



Towards a Green Al

Evolutionary solutions for an ecologically-viable artificial intelligence

Nayat Sánchez-Pi and Luis Martí Inria Chile Research Center, Santiago, Chile.

GECCO '21 Companion, July 10–14, 2021, Lille, France © 2021 Association for Computing Machinery. ACM ISBN 978-1-4503-8351-6/21/07 \$15.00 https://doi.org/10.1145/3449726.3461428

Instructors

3

Nayat Sánchez-Pi is currently the Director and CEO of the Inria Chile Research Center, created in 2012 by Inria, the French National Research Institute for Digital Sciences to facilitate scientific and industrial cooperation between France, Chile, and Latin America. Before that, she was a professor of Artificial Intelligence and Human-Computer Interaction at the Department of Informatics and Computer Science of the Institute of Mathematics and Statistics of the Rio de Janeiro State University. Prof. Sánchez-Pi's research interests have broadened over the years and span topics that range from artificial intelligence, machine learning, and data mining to ambient intelligence, ubiquitous computing, and multi-agent systems. She received a degree in Computer Science in 2000 from the University of Havana and a Ph.D. degree in Computer Science in 2011 from the Universidad Carlos III de Madrid. She has led numerous research projects applying evolutionary computation, machine learning, and other artificial intelligence methods.

Luis Martí is currently the scientific director of Inria Chile, the Chilean Center of Inria, the French National Institute for Computational Sciences. Before that, he was a senior researcher of the TAU team at Inria Saclay since 2015. He was also an Adjunct Professor (tenured) at the Institute of Computing of the Universidade Federal Fluminense. Previous to that, Luis was a CNPq Young Talent of Science Fellow at the Applied Robotics and Intelligence Lab of the Department of Electrical Engineering of the Pontificia Universidade Católica do Rio de Janeiro, Brazil. Luis did his Ph.D. at the Group of Applied Artificial Intelligence of the Department of the Universidad Carlos III de Madrid, Madrid, Spain, and got his Computer Science degree from the University of Havana. He is mainly interested in artificial intelligence, and, in particular, machine learning, neural networks, evolutionary computation, optimization, machine learning, hybrid systems, and all that.





Ínría_



Computing (AI) for efficient green stuff



OPTIMIZATION *Plant and export layouts and robustness.*



MODELING Wind and wave modeling for solar, wind and tidal energies.



PREDICTION Being able to predict the demand and production leads to efficient production.



<section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><list-item><list-item><list-item><section-header><list-item><list-item><list-item><list-item><list-item>







Common carbon foo	tprint benchmarks
in lbs of CO2 equivalent	
Roundtrip flight b/w NY and SF (1 passenger)	1,984
Human life (avg. 1 year)	11,023
American life (avg. 1 year)	36,156
US car including fuel (avg. 1 lifetime)	126,000
Transformer (213M parameters) w/ neural architecture search	626,155











Better	It is unlikely that we get right of GPUs (or TPUs) at training time.
Hardware	There are hardware alternatives at use time :
	→ Field Programmable Gate Arrays (FPGAs), Application-Specific Instruction-set Processors (ASIPs), etc.
	We should also keep exploring the use of <i>low-precision</i> computing:
	Reducing the quality (and therefore length) of the floating-point representation of numbers.





Roofline model for evaluating AI applications

Observed characteristics DT:

- ➡ Balanced in memory and processing requirement
- ↓ The bottleneck is DRAM memory
- They have low arithmetic intensity and GFLOPs



Vitor de Sá, V, Klöh, V, Schulze, B, and Ferro, M. (2020). Análise de desempenho e de requisitos computacionais utilizando o modelo Roofline: Um estudo para aplicações de inteligência artificial e do NAS-HPC. In WSCAD 2020 – WIC



18

Cloud computing for AI?



19

AI and ML pipelines are very computationally intensive but not continuously run:

Training: when we fit a model to the data available, very costly but run infrequently -or maybe once.

- Punctual use of expensive computing equipment.
- GPUs are both expensive to acquire and expensive to operate (high energy demands, cooling, etc.).

Prediction: using the model to make decisions, low computing requirements and used very frequently.

→ High-availability low cost computing.

Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L.-M., Rothchild, D., So, D., Texier, M., & Dean, J. (2021). Carbon Emissions and Large Neural Network Training. In arXiv [cs.L6]. arXiv. http://arxiv.org/abs/2104.10350



Not all cloud locations are the same

For example, different cloud locations and their CO₂ impact (indirectly, cost)



Self-scaling and cloud computing



Self-scaling computing facilities make available a pool of shared resources.

Optimally schedule computing time.

Cloud computing allows **to pick the location** where programs will be run.

Code is mobile!

We can, for example, "track the sun" and ensure that the AI/ML processes use renewable sources.

"the choice of DNN, datacenter, and processor can reduce the carbon footprint up to ~100x - 1000x."¹

1. Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L.-M., Rothchild, D., So, D., Texier, M., & Dean, J. (2021). Carbon Emissions and Large Neural Network Training. In arXiv [cs.LG]. arXiv. http://arxiv.ors/abs/210410350



Ínría_

Smarter experimentation: Better AutoML



Finding the **right configuration** of the **hyperparameters** probably where more energy is consumed.

- This is an optimization problem => NP-Hard problem
- ➔ Neural architecture search, AutoML, AutoDL, etc.
- → Go beyond 'regular' hyperparameter search
- However, better approaches like evolutionary computing is here to help!
 ...but they need populations of individuals, hence
 - more energy.
 - This is a multi-objective optimization problem!

21





Self adaptation



25

To look for **methods** that **adapt** their complexity **automatically** to the **complexity** of the **problem** being solved.

Neural networks based on **adaptive resonance theory (ART)** and growing neural gas (GNG) have rules to adapt themselves to the complexity of the problem.

This self-adaptation is best profited when using cloud-based infrastructure.















Thank you! Obrigados! Merci ! - ¡Gracias!

Questions?

Find more in greenai.inria.cl

Ínría_