# WHOAMI

- Assistant Professor in Computer Science @ University of Salento, Italy
- https://goo.gl/maps/dnANaZjMLRugJ5k2A
- Visiting researcher @ Smart Grid Energy Research Center (UCLA - Los Angeles)
- Main research topics: enhanced CX models in the Fintech field using IoT, AI and DLT, and, obviously, green software
- Co-founder of University spin-off (Vidyasoft s.r.l.)
- Open to collabs (thesis, host students in lab, …)
- https://software.green/

# CLIMATE CHANGE

CHANGES IN

SEA LEVEL CHANGE

GLACIERS AND ICE SHEETS

SEA ICE

CHANGES IN ECOSYSTEMS

HURRICANES

RISE IN TEMPERATURE

# GLOBAL ENERGY CONSUMPTION

| Energy topic ⑦ | Indicator ⑦ | Country or region |
|---|---|---|
| **Energy consumption** ⌄ | **Electricity consumption** ⌄ | **World** ⌄ |

Electricity consumption, World, 1990-2021
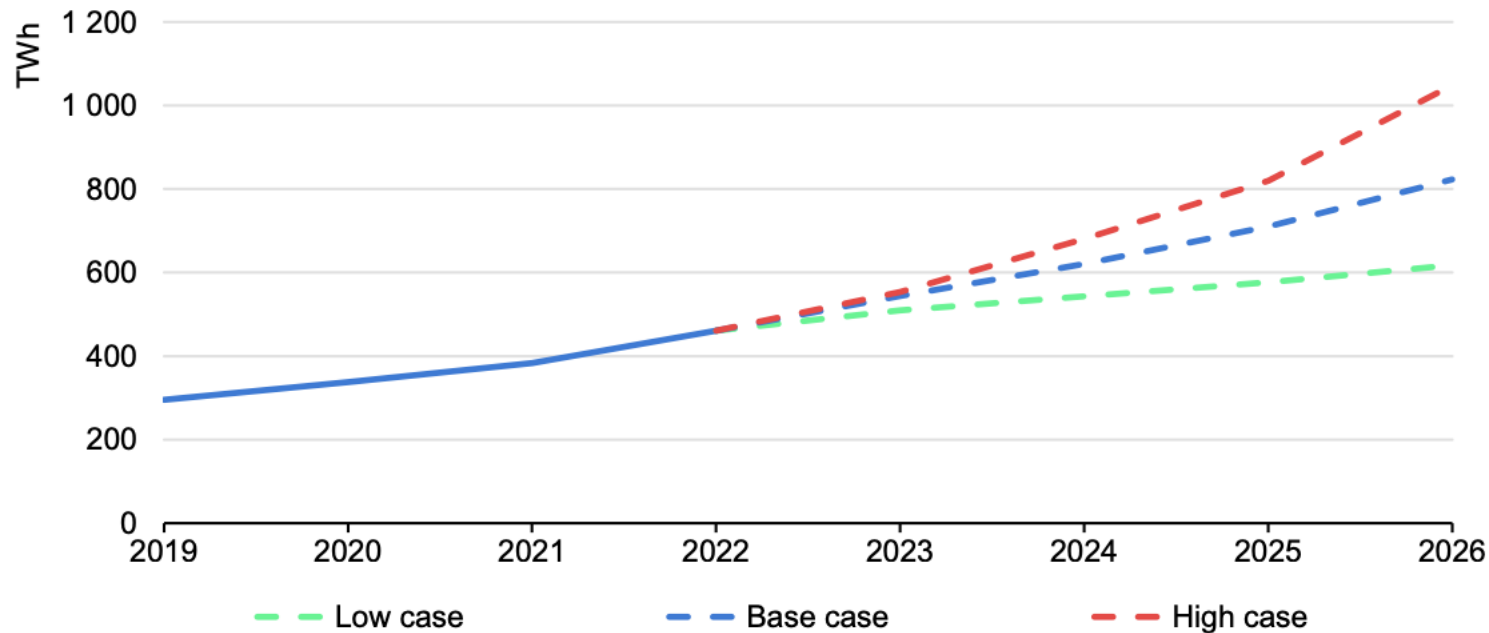


2021
● Electricity consumption: **26 467.2 TWh**

Licence: CC BY 4.0

The world's electricity consumption has continuously grown over the past half a century, reaching approximately 26,500 terawatt-hours in 2021

# GLOBAL ENERGY CONSUMPTION

**Global electricity demand from data centres, AI, and cryptocurrencies, 2019-2026**



IEA. CC BY 4.0.

