

# Experiments for validating our tool

Clément Morand

May 17, 2023

## Contents

<b>1</b>	<b>Checking that we can get the same dynamic consumption estimate as Green Algorithms</b>	<b>3</b>
<b>2</b>	<b>replicating results from [?]</b>	<b>4</b>
2.1	detailing the Hardware configurations . . . . .	4
2.2	Problems with the provided data . . . . .	4
2.2.1	Trying to understand the problem . . . . .	5
2.3	experiments . . . . .	5
<b>3</b>	<b>replicating results from [?]</b>	<b>6</b>
<b>4</b>	<b>replicating results from [?]</b>	<b>8</b>
4.1	Trying to find information about the hardware setup . . . . .	8
4.1.1	Hardware for the fine-tuning . . . . .	8
4.1.2	Hardware for training the models . . . . .	9
4.2	coherency of the results . . . . .	9
4.3	Estimating energy consumption . . . . .	9
4.3.1	fine tuning of the SSL model . . . . .	9
4.3.2	Table 1 . . . . .	10
<b>5</b>	<b>results from [?]</b>	<b>10</b>
5.1	Hardware configuration . . . . .	11
5.2	running experiments . . . . .	11
5.3	Explaining the massive differences between our estimates and the expected results . . . . .	11
5.4	table from [?] . . . . .	12
5.4.1	Carbon intensity used . . . . .	12
5.4.2	results . . . . .	12

5.5	New experiment : . . . . .	13
<b>6</b>	<b>estimations from [?]</b>	<b>14</b>
6.1	Information about the hardware configuration . . . . .	14
6.2	Reproducing figures from table 3 . . . . .	15
6.2.1	Checking the Coherency of the presented results . . .	15
6.2.2	running our estimations . . . . .	15
6.2.3	integrating Life cycle to previous analyses . . . . .	17
<b>7</b>	<b>comparing manufacture impacts with Dell LCAs</b>	<b>18</b>
7.1	Dell R740 . . . . .	18
7.2	Dell R6515, R7515, R7525 . . . . .	19
<b>8</b>	<b>replicating the Bloom estimates from [?]</b>	<b>21</b>
8.1	Gathering information about the setup . . . . .	21
8.2	comparing the server footprint with the PCF sheet. . . . .	22
8.3	comparing the GPU footprint with the chosen value . . . . .	22
8.4	Estimating the total impacts . . . . .	23
<b>9</b>	<b>Conclusions</b>	<b>23</b>
9.1	about the replication of results . . . . .	23
9.2	about the validity of the tool . . . . .	24
<b>10</b>	<b>TODO :noexport:</b>	<b>24</b>

In order to validate our tool we first need to conduct some experiments to ensure that it produces results consistent with the state of the art. Our tool estimates manufacture impacts of the hardware used and energy consumption over the usage duration (typically during the training phase of a model). We therefore want to test those two parts. Firstly, we will present experiments aimed at testing the estimates for the dynamic energy consumption (and the results we will present will therefore focus only on the energy consumption estimated and on the Carbon footprint induced by this energy consumption). These experiments will start by reproducing the same exact results as the Green Algorithms tool, by first choosing a scenario where we know that our tool and Green Algorithms use the same data. Then, we will focus on reproducing results presented in the two surveys of existing tools ([?] and [?]). Finally, we will try to reproduce results obtained with a different method than the one used by our tool. This will be done by trying to reproduce results from [?], [?] and [?]. We chose those articles because they all use a different tool (even if measures presented in these three articles are based on

the RAPL and NVLM tools), because we were able to contact the authors of [?, ?] to obtain further details about the hardware configuration they use ; Because [?] presents results about the inference phase of models which is a phase that is rarely studied and because [?] was the paper that made NLP researchers consider the impacts of the models they produced. Secondly, we will compare the results our tool produces with LCA results produced by Dell about the impacts of the servers they sell. This will allow us to validate the estimations of embeded impacts our tool generates. Finally, we will try to reproduce the results from [?], this step is really important because this paper conducts an analysis of the global warming potential induced by the Bloom model. This analysis takes into account embodied emissions and we use figures they present to define the default dynamic ratio our tool uses.

