

# Sustainable AI: Assessing the environmental impact of AI

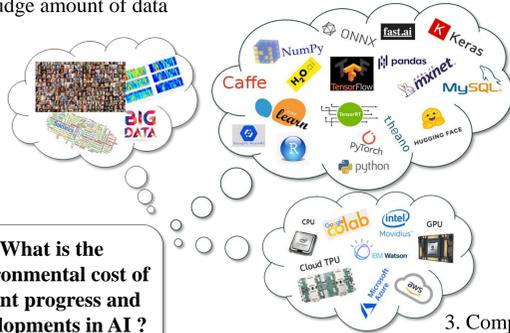
M. Guillaumont\*, T. Cuvilliers\*, P. Greullet\*, O. Matz\*, B. Deguilhem\*

\* Capgemini Engineering - DEMS, département R&I, Toulouse, France



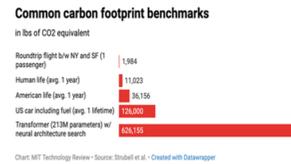
## General context overview, main motivation and use case

### 1. Huge amount of data



### 2. Many languages and framework

➤ The carbon footprint and environmental impact of AI is not negligible:



Training of Transformer algorithm (NAS) in 2019<sup>1</sup> has emitted more CO<sub>2</sub> than 5 cars in their entire life cycle (more than 300 000 kg eq CO<sub>2</sub>). Due to **increase** in the total number of optimized parameters (65M in 2017 vs 213M in 2019) and the total number of data used for training the CO<sub>2</sub> emission for this algorithm has been multiplied by **24 in two years**.

➤ Recent developments in IA tend to worsen this impact:

NLP algorithm	Parameters	Data (number of token)	Year
BERT <sup>2</sup>	213M	3 B	2019
GPT-2	1.5 B	40 B	2019
GPT-3 <sup>3</sup>	175B	500 B	2020

Wolff *et al.* has evaluated that the training of GPT-3 algorithm has emitted more than 80 000 kg eq. CO<sub>2</sub> based on average intensity carbon in USA in 2017<sup>3</sup>.

➤ Yet environmental impact is not considered as valuable metric to evaluate performance of deep learning algorithms.

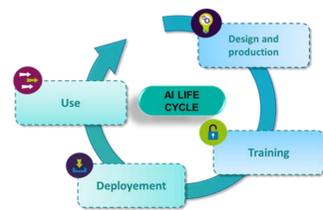
**Our first objective:** Recent studies report that training of NLP algorithms have a significant impact, what about this impact in the computer vision domain and for other processes than training ?  
**Project use case:** Object detection in **computer vision** using convolutional neural network.

Evaluation of classification, segmentation and detection tasks. Deployment of IA on optimized embedded hardware.

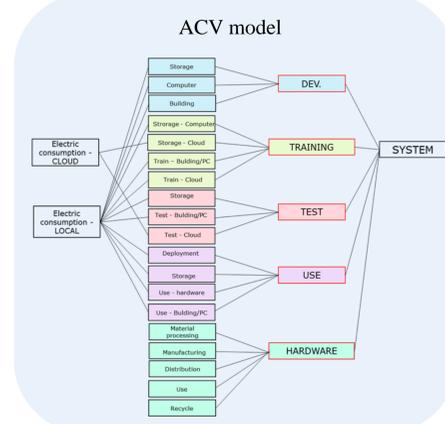


## Life Cycle Analysis

➤ Our goal is to adapt LCA methodology to the life of an AI taking in account Design and production, Training, Deployment and Use process.

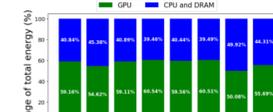


ACV is a normalized methodology for **systemic quantification** of the environmental impacts of a product, a process or a service.



## Asses energy consumption during training

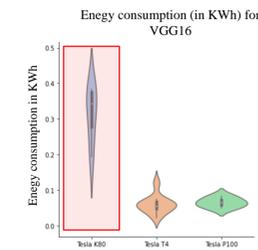
➤ On the training process, the nature of the **hardware**, algorithm **architecture** and the **localization** of the datacenter are depicted as key parameters to assess energy use and carbon footprint. CarbonTracker<sup>4</sup> and CodeCarbon libraries are used to monitor energy use during training.



Nature of GPU can lead to a **five times increase** in energy consumption.