Low case    Base case    High case

**Electricity consumption from data centres, artificial intelligence (AI) and the cryptocurrency sector could double by 2026.**
After globally consuming an estimated 460 terawatt-hours (TWh) in 2022, data centres' total electricity consumption could reach more than 1 000 TWh in 2026.

# THE CARBON FOOTPRINT OF TRAINING AI

How an AI training can be perfomed with a low-environmental impact?

Is focusing on training meaningful?

[1] Dodge, J., Prewitt, T., Tachet des Combes, R., Odmark, E., Schwartz, R., Strubell, E., ... & Buchanan, W. (2022, June). Measuring the carbon intensity of AI in cloud instances. In Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (pp. 1877-1894).

How an AI training can be perfomed with a

...ingful.

...ut low-carbon
...echnology is

...ensity of AI in
cloud instances»[1] is our baseline.



Google Scholar

Measuring the carbon intensity of AI in cloud instances

Articoli

In qualsiasi momento
Dal 2024
Dal 2023
Dal 2020
Intervallo specifico...

Ordina per pertinenza
Ordina per data

Qualsiasi lingua
Pagine in Italiano

Qualsiasi tipo
Articoli scientifici

Measuring the carbon intensity of AI in cloud instances
J Dodge, T Prewitt, R Tachet des Combes, E Odmark, R Schwartz, E Strubell, AS Luccioni…
Proceedings of the 2022 ACM Conference on Fairness, Accountability, and …, 2022 · dl.acm.org

The advent of cloud computing has provided people around the world with unprecedented access to computational power and enabled rapid growth in technologies such as machine learning, the computational demands of which incur a high energy cost and a commensurate carbon footprint. As a result, recent scholarship has called for better estimates of the greenhouse gas impact of AI: data scientists today do not have easy or reliable access to measurements of this information, which precludes development of

MOSTRA ALTRO ⌄

☆ Salva   99 Cita   Citato da 90   Articoli correlati   Tutte e 6 le versioni   ≫

*Visualizzazione del risultato migliore di questa ricerca. Mostra tutti i risultati*

