfiniBand, RDMA) and computation (concurrency, scheduling); global memory addressing; dynamic load balancing; vectorization; cache optimization.</li><li>Teaching assistant (400h) for <i>Unsupervised and reinforcement learning</i>, <i>Project in neuroinformatics</i> and <i>In silico neuroscience</i>.</li><li>Scientific reviewer for <i>SuperComputing</i>, <i>IPDPS</i>, and <i>ISC</i> conferences. As postdoc: supervision of PhD students and engineers.</li></ul>
Feb 2015 Mar 2011	<b>Research Engineer for High Performance Computing, Blue Brain Project, EPFL, Lausanne, Switzerland</b> <ul style="list-style-type: none"><li>Research, development (C, C++, MPI, OpenMP) and publication of methods for parallel/distributed volumetric spatial decomposition, load balancing, spatial indexing, sorting, I/O, sparse matrix transpose, and graph navigation, that underlie an efficient storage and processing of neural networks on SGI and IBM BlueGene supercomputers with 16K+ compute nodes.</li></ul>
Feb 2011 Sep 2009	<b>Junior Architect for IT infrastructures, Noble Group, London, New York, &amp; São Paulo</b> <ul style="list-style-type: none"><li>Design and configuration of Linux servers, CISCO networks, and backup/redundancy sites for physical trading of commodities.</li></ul>
Oct 2008 Mar 2007	<b>Analyst programmer, Investment Property Databank (now MSCi Real estate), London, UK</b> <ul style="list-style-type: none"><li>Development of a search engine and web/windows app (C#, C++) for efficient storage and analytics of financial data.</li></ul>

## 📖 Education

Jun 2019 Mar 2015	<b>PhD Computational Neuroscience, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland</b> <ul style="list-style-type: none"><li>Thesis <i>Asynchronous Simulation of Neuronal Activity</i> nominated for the EPFL PhD excellency award and the IBM Research best thesis in computational sciences award. Visiting researcher at CREST at Indiana University (US), Summers 2015, '16 and '17.</li></ul>
Sep 2009 Oct 2008	<b>MSc Advanced Computing, Imperial College London, UK</b> <ul style="list-style-type: none"><li>Grade: Merit. Thesis <i>GPU-enabled steady-state solution of large Markov models</i> awarded distinction and published at <i>NSMC'10</i>.</li></ul>
Jul 2007 Oct 2002	<b>BSc Systems Engineering and Computer Science (5 year degree), University of Minho, Portugal</b> <ul style="list-style-type: none"><li>Grade: A, top 10%. Exchange student at the University of Maribor, Slovenia, 2005/06. Intern at IBM and CERN.</li></ul>

## 📄 Selected Publications [full list on scholar.google.com/citations?user=pirWLLgAAAAJ](https://scholar.google.com/citations?user=pirWLLgAAAAJ)

- |      |  |
|------|--|
| 2023 | Multi-dimensional optical data writing techniques for cloud-scale archival storage, Proc. Laser Applications in Microelectronic and Optoelectronic Manufacturing (LAMOM) XXVIII                        |
| 2022 | Cloud-Scale Archival Storage Using Ultrafast Laser Nanostructuring, Conf. Lasers and Electro-Optics Technical Digest Series 2022   |
| 2020 | Fully-Asynchronous Fully-Implicit Variable-Order Variable-Timestep Simulation of Neural Networks, Proc. International Conference on Computational Science (ICCS 2020), Amsterdam, Holland              |
| 2019 | Asynchronous SIMD-Enabled Branch-Parallelism of Morphologically-Detailed Neuron Models, Frontiers in Neuroinformatics  |
| 2019 | Exploiting Implicit Flow Graph of System of ODEs to Accelerate the Simulation of Neural Networks, Proc. International Parallel & Distributed Processing Symposium (IPDPS 2019), Rio de Janeiro, Brazil |
| 2019 | Fully-Asynchronous Cache-Efficient Simulation of Detailed Neural Networks, Proc. International Conference on Computational Science (ICCS 2019), Faro, Portugal   |
| 2015 | Reconstruction and Simulation of Neocortical Microcircuitry, Cell 163, 456–492.  |