## 1 Checking that we can get the same dynamic consumption estimate as Green Algorithms

To do a first sanity check, we verify that we are able to reproduce the same results as GA on the dynamic consumption part :

We choose a configuration that we know is available in both databases (GA version 2.2 at the time of this experiment):

- 1 CPU A8-7680 (4 cores)
- 1 GPU NVIDIA GTX 1080 Ti
- 64 GB Memory
- Use time of 12h 0min
- no PUE / dynamic ratio
- carbon intensity of France is used (51.28 g CO<sub>2</sub> e/kWh)

We are using Green Algorithms v2.2 for an expected result of 196.32g of CO<sub>2</sub> e and 3.83 kWh of dynamic consumption (this link should in theory get you to the page with this exact setup and results but it seems like GA sharing feature is broken right now).

If we now run the experiment with our tool : we see that we indeed obtain the same results of 196gCO<sub>2</sub> e and 3.83 kWh of dynamic energy consumption.

## 2 replicating results from [?]

In order to replicate results, we first need to gather some information about the hardware configuration used to run the experiments. Then, we will face the challenge of incoherences in the data presented in the paper. Finally, we will be able to run experiments that give the same energy consumption estimates as those presented in the paper.

### 2.1 detaillig the Hardware configurations

The authors provided us with informations about the hardware configurations used to run the experiments.

the facility setup is the LaBia cluster. We can see that the only nodes using a 20 core CPU are: n[101-102]:

- 2 x Intel Xeon Gold 6148 20 cores / 40 threads @ 2.4 GHz (Skylake)
- 384 GiB of RAM
- 4 x NVIDIA Tesla V100 with 32 GiB of RAM (NVLink)

using 32 GB of RAM and not the full 384.

The lab server on the other hand is using one GTX 1080 Ti with 11GB of memory. it is a Dell PowerEdge R730 with 2 GTW 1080 Ti, 2 Intel Xeon E5-2620 v4 CPU and 125 GB memory (only 11 of whihch are requested).

while we do not have the Intel Xeon Gold 6148 in our CPU database, we can see on Intel’s website that it has a TDP of 150W, was realeased in 2017 with a process of 14nm with the Skylake architecture, this is sufficient information to add one entry to our database, knowing the information about the Skylake architecture from WikiChips.

### 2.2 Problems with the provided data

Results presented in the paper do not seem coherent from one table to the other (tables 3 and 4). If we try to convert from energy consumption to carbon emissions using the presented carbon intensity of 39 gCO<sub>2</sub> e/kWh we do not at all find the same results as the ones presented. For instance, for the first method (Yu2020) for the French Press benchmark, it is indicated 1.38kWh consumption and 350.15g CO<sub>2</sub> e.

We can see that if we are to use the presented carbon intensity, we get emissions of 53.8 gCO<sub>2</sub> e for a 1.38kWh energy consumption. This is really far from the 350 gCO<sub>2</sub> e presented in the paper.

### 2.2.1 Trying to understand the problem

Let us check if the factor to convert from table 4 to table 3 is constant. If it is, it would maybe explain the problems. When filling the table the authors might have missclicked on the location and the Carbon Intensity used would just be the one of another country.

We obtain results around 250 gCO<sub>2</sub> e/kWh with some non negligible variations (The smallest conversion factor is of 191.5 gCO<sub>2</sub> e/kWh while the highest is of 283.5 gCO<sub>2</sub> e/kWh)

according to GA's v2.2 database, this carbon intensity of around 250gCO<sub>2</sub> e/kWh would approximately correspond to Lithuania's one. According to the version 1.1 of the data (version seemingly used in the article), the closest one would be Hungary.

Still, we can observe quite important variations in carbon intensity to convert from the presented energy consumption to the presented carbon emissions, this would tend to infirm the hypothesis of just an error of selection in the carbon intensity used.

Even if there are obviously problems with the presented data, we still want to try and replicate the presented results. Indeed, if the data is flawed only in the table presenting the energy consumption or only in the table presented the carbon footprint, we might be able to reproduce the results of one of the tables (i.e. either the consumption or the carbon footprint)

## 2.3 experiments

It is said that the default PUE used is 1.67. In order to replicate the results, and even if the dynamic ratio and the PUE do not have the same meaning. Since they are both used in the same way we will use a dynamic ratio of 1.67

we can see in the latest version of Green Algorithms' GPU TDP database that they have a TDP value of 300W for a Tesla V100 GPU whereas we have a TDP of 250W for the same card in our database. In order to see if we can replicate the same consumption and see the difference resulting from this data-point incoherency we will try two versions. One with a V100 and one with a card with a TDP of 300W in our database: the NVIDIA A100 PCIe 80 GB. This will of course also impact the manufacture impacts but we are here only focusing on reproducing the same direct impacts