*[1] Dodge, J., Prewitt, T., Tachet des Combes, R., Odmark, E., Schwartz, R., Strubell, E., ... & Buchanan, W. (2022, June). Measuring the carbon intensity of AI in cloud instances. In Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (pp. 1877-1894).*
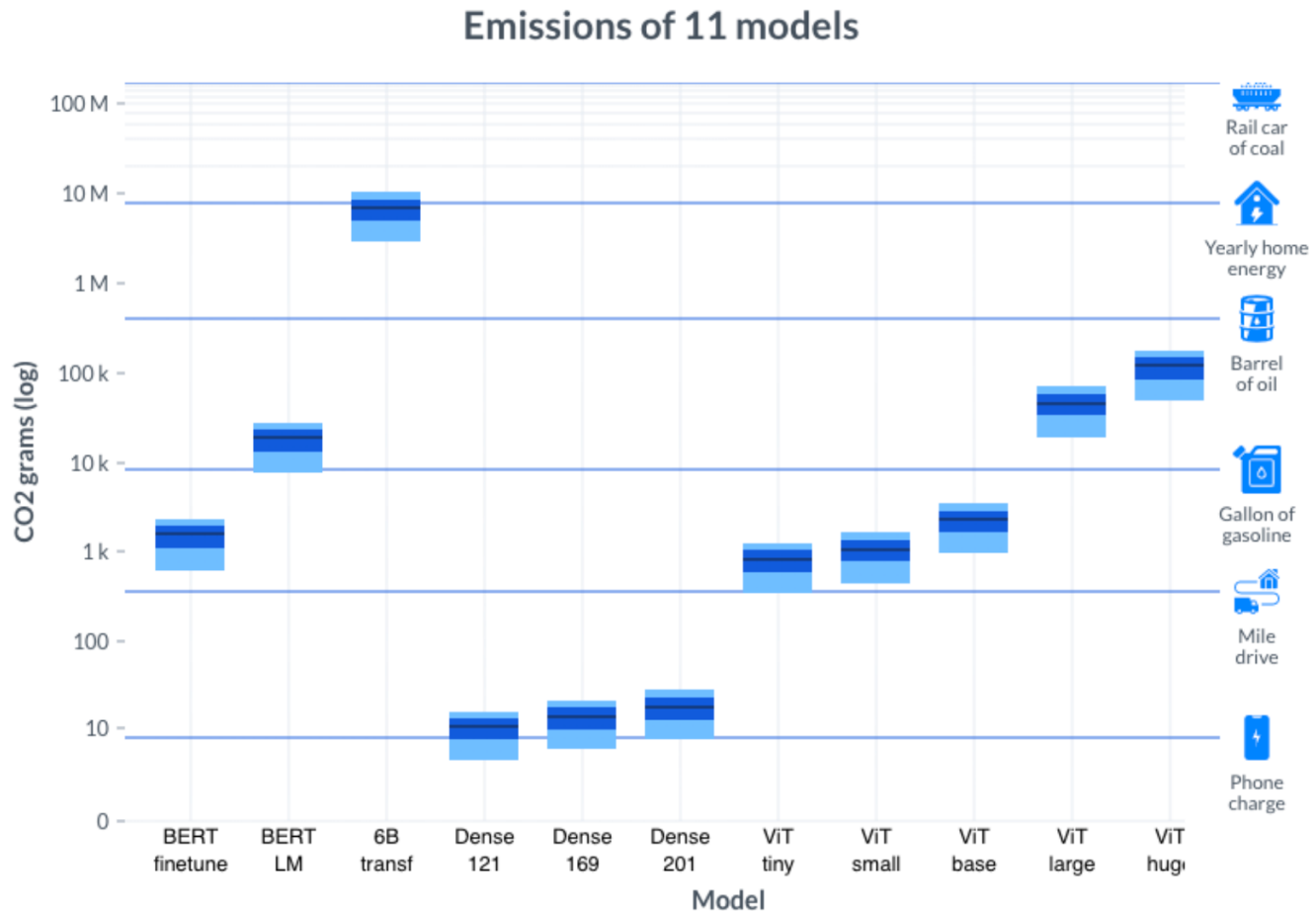
# BENCHMARKING ALGORITHMS



Emissions of 11 models

| Model | BERT finetune | BERT pretrain | 6B Transf. | Dense 121 | Dense 169 | Dense 201 | ViT Tiny | ViT Small | ViT Base | ViT Large | ViT Huge |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GPU | 4·V100 | 8·V100 | 256·A100 | 1·P40 | 1·P40 | 1·P40 | 1·V100 | 1·V100 | 1·V100 | 4·V100 | 4·V100 |
| Hours | 6 | 36 | 192 | 0.3 | 0.3 | 0.4 | 19 | 19 | 21 | 90 | 216 |
| kWh | 3.1 | 37.3 | 13,812.4 | 0.02 | 0.03 | 0.04 | 1.7 | 2.2 | 4.7 | 93.3 | 237.6 |

# BENCHMARKING ALGORITHMS



CO2 Grams Emitted, BERT Language Modeling

# MITIGATION STRATEGY 1

- FLEXIBLE-START

Launch the training at the starting time that would result in the lowest emissions.

# MITIGATION STRATEGY 2

- PAUSE AND RESUME

Run the training only during the 5-minute slots with the lowest marginal emissions.