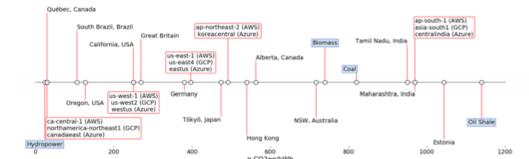
GPU usage is accounting for more than a half of the total energy<sup>4</sup>

	Efficacité (GOPS/W)	Précision de calcul
T4	804,49	Int 8
P100	63,85	Float 16
K80	23,89	Float 32



Ratio CO <sub>2</sub> Emissions/Energy consumption (kgCO <sub>2</sub> /kWh)
1. USA, Oregon (0,139)
2. USA, District of Columbia (0,219)
3. USA, South Carolina (0,286)
4. USA, Nevada (0,349)
5. USA, Iowa (0,453)
6. USA, Maryland (0,459)
7. Belgium, Brussel capital (0,627)
8. Netherlands, Groningen (0,770)
9. Taiwan, New Taipei (0,790)

Depending on the localization of the datacenter CO<sub>2</sub> emissions can be **reduced up to 90%**.



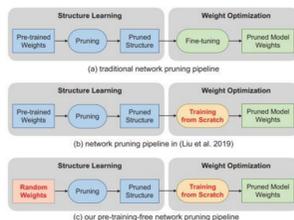
## How to reduce the environmental impact of AI ?

Several methods can be used to reduce environmental impact of IA:

- During Training and inference: choose the best parameters among nature of **hardware** and **software**, hardware's **localization**, algorithms **architectures**...
- During inference: Use of **model compression of model reduction methods**

### ➤ Pruning

This methods is efficient to develop smaller network by eliminating unnecessary values in the weight tensor

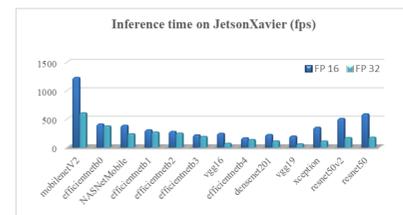


**Winning Tickets method:** Pruned subnetworks can reach test accuracy comparable to the original network in the same number of iterations<sup>5</sup>.

This figure is extract from Y. Wang *et al.*, « Pruning from Scratch », 2019 <http://arxiv.org/abs/1909.12579>

### ➤ Quantization

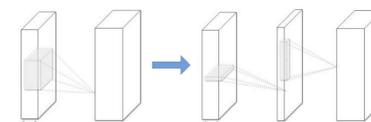
- This method is suitable for **convolutional** and **dense** layers. It consist in reducing the bit size of the model weights.
- Several levels of quantization: FP32, FP16, INT8.
- Hardware specific method: Nvidia Maxwell GPU supports FP16 but Nvidia Volta GPU supports INT8
- Integer quantization may require dataset calibration.



Going from FP32 to FP16 can increase the inference time (fps) up to 70%. The efficiency of quantization method highly depends on algorithm nature.

### ➤ Low-rank Factorization<sup>6</sup>

- This method is suitable for **convolutional** and **dense** layers. It consist in using tensor/matrix decomposition to estimate the informative parameters. (e.g. Truncated SVD)
- Cheng *et al.* report results achieved with no loss in accuracy<sup>6</sup>
- Compression: less parameters to store (2-5x less)
- Acceleration: architecture-dependant (1-2x faster)



On the left : original convolutional layer. The figure on the right show application of low-rank constraints to the convolutional layers with rank-K.

## Next steps and bibliography

- Generation of a **large dataset** in Computer Vision domain for Classification, Segmentation and Detection tasks : fist focus on **carbon footprint**
- Run ACV models to evaluate the **global environmental impact** of IA with an ensemble of criteria
- Complexification of **ACV models** by adding other processes than training and hardware life cycle
- Reduction of environmental impact during training and inference using **model compression of model reduction methods** on specific hardware

- E. Strubell *et al.*, « Energy and Policy Considerations for Deep Learning in NLP », 2019, <http://arxiv.org/abs/1906.02243>
- J. Devlin, *et al.*, « BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding », 2019, <https://arxiv.org/pdf/1810.04805.pdf>
- A. Brown *et al.*, « Language Models are Few-Shot Learners », 2020, <https://arxiv.org/abs/2005.14165>
- A. Wolff *et al.*, « CarbonTracker : Tracking and Predicting the Carbon Footprint of Training Deep Learning Models », 2020, <https://arxiv.org/abs/2007.03051>
- J. Frankle and M. Carbin, « The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks », 2018, <http://arxiv.org/abs/1803.03635>
- Y. Cheng, *et al.*, « A Survey of Model Compression and Acceleration for Deep Neural Networks », 2020, <http://arxiv.org/abs/1710.09282>