Method	Task	Hardware	Expected Energy (kWh)	Estimated Energy (kWh)	Estimation trying to match Facility only (kWh)	Expected Carbon (gCO2e)	Estimate Carbon (gCO2e)
Yu2020	French Press	Server	1.38	1.16	nan	350.15	45.15
Yu2020	French Press	Facility	1.03	0.861	1.03	260.26	33.26
Yu2020	EMEA	Server	0.07	0.0673	nan	16.67	2.67
Yu2020	EMEA	Facility	0.06	0.0499	0.0595	14.31	1.95
Yu2020	MEDLINE	Server	0.08	0.0843	nan	20.68	3.28
Yu2020	MEDLINE	Facility	0.08	0.0669	0.0797	20.03	2.69
Ma2016	French Press	Server	0.41	0.414	nan	104.4	16.44
Ma2016	French Press	Facility	0.4	0.341	0.406	102.08	13.08
Ma2016	EMEA	Server	0.02	0.0158	nan	3.8	0.61
Ma2016	EMEA	Facility	0.02	0.0179	0.0213	4.99	0.69
Ma2016	MEDLINE	Server	0.02	0.0225	nan	5.57	0.87
Ma2016	MEDLINE	Facility	0.02	0.0216	0.0258	5.67	0.84

We can see that we are able to obtain the same exact energy consumption estimates up to rounding (when we do the modifications to the inputted setup for the facility) except for Yu2020, French Press, Server where we have a slightly lower estimation than the one proposed in the paper. We can also see that, as expected, the estimates we do when considering the "real" setup are lower than the ones presented in the paper and this can be entirely explained by the difference in TDP in the database. We can also conclude that the problem in the presented data lies in the estimates of the carbon footprint and not in the estimates of energy consumption.

### 3 replicating results from [?]

In order to replicate the results from the paper, we first need to gather some information from the paper and its supplementary material which is designed to allow for reproducible experiments.

- The hardware used is a Nvidia DGX-1 with two Intel Xeon E5-2698 v4,

512 GB of memory and 8 NVIDIA Tesla V100-SXM2-32GB.

- The Carbon Intensity for France used in Green Algorithms V2.2 is

51.28gCO<sub>2</sub> e/kWh (latest version of Green Algorithms' Carbon Intensity Database)

- To convert from kWh to kJ, one must multiply the result by 3.6E+3.

we can see in the latest version of Green Algorithms' GPU TDP database that they have a TDP value of 300W for an NVIDIA V100 GPU whereas we have a TDP of 250W for the same card in our database. As a first version, just to see if we are able to obtain the same exact results as those presented in the paper, we will use as GPUs a card with a TDP of 300W in our database: the NVIDIA A100 PCIe 80 GB.

We can also see that the CPU model used is the Xeon E5-2698 v4 with a tdp 135. However, it isn't available in Green Algorithm, the model used is the Xeon E5-2697 v4 with a TDP of 145W and 18 cores. In order to reproduce the results presented in the paper, we will use in our setup one CPU with 40 cores, a TDP of 324W (145/18\*40) and a die size of 9.12cm<sup>2</sup> (2\*the die size of a Xeon E5-2698 v4, not relevant for the computation of energy)

In the notebook accompanying the paper, we can see that the link explaining the configuration used for the CPU benchmarks are exact copies of the ones for GPU benchmarks. We will therefore assume that the cpu usage was 1 and gpu usage was 0. This configuration leads to an energy consumption of 8.58Wh for one minute. Since this value is strangely similar to the value of 7.58Wh/min used in the paper, we will also assume that there was a mistake when copying results from the Green Algorithm website and therefore use the value of 8.58Wh/min instead of the value of 7.58Wh/min to compute the expected results.