# TUNING STRATEGY PARAMETERS

- Find the 5 minute intervals with the lowest marginal emissions during the (N + job_duration) hour window, and select enough intervals to add up to the job duration.
- Then simulate running the job only during those intervals and compute the corresponding emissions
- They explored two sets of values for N:
  - Absolute: N ∈ {6, 12, 18, 24} (hours)
  - Relative: N ∈ {25%, 50%, 75%, 100%} x job_duration

# CHOOSE REGIONS WISELY

The region that the algorithms are evaluated in has a significant impact for both strategies

For example, the West US region varies frequently throughout a single day, and thus **Pause and Resume** can lead to significant reductions.



(a) *Flexible Start* optimization for Dense 201.

(b) *Flexible Start* optimization for 6B parameters Transformer.

(a) *Pause and Resume* optimization for Dense 201.

(b) *Pause and Resume* optimization for 6B parameters Transformer.

# EVALUATION

- In order to account for daily variations (weather, electricity demand, etc.), they report the average emissions decrease computed over 5 different start times in each month, giving a total of 60 data points.
- FLEXIBLE START
  - Significant emissions reductions for shorter jobs (e.g., the DenseNet experiments
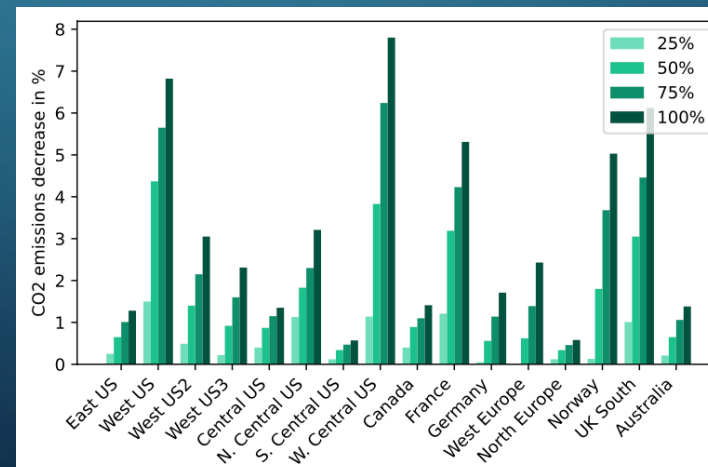  - Minimal savings for jobs longer than a day
  - Useful for use cases where an AI workload needs to run regularly, but the practitioner has some flexibility on when it runs (e.g. where models are re-trained on a regular schedule to incorporate new data over time)

# EVALUATION

- PAUSE AND RESUME
  - short experiments only see emissions reductions smaller than 10%
  - the 6 billion transformer training sees the largest decrease in emissions
  - Useful for use cases where an AI workload can be increased in duration by some proportion of the original run time

# EVALUATION

| Model | BERT finetune | BERT LM | 6B Transf. | Dense 121 | Dense 169 | Dense 201 | ViT Tiny | ViT Small | ViT Base | ViT Large | ViT Huge |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FS | 14.5% | 3.4% | 0.5% | 26.8% | 26.4% | 25.9% | 5.6% | 5.3% | 4.2% | 1.3% | 0.5% |
| P&R | 19.0% | 8.5% | 2.5% | 27.7% | 27.3% | 27.1% | 12.5% | 12.3% | 11.7% | 4.7% | 2.4% |
| Pauses / hr | 0.23 | 0.3 | 0.15 | 0.06 | 0.07 | 0.08 | 0.3 | 0.3 | 0.3 | 0.23 | 0.14 |

(14.5+3.4+0.5+26.8+26.4+25.9+5.6+5.3+4.2+1.3+0.5)/11 = 10.4%

(19+8.5+2.5+27.7+27.3+27.1+12.5+12.3+11.7+4.7+2.4)/11 = 14.5%

# OPEN ISSUES

## FLEXIBLE-START

- **more efficient for short workloads.**
  - **3.4%** emission reductions on BERT LM

## PAUSE AND RESUME

- **more efficient for workloads longer than a day.**
  8.5% emission reductions on BERT LM

Both strategies are based on the temporal management of training during

Training completion can be delayed by up to 24 hours or even more: you have to be patient!

