







Spoke 7 - Edge and Exascale AI

Giuseppe Averta, Barbara Caputo (PI) Politecnico di Torino

Rome, Oct 20, 2023



Politecnico di Torino









Mainstream AI is inefficient



Energy

"Training GPT-3 would cost over \$4.6M. [...] With the increase in parameters, there's an exponential increase in energy."^{Lambda Labs}

¹Source: <u>https://lambdalabs.com/blog/demystifying-gpt-3</u> ²Source: <u>https://www.nature.com/articles/s41928-018-0068-2</u>



Hardware

"The computational demands of AI presents an emerging problem for its implementation on different hardware platform"² Nature



Data

"Data labeling takes anywhere from 35 to 80% of project budgets."³ Forbes



Environment

"Training a single AI model can emit as much carbon as five cars in their

> lifetimes"⁴ MIT Technology Review

³Source: <u>https://www.forbes.com/sites/cognitiveworld/2022/08/06/no-youre-not-alone-google-is-also-making-this-big-mistake-on-ai</u> ⁴Source: <u>https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/</u>



















Spotlight works









We study how to better build neural models

The Multiply and Max/Min (MAM) -based neuron

MAM-based neurons do not accumulate all the weighted inputs, but **sum together only the maximum and minimum contributes.**



Neural networks built with MAM learn to use **only a small subset of interconnections during inference**.

Prono, Luciano; Bich, Philippe; Mangia, Mauro; Pareschi, Fabio; Rovatti, Riccardo; Setti, Gianluca (2023). A Multiply-And-Max/min Neuron Paradigm for Aggressively Prunable Deep Neural Networks. TechRxiv. Preprint. https://doi.org/10.36227/techrxiv.22561567.v1

Case study: MAM for ECG autoencoder tail



By using all the available memory on device, MAM achieves 33 dB reconstruction performance (where a standard DNN would achieve 18 dB)









We assess the reliability of existing architectures

Robustness w.r.t. external disturbance (e.g. neutron strikes)







Robustness w.r.t. information representation (e.g. POSIT)

- V.Turco, A. Ruospo, G. Gavarini, E. Sanchez, M. Sonza Reorda, "Uncovering hidden vulnerabilities in CNNs through evolutionary-based Image Test Libraries" IEEE Int. Symp on DFT 2023

- Cavagnero, Niccolò, et al. "Transient-Fault-Aware Design and Training to Enhance DNNs Reliability with Zero-Overhead." IEEE IOLTS 2022.







Future Artificial Intelligence Research

amazon alexa

And study how and what models are learning



Divergence as a measure of anomalous behavior of a data subgroup S w.r.t. overall dataset D for a function f

Subgroup	Sup	acc	Δ_{acc}	t
{age=22-40, gender=male, loc=none, speakRate=high, tot_silence=high}	0.03	74.79	-18.38	4.7
{action=increase, gender=male, speakRate=high}	0.03	74.81	-18.36	4.9

Miss-classification (in red) rate of speech models (intent classif.) also depend on subgroups – how to detect low-accuracy subgroups without supervision?



$$gain_f(S, M_1, M_2) = f(S, M_2) - f(S, M_1)$$



Koudounas, Alkis, et al. "Exploring subgroup performance in end-to-end speech models." ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2023.









We care about applications: Efficient 2D and 3D computer vision on edge and very large scale











We care about applications: Understanding human behavior from egocentric data

- Mobile models for first person action recognition
- Challenges:
 - Untrimmed videos
 - Robust to Multi-modal domain gap (time, space)
 - Model footprint for edge deployment



3rd place EPIC KITCHENS UDA challenge at <u>CVPR2022</u> & <u>CVPR2023</u>

Goletto, G., et al. 2023 RA-L, Peirone et al. 2023, under review









Task



Not only "supervised": RL for (soft) intelligent manipulators



(B) Training efficiency (up to 8x faster)

Tiboni, G., Protopapa, A., Tommasi, T., & Averta, G. (2023). Domain Randomization for Robust, Affordable and Effective Closed-loop Control of Soft Robots. In 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS).









AI and its societal impact



Semantic segmentation of cultural heritage point clouds

Predictive maintenance on district heating networks



Dependence	SMALL	Inclusiveness		Training likelihood	od
	Pange [01]		Probability		Probability
Contingency coefficient	01/17		Range [0,1]		Range [0,1]
Effect size	0.1413	P(Asian ∩ 0)	0.0023	P(Caucasian 1)	0.293
	0.1427	P(Asian ∩ 1)	0.0008	P(Caucasian 0)	0.381
		P(Black ∩ 0)	0.1514	P(0 Caucasian)	0.609
D'		P(Black ∩ 1)	0.1661	P(1 Caucasian)	0.391
Diverseness		P(Caucasian ∩ 0)	0.1281	P(Black 1)	0.591
	Probability	P(Caucasian ∩ 1)	0.0822	P(Black 0)	0.450
Target variable	Range [0,1]	P(Hispanic ∩0)	0.0320	P(0 Black)	0.477
0	1 05/5	P(Hispanic ∩ 1)	0.0189	P(1 Black)	0.523
1	1 0.455	P(Native american ∩ 0)	0.0006		
	0.100	P(Native american ∩ 0)	0.0005		
Protected attribute		P(other ∩ 0)	0.0219		
Asian 🗠	0.005	P(other ∩ 1)	0.0124		
Black	0.514				
Caucasian	0.341 C]	

Ethically-sensitive dataset labeling

Balancing the complexity and interpretability of Al-based energy management strategies











Spoke 7 in Numbers

- 30+ professors
- 8 Assistant Professors enrolled (more coming soon)
- 4 PhD students enrolled (more coming soon)
- 6 research assistants (all on board)
- 9+ journal papers with peer review
- 40+ conference papers with peer review









Dissemination activity











Spoke 7 in the world











Not only rocket science

- We do have extensive experience in academia-industry collaborations
- Many already existing collaborations
- Keen to use FAIR as flywheel to enhance current relationships and open new opportunities













Future Artificial Intelligence Research

Thanks

- Happy to chat in the networking session
- Come and visit us in Turin!
- giuseppe.averta@polito.it
- <u>barbara.caputo@polito.it</u>

This study was carried out within the FAIR - Future Artificial Intelligence Research and received funding from the European Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) – MISSIONE 4 COMPONENTE 2, INVESTIMENTO 1.3 – D.D. 1555 11/10/2022, PE00000013). This manuscript reflects only the authors' views and opinions, neither the European Union nor the European Commission can be considered responsible for them.