When trying to obtain the exact same results (same hardware setup as used for obtaining values with Green Algorithms) for the GPU benchmark

Benchmark	Value (kJ)	Difference (kJ)
EP	176.04	-0.134
LU	381.6	-0.043
MG	134.28	-0.049

for the CPU benchmarks

Benchmark	Value (kJ)	Difference (kJ)
EP	25.74	0
LU	15.444	0
MG	64.44	0.09

When using the hardware setup really used: for the GPU benchmark

Benchmark	Value (kJ)	Difference (kJ)
EP	149.04	-27.134
LU	324.36	-57.283
MG	117	-17.329

for the CPU benchmarks

Benchmark	Value (kJ)	Difference (kJ)
EP	23.04	-2.7
LU	13.824	-1.62
MG	57.6	-6.75

We can see that we are able to obtain results that are exactly the same as the expected ones up to rounding errors (difference 3 orders of magnitude lesser than the value). We can also see that even though the input value to Green Algorithms does not exactly correspond to the hardware setup used, the difference to the expected results isn't too high. The difference between our estimate using 'correct' data and the expected values is around 10 times less than the estimated value.. These results demonstrate the importance of inputting the right hardware if one wants precise results.

## 4 replicating results from [?]

### 4.1 Trying to find information about the hardware setup

The authors gave us some insight on the hardware used for running their experiments. Without their help, we would not have been able to produce a single estimate

#### 4.1.1 Hardware for the fine-tuning

They said that a node from the Jean Zay supercomputer with 4 GPUs with 32GB memory was used for the fine tuning of the wave2vec model. if we look at the Idris' website we think that the nodes used were from the **v100-32g**, it is the only node with matching requirements in terms of number of GPU and memory per GPU.

these nodes have the following hardware configuration :

- 2 Intel Cascade Lake 6248 (20 cores at 2,5 GHz)
- 192 GB de memory per node



- 4 GPU Nvidia Tesla V100 SXM2 32 GB

Because we do not have the Intel Cascade Lake 6248 in our database, we need to find some information about it. We can see on Intel's webpage that it is a processor of the Cascade Lake architecture. On Wikichip, we can see that Cascade Lake Processors use dies largely similar to those of the Skylake cores. Combining all of these pieces of information, we can get an estimation of the details of an Intel Cascade Lake 6248 : model: "Xeon Gold 6248" manufacture date: "2019" process: 14nm number of cores: 20 die size: 694 mm<sup>2</sup> (XCC configuration)

#### 4.1.2 Hardware for training the models

We are told that training uses only one GPU at a time and that it uses roughly half of the time a RTX 2080 Ti and the other half a GTX 1080 Ti, to represent this, we will put the two different models in the list of GPUs and use a 'gpu usage' of .5. We are also told that the training uses 80 GB memory with no additional information on the hardware used. Since we do not know any more precise information, we will use the default values of our tool to complete the missing pieces of information

### 4.2 coherency of the results

One first good news is that information are coherent with themselves. Using the indicated (in the paper) carbon intensity of 51gCO<sub>2</sub> e/kWh used and indicated energy consumption, we are able to find back the carbon emissions indicated in the table. The only problem is that for table 1, it seems that there was a translation error when filling the table. The figures are written in the french notation with "," separating units from decimals and not the usual ".". For instance, if we look at the first line of table 1, we can read a consumption of 4,473 kWh, that we can translate to 4.473 kWh. We obtain =228.123=g CO<sub>2</sub> e, the same value as indicated in the paper.

We then only need to be able to find coherent energy consumption values to obtain comparable results.

### 4.3 Estimating energy consumption

#### 4.3.1 fine tuning of the SSL model

We can see that we obtain an estimate of 5.46kg CO<sub>2</sub> e for the direct impacts and a dynamic consumption of 107 kWh, which is close to the 4.729kg CO<sub>2</sub> e and 97.720 kWh presented in the paper. The fact that results aren't a perfect

match and slightly higher than presented can be explained by the fact that measures presented were carried out based on a measurement tool (Carbon-Tracker). (results presented are borrowed from [?] using the methodology from [?])

**4.3.2 Table 1**

model	expected power (kWh)	estimated power (kWh)	expected carbon (gCO2e)	estimated carbon (gCO2e)
spectro 3 steps	4.473	10.1	228	53
spectro 2 steps	2.989	6.78	152	34
spectro 1 step	1.708	4.44	87	23

We can see that we obtain carbon emission estimates around 3 times higher than those presented in the paper. It is expected that we obtain higher estimates than the measurements as presented in [?]