Strategy of moving workloads to regions or times when resources are constrained.

Vergallo, R., Errico, A., & Mainetti, L. On the Effectiveness of the'Follow-the-Sun'Strategy in Mitigating the Carbon Footprint of AI in Cloud Instances. *Available at SSRN 4566638.*

Green Software Foundation
greensoftware.org

# MOTIVATIONS



Emissions depend not only on the time of day but also on the grid region where training is performed

"Follow the Sun"[2] is an approach applied to various problems, but there is no scientific validity or evidence regarding its effectiveness

Current state-of-art strategies don't preserve time

[2] Follow the Sun, GitHub. https://follow-the-sun.github.io

# OBJECTIVES

Develop a new green AI training approach by leveraging the benefits of Cloud technology.

Compare the proposed strategy with other strategies related to the same problem

The research is part of AMEDEA project (Assessment and Mitigation of the Environmental impact of DL Algorithms ) Cod. IsCa7_AMEDEA 107C

# BENCHMARK



FRAUD-DETECTION WORKLOADS

-AUTOENCODER

-HF-SCA [3]

-SVM

-ISOLATION FOREST



BANK TRANSACTION DATASET PROVIDED BY AN ITALIAN BANK WITHIN THE REGULAMENTARY SANDBOX INITIATIVE

[3] *Distante, C., Fineo, L., Mainetti, L., Manco, L., Taccardi, B., & Vergallo, R. (2022). HF-SCA: Hands-Free Strong Customer Authentication Based on a Memory-Guided Attention Mechanisms. Journal of Risk and Financial Management, 15(8), 342.*

# CARBON INTENSITY DATA

HISTORICAL CARBON INTENSITY DATA FOR YEAR 2021 FOR THE FOLLOWING REGIONS:

