

# Bruno Magalhaes    High Performance Computing and Machine Learning

✉ brunomaga@gmail.com    🌐 brunomaga    💬 brunomaga    🔗 <https://brunomaga.github.io>    🏠 Lausanne, Switzerland  
🔧 Python, PyTorch, Triton, DeepSpeed. C, C++, CUDA, HPX, MPI    🇵🇹 Portuguese-Swiss    🌍 fluent in English, French, Portuguese, Spanish

## Work Experience

present	<b>Principal Researcher for AI Systems and Networks, Huawei Research, Zurich, Switzerland</b>
July 2025	<ul style="list-style-type: none"><li>➢ Started, grew and led the AI Data Center Networks team (7). Defined the research direction and implementation of distributed kernel fusion (scale-up) and in-network computing (scale-out), towards next-gen AI network devices on HPC networks.</li><li>➢ Developed (C++, CUDA, nvshmem) fused asynchronous communication-computation mega-kernels for Ring/Tree Attention, Megatron-LM Tensor Parallelism, LLMs, MoEs dispatch+combine, and expert sharding. Many DeepSeek-like improvements.</li></ul>
May 2025	<b>Senior Research Engineer for Large-scale AI, Synthesia, Zürich, Switzerland</b>
Nov 2023	<ul style="list-style-type: none"><li>➢ Started, grew and led the ML performance team (4), aimed at large-scale training and efficient inference of diffusion and transformer models for text, audio and video (Wan, CogVideoX, DALL-E, DDPM, AnymateAnyone). Implemented data, context and sequence parallelism (Ulysses, Ring attention), sharding and activation checkpointing. Performance modelling (scaling laws).</li><li>➢ Model speedup via kernels fusion (triton, CUDA, torch compile), quantisation (int8 Sage attention), and channels-last format.</li><li>➢ Created the DeepSpeed modules for distributed curriculum learning setup (PR 5129) and variable batch size and LR (PR 7104).</li></ul>
Oct 2023	<b>AI Resident » Researcher » Senior Researcher, Microsoft Research, Cambridge, UK</b>
Sep 2019	<ul style="list-style-type: none"><li>➢ as Sr Researcher, 2022-23: Efficient inference of vision models on optical hardware with pipeline and model parallelism (C++). Large-scale ML training via data parallelism, sharding, gradient accumulation, activation checkpointing, IO offloading and distillation (DeepSpeed/ZeRO and Torch Distributed/RPC). Distributed Mixture of Experts. Likelihood estimators, information encoding and error correction (LDPC) for storage systems. Gaussian Processes for hyperparameter finetuning. Mentoring of junior members and PhD interns. MLOps on cluster and cloud environments (docker, builds, CI).</li><li>➢ as Researcher, 2021: computer vision models for thousand-object classification on 3D glass at Project Silica. Presenter of talks on the topics of <i>CPU/GPU optimization, distributed algorithms and 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 database.</li></ul>
Aug 2019	<b>HPC engineer » PhD candidate » postdoc researcher, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland</b>
Mar 2011	<ul style="list-style-type: none"><li>➢ as postdoc: supervision of PhD students and engineers. Scientific reviewer for <i>SuperComputing, IPDPS, and ISC</i> conferences.</li><li>➢ as PhD candidate, 2015-19: research, development (C, C++, HPX) and publication on the field of distributed-parallel asynchronous variable-order variable-step simulation of spiking neural networks, on Cray and SGI supercomputers (10K+ nodes).</li><li>➢ Technologies: asynchronous runtime system (HPX) with distributed memory with global addressing (InfiniBand, RDMA, PGAS), and distributed control objects (concurrency, scheduling); dynamic load balancing; vectorization; cache optimization.</li><li>➢ Teaching assistant (400h) for <i>Unsupervised and reinforcement learning, Project in neuroinformatics and In silico neuroscience</i>.</li><li>➢ as HPC Engineer, 2011-15: research, development (C, C++, MPI, OpenMP) and publication of methods for parallel/distributed volumetric decomposition, load balancing, spatial indexing, sorting, IO, sparse matrix transpose, and graph navigation, that underlie an efficient processing of neural networks on SGI and IBM Blue Gene Q supercomputers with 16K+ compute nodes.</li></ul>
Feb 2011	<b>Junior Architect for IT infrastructures, Noble Group, London, New York, &amp; São Paulo</b>
Sep 2009	<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	<b>Analyst programmer, Investment Property Databank (now MSCi Real estate), London, UK</b>
Mar 2007	<ul style="list-style-type: none"><li>➢ Development of a search engine app (C#) and backend (C++) for efficient storage and analytics of financial data.</li></ul>

## Education

Jun 2019	<b>PhD Computational Neuroscience, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland</b>
Mar 2015	<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 award for best thesis in computational sciences. Visiting researcher at Indiana University (US) during Summers of 2015, '16 and '17.</li></ul>
Sep 2009	<b>MSc Advanced Computing, Imperial College London, UK</b>
Oct 2008	<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	<b>BSc Systems Engineering and Computer Science (5 year degree), University of Minho, Portugal</b>
Oct 2002	<ul style="list-style-type: none"><li>➢ Grade: A, top 10%. Exchange student at the University of Maribor, Slovenia, 2005/06. Intern at IBM. Part-time project at CERN.</li></ul>

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

2024	Project Silica: sustainable cloud archival storage in glass, <i>Frontiers in Ultrafast Optics: Biomedical, Scientific, and Industrial Applications XXIV</i>
2023	Project Silica: Towards Sustainable Cloud Archival Storage in Glass, <i>SOSP '23: Proc. of the 29th Symposium on Operating Systems Principles</i>
2022	Cloud-Scale Archival Storage Using Ultrafast Laser Nanostructuring, <i>Conf. Lasers and Electro-Optics Technical Digest Series 2022</i>
2020	Fully-Asynchronous Fully-Implicit Variable-Order Variable-Timestep Simulation of Neural Networks, <i>Proc. ICCS 2020, Amsterdam, Holland</i>
2019	Exploiting Implicit Flow Graph of System of ODEs to Accelerate the Simulation of Neural Networks, <i>Proc. IPDPS 2019, Rio de Janeiro, Brazil</i>
2015	Reconstruction and Simulation of Neocortical Microcircuitry, <i>Cell</i> 163, 456–492.