## 5 results from [?]

We try to replicate the following results:

Coûts écologiques et énergétiques passés à l'échelle Steps Inférences sur 1 journée (27 Millions d'appels)

Tasks Models	MEDIA			ATIS-FR		
	Time (Heures)	Energy (MWh)	CO2 (Kg)	Time (Heures)	Energy (MWh)	CO2 (Kg)
FlauBERTbase	20.19	204.24	147.84	3.08	30.88	22.33
CamemBERTlarge, CCNet 135 Gb	50.63	512.67	371.14	7.36	74.23	53.75
CamemBERTbase, OSCAR 138 Gb	20.23	204.67	148.15	3.27	32.57	23.56
CamemBERTbase, CCNet 135 Gb	15.57	157.39	113.96	2.55	24.79	17.94
CamemBERTbase, OSCAR 4 Gb	15.89	160.70	116.35	2.52	25.18	18.25
CamemBERTbase, CCNet 4 Gb	15.64	158.08	114.42	2.59	25.49	18.48
CamemBERTbase, Wiki 4 Gb	15.38	155.46	112.57	2.50	24.95	18.10
FrALBERTbase, Wiki 4 Gb	9.11	92.02	66.61	1.39	13.71	9.93
XLM-Rbase	17.20	173.94	125.90	2.40	25.72	18.63
XLM-Rlarge	55.68	563.95	408.25	8.02	76.08	58.60
mBERTbase	17.95	181.41	131.36	2.48	24.72	17.94
distill-mBERTbase	15.06	152.08	110.11	2.35	23.25	16.79
small-mBERTbase-fr	16.45	166.24	120.35	2.46	24.56	17.79

## 5.1 Hardware configuration

We were told that the hardware used was an NVIDIA DGX equipped with 8 NVIDIA Tesla V100 SMX2 16GB. I was not able to find such a configuration on NVIDIA’s Website but since the Tesla V100 SMX2 32GB GPU present in an NVIDIA DGX-1 server have the same exact TDP, we will suppose that this is the hardware used.

## 5.2 running experiments

We can run the experiments and compare the results with the expected results:

Tasks Models	MEDIA Time (Heures)	MEDIA Expected Energy (MWh)	MEDIA Estimated Energy (MWh)	MEDIA Expected C (P
FlauBERTbase	20.19	204.24	0.0442	147
CamemBERTlarge, CCNet 135 Gb	50.63	512.67	0.111	371
CamemBERTbase, OSCAR 138 Gb	20.23	204.67	0.0443	148
CamemBERTbase, CCNet 135 Gb	15.57	157.39	0.0341	113
CamemBERTbase, OSCAR 4 Gb	15.89	160.7	0.0348	116
CamemBERTbase, CCNet 4 Gb	15.64	158.08	0.0343	114
CamemBERTbase, Wiki 4 Gb	15.38	155.46	0.0337	112
FrALBERTbase, Wiki 4 Gb	9.11	92.02	0.02	66
XLM-Rbase	17.2	173.94	0.0377	12
XLM-Rlarge	55.68	563.95	0.122	408
mBERTbase	17.95	181.41	0.0393	131
distill-mBERTbase	15.06	152.08	0.033	110
small-mBERTbase-fr	16.45	166.24	0.036	120

We can see that we obtain results as low as 4 orders of magnitude lower than the expected results. This massive difference cannot be easily explained and is a really surprising result.

## 5.3 Explaining the massive differences between our estimates and the expected results

In our estimates, the consumption of one DGX-1 is estimated at 2460W (if we were to suppose that CPUs are running at full capacity) and 2190W if we suppose that CPUs do not run. This is significantly lower than the 3500W

provided by NVIDIA and can be due at least in part to the fact that we do not account for storage in our estimation.

Results are way lower than those presented. however, the presented results seem at least surprising. If we use the consumption value provided by NVIDIA of 3500W for one DGX-1 DGX-1 datasheet. If used for 8 hours like for ATIS-FR with XLM-Rlarge, we would expect a consumption of 28kWh. This is extremely far from the 76MWh presented. There is therefore a problem in the expected data or (more probably) in the hardware configuration used.