-MILAN

-PARIS

-FRANKFURT

-ZARAGOZA

-LONDON

-DUBLIN

-STOCKHOLM

# EXPERIMENTAL SET-UP

**Workload Runner**
- Perform training and track energy consumptions

**Strategy Launcher**
- Compute emissions for different strategies

**ISCRA** CINECA — NVIDIA DGX A100
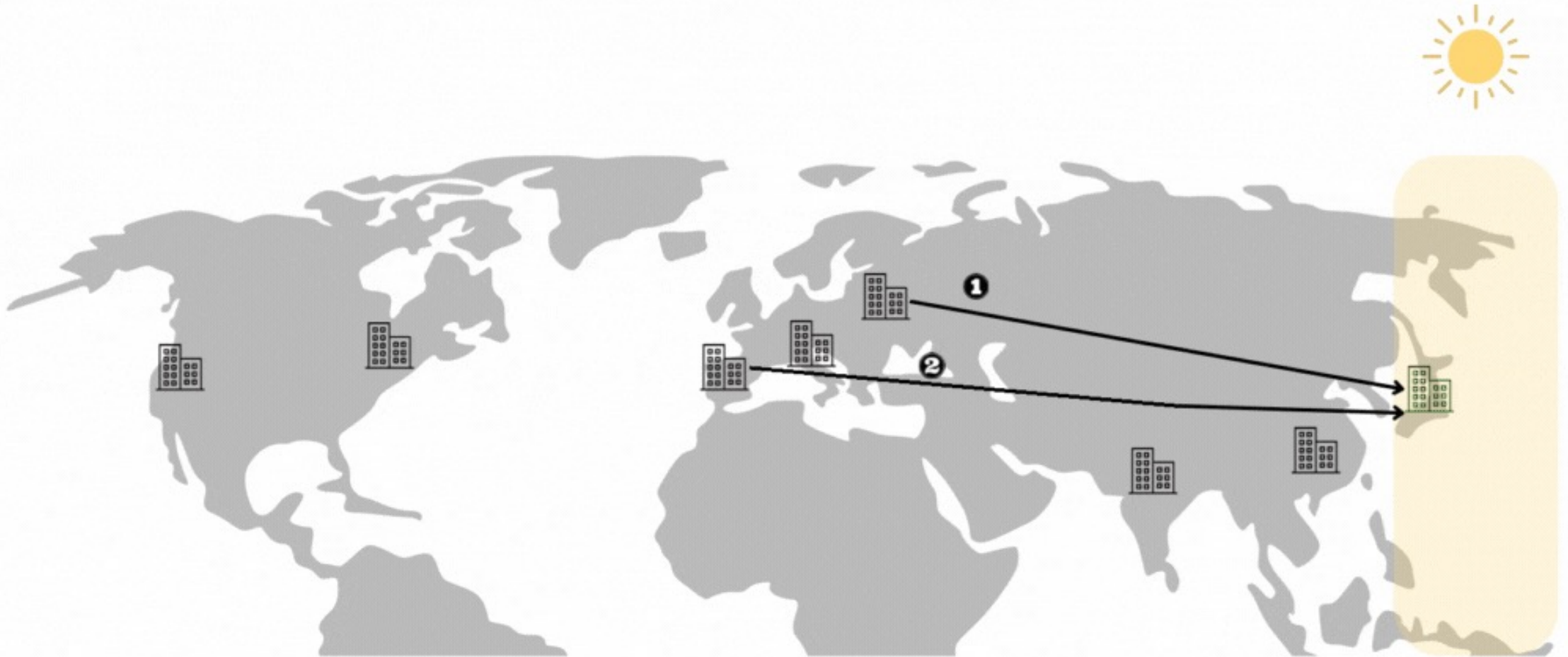
CodeCarbon python library for energy consumptions each 5 minutes

PROPOSED STRATEGY: FOLLOW THE SUN

# PROPOSED STRATEGY: FOLLOW THE SUN

## STRATEGY VERSIONS

### STATIC-START FOLLOW THE SUN

### FLEXIBLE-START FOLLOW THE SUN

## DATA TRANSFER VERSIONS

### UPSTREAM
- Transfer dataset to all Cloud Instances before the training start

### IN-TRAINING
- Transfer data during the training

# ASSUMPTIONS

Energy consumed for data transfer: 0.023 kWh/Gb from 2015 [4]

COMPRESSED DATASET SIZE: 0.320GB

WORKLOADS NOT PARTICULARLY LARGE

NEGLIGLIBLE TRANSFER TIME

NEGLIGIBLE EMISSIONS FOR TRAINING STATE DATA TRANSFER

[4] Malmodin, J., & Lundén, D. (2018). The energy and carbon footprint of the global ICT and E&M sectors 2010-2015. Sustainability, 10(9), 3027.

# PROPOSED STRATEGY: FOLLOW THE SUN

## GENERAL IDEA

- Checking-time: how often to designate the new region to transfer the training to
- The workload is divided into k slots based on the selected checking time
- Each slot corresponds to a training segment that will be executed on the region with the least environmental impact

# PROPOSED STRATEGY: FOLLOW THE SUN

## 1° version

### Static-start Follow The Sun

- The workload will start at the specified starting time, and at each checking time, it will be moved to the greenest region to continue with the training.

| REGION 1 | 00:00 | 00:05 | 00:10 | 00:15 | 00:20 | 00:25 | 00:30 | 00:35 | 00:40 | 00:45 | 00:50 | 00:55 | 01:00 | 01:05 | 01:10 | 01:15 | 01:20 | 01:25 | 01:30 | 01:35 | 01:40 | 01:45 | 01:50 | 01:55 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| REGION 2 | 00:00 | 00:05 | 00:10 | 00:15 | 00:20 | 00:25 | 00:30 | 00:35 | 00:40 | 00:45 | 00:50 | 00:55 | 01:00 | 01:05 | 01:10 | 01:15 | 01:20 | 01:25 | 01:30 | 01:35 | 01:40 | 01:45 | 01:50 | 01:55 |
| REGION 3 | 00:00 | 00:05 | 00:10 | 00:15 | 00:20 | 00:25 | 00:30 | 00:35 | 00:40 | 00:45 | 00:50 | 00:55 | 01:00 | 01:05 | 01:10 | 01:15 | 01:20 | 01:25 | 01:30 | 01:35 | 01:40 | 01:45 | 01:50 | 01:55 |