Furthermore we can see that conversion from energy consumption to carbon emissions make us remark that the carbon intensity seemingly used is approximately 1.38 gCO<sub>2</sub> e/kWh. This is extremely low as the Carbon Intensity for France is estimated between 50 and 200 gCO<sub>2</sub> e/kWh

## 5.4 table from [?]

It is said that only one V100 GPU is used for training the different models. (we will suppose that it was done on one DGX-1 server)

### 5.4.1 Carbon intensity used

```
import numpy as np
energy = [1.08,3.10,.57,1.14,3.30,1.07,1.09,1.06]
emission = [317.87,914.27,167.8,337.70,973.29,317.02,321.42,314.17]
print(np.mean([em / en for en, em in zip(energy, emission)]))

295.2935224349162
```

We can see that the carbon intensity used seems to be of 295 gCO<sub>2</sub> e / kWh.

### 5.4.2 results

we can see on Experiment-Impact-Tracker's repository that they by default use a PUE of 1.58, in order to replicate their results. We will choose to use this value of 1.58 as dynamic ratio.

We can suppose that during training only the GPU is used at full capacity. we can also try a scenario where one core of one GPU is used during training. This would lead to including a cpu usage of 1/20 (since the CPU has 20 cores).

Table ?? presents the results of these experiments :

model	estimate	time (s)	expected energy (kWh)	estimated energy (kWh)	expected carbon (kgCO2e)	estimated carbon (kgCO2e)
CamemBERT <sub>base</sub>	lower	7207	1.08	1.41	0.317	0.415
CamemBERT <sub>base</sub>	upper	7207	1.08	1.43	0.317	0.421
CamemBERT <sub>large</sub>	lower	19445	3.1	3.77	0.914	1.11
CamemBERT <sub>large</sub>	upper	19445	3.1	3.83	0.914	1.13
FrALBERT <sub>base</sub>	lower	3816	0.57	0.75	0.167	0.221
FrALBERT <sub>base</sub>	upper	3816	0.57	0.761	0.167	0.225
XLM-R <sub>base</sub>	lower	7676	1.14	1.5	0.337	0.441
XLM-R <sub>base</sub>	upper	7676	1.14	1.52	0.337	0.448
XLM-R <sub>large</sub>	lower	21137	3.3	4.1	0.973	1.21
XLM-R <sub>large</sub>	upper	21137	3.3	4.16	0.973	1.23
mBERT <sub>base</sub>	lower	7333	1.07	1.43	0.317	0.422
mBERT <sub>base</sub>	upper	7333	1.07	1.45	0.317	0.428
samll-mBERT <sub>base</sub>	lower	7190	1.09	1.4	0.321	0.414
samll-mBERT <sub>base</sub>	upper	7190	1.09	1.42	0.321	0.42
distil-mBERT <sub>base</sub>	lower	6466	1.06	1.26	0.314	0.372
distil-mBERT <sub>base</sub>	upper	6466	1.06	1.28	0.314	0.378

We can see that for upper and lower estimates we obtain results slightly higher than those presented in the paper but in the same order of magnitude. This is expected since estimation tools tend to provide higher (and closer to reality) estimates than measurement tools. However, we can also see that the estimation tool ([?]) does not capture some subtulties. For instance small-mBERT<sub>base</sub> training is quicker than mBERT<sub>base</sub> one. However this does not translate to smaller energy consumption most probably because one model training uses more ressources than the other one. Without fine knowledge of the processing units usage, we cannot provide very precise estimations and track small changes such as this one.

All of these results tend to confirm that there are problems with the data available in [?] but that the data from [?] confirms us the hardware configuration used.

## 5.5 New experiment :

A new experiment was run on the Segur machine. The following results were obtained using Experiment-Impact-Tracker :

cpu <sub>hours</sub>	1.0428555555555554
gpu <sub>hours</sub>	0.9933892874755572
estimated <sub>carbonimpactkg</sub>	0.024094323442314113
total <sub>power</sub>	0.4302695971583645
kw <sub>hrgpu</sub>	0.2516560133562949
kw <sub>hrcpu</sub>	0.02066651649077121
explen <sub>hours</sub>	0.5388999266756905

from these results, and knowing that the Segur machine is equipped with 20 core CPUs with 125 GB RAM and 2 GTX 1080 Ti, we can estimate that approximately 2 cores (1.04/.53) were used at full capacity during training, which equates to 1/20 usage. The two GPU also seem to have been used at full capacity. we can deduce the used Carbon Intesity by dividing the estimated carbon by the measured power

this result of 56 gCO<sub>2</sub> e/kWh lead us to think that the Carbon Intensity of France was used. (which would be logical since the experiment was run in France)

We also know that Experiment Impact Tracker uses a PUE of 1.58, in order to try and reproduce these results, we will use a dynamic ratio of 1.58. We will also try with the base dynamic ratio and see the difference

All of this allows us to run the following experiment to try and reproduce these results

	Expected	Estimated	Match
energy (kWh)	0.43	0.855	0.436
Carbon (kgCO <sub>2</sub> e)	0.0241	0.0479	0.0244

We can see that we obtain very close results (a little bit higher just as expected) when trying to get an exact match by using a dynamic ratio of 1.58 and estimates are approximately doubled when using the base dynamic ratio which stands around 3.