# PROPOSED STRATEGY: FOLLOW THE SUN

## 2° version

## Flexible-start Follow The Sun

- **Starting at the time that minimizes emissions without interrupting the execution**

# EVALUATION

Reductions increase as the time window increases

As the checking time decreases, the reductions increase

# EVALUATION

## STATIC-START
## FOLLOW THE SUN

- Average reduction percentages ranging from 5% to 7%

## FLEXIBLE-START
## FOLLOW THE SUN

- Average reduction is between 14-16% for the shorter workloads and almost 10% for the longest one
- Showed peaks of reductions beyond 81%
- Flexible-Start is a lower-bound: we can only do better!

# EVALUATION

## STATIC-START FOLLOW THE SUN

⬇

## PRESERVE THE ENTIRE WORKLOAD DURATION

| Strategy | HF-SCA | Autoencoder | SVM | Isolation Forest |
|---|---|---|---|---|
| No-Strategy | 16h | 3:30h | 2:30h | 4:15h |
| Flexible Start | 19:15h | 12:56h | 12:45h | 13:25h |
| Pause and Resume | 21:33h | 16:57h | 15:50h | 17:34h |
| Static-Start Follow the Sun | 16h | 3:30h | 2:30h | 4:15h |
| Flexible-Start Follow The Sun | 37:32h | 17:31h | 17:47h | 16:22h |

# EVALUATION

## EMISSION REDUCTION

- Flexible-Start Follow the sun has the best average percentage reduction between all strategies

## TIME SAVING

- Static-Start version perserve the workloads length. Other strategies does not consider this opportunity at all

| Strategy | Avg time dilatation | Avg carbon reduction |
|---|---|---|
| Flexible-Start[1] | 7:54h | 5.72% |
| Pause and Resume[1] | 11:15h | 6.51% |
| Static-Start FtS | No dilatation | 5.925% |
| Flexible-Start FtS | 15:36h | 13.85% |

# PROBLEMS?

# EVALUATION

**ROBUSTNESS WITH RESPECT TO THE SET OF REGIONS**

- the Follow the sun strategy gives the same result no matter what the starting region is

**GDPR LIMITATIONS**

- The proposed strategy could be subjected to these kind of constrain

**ARCHITECTURE COMPLEXITY**

- The strategy proposed requires a more complex infrastructure

| Strategy | GDPR limitations | Complex architecture required | Robustness wrt regions |
|---|---|---|---|
| Flexible-Start[1] | No | No | No |
| Pause and Resume[1] | No | No | No |
| Static-Start FtS | Yes | Yes | Yes |
| Flexible-Start FtS | Yes | Yes | Yes |

# ACCURACY vs SUSTAINABILITY

| HF-SCA | Autoencoder | SVM | Isolation Forest |
|--------|-------------|-----|------------------|
| 0.97 | 0.73 | 0.51 | 0.56 |

## Consideration on the experiment

- In this case HF-SCA has much better AUC score and it is preferable despite high emissions
- In cases of more comparable performance, prefer the greener model

# CONCLUSION

## EMISSIONS REDUCTION

- 13.85% of average reductions between workloads, vs 6.51% for state of art strategies

## TIME SAVING

- You can avoid wasting time for reducing emissions

## ROBUSTNESS WRT STARTING REGION

- Same results regardless the starting region

# THANK YOU!

# Roberto Vergallo

# r.vergallo@tudelft.nl