## 6 estimations from [?]

### 6.1 Information about the hardware configuration

It is described in the paper that estimates are conducted by training all models for a maximum of 24h. They use RAPL and NVIDIA System Management Interface to measure the average consumption of the CPUs and GPUs. All models are trained on one NVIDIA TITAN X except for ELMo which is trained on 3 GTX 1080 Ti. They then transcribe these results to

estimates by using the training time given in the paper and the description of the hardware given in the paper.

No figures are presented regarding the average consumption of the memory, CPU and GPU (separated). We only know about the model of GPU used for estimating the consumption and the total estimated consumption for training each model. We will therefore not give any value for the CPU and ram and run our estimates as is. We will see what results we obtain. We would like, not to obtain exact results since it wont be possible given the informations missing. Since they use measurement tools, we can think that using a modelisation using the TDP will give an higher result but since we do not know the quantity of memory used and the CPU used, we are not sure that the results will be higher (even if we can hypothesize that the CPU average consumption is negligible compared to the GPU consumption.)

One reassuring point is that CTX 1080 Ti, V100, P100 and Titan X GPUs have the same TDP so the consumption estimated should make sense.

They use a PUE of 1.58 and a Carbon Intensity of 0.954 pounds CO<sub>2</sub> e/kWh for American electricity production which is equivalent to 432.72 g CO<sub>2</sub> e/kWh.

## 6.2 Reproducing figures from table 3

### 6.2.1 Checking the Coherency of the presented results

Since there are no estimates given for models trained on TPUs, we will in the first time at least ignore these models.

Since table 3 presents the estimated consumption used, we can first check the coherency of the table by seeing if we can reproduce the same energy consumption by multiplying the power by the training time and the PUE

We can see that, up to rounding we obtain the same results. We can also check that we obtain the same carbon emissions.

Also the same up to rounding errors

### 6.2.2 running our estimations

For a first check, we will compare the estimated energy consumption of just the GPUs with the presented hardware consumptions. The TDP of a P100 GPU is 250W, also the same as the one of a GTX 1080 ti.

model	estimated	measured
Transformer <sub>base</sub>	2000	1415.78
Transformer <sub>big</sub>	2000	1515.43
ELMo	750	517.66
BERT <sub>base</sub>	16000	12041.51
NAS	2000	1515.43

We can see that, as expected since the provided consumption result from using measurement tools, the estimated consumption is bigger (approximately + 1/3) than the measured consumption. Still, it remains in the same order of magnitude

First let us convert the expected results from pounds to kg.

model	estimated pounds	estimated kg
Transformer <sub>base</sub>	26	11.79
Transformer <sub>big</sub>	192	87.09
BERT <sub>base</sub>	1438	652.27
NAS	626155	284018.9
ELMo	262	118.84

then run our estimates

model	expected energy (kWh)	estimated energy match (kWh)	estimated energy base (kWh)	expected CO2e (kg)	estimated CO2e match (kg)	estimated CO2e base (kg)	ex
Transformer <sub>base</sub>	27	38	74	11.79	16	27	
Transformer <sub>big</sub>	201	267	523	87.09	116	194	
BERT <sub>base</sub>	1507	2000	3920	652.17	865	1450	
NAS	656347	871000	1.71e+06	284018	377000	632000	
ELMo	275	404	793	118.84	175	293	

We can see that we obtain estimates that are, as expected, a little bit higher than those presented, the differences between the match and base setups can be explained by two things: the used Carbon Intensity for the USA in the base values is 370gCO<sub>2</sub> e/kWh instead of the 432gCO<sub>2</sub> e/kWh when trying to match. The dynamic ratio is roughly twice as high when using the base value compared to using the indicated PUE.

To complement the case study on hyperparameter search and costs not only on training one model but of the whole process, let us try and reproduce similar results, which we would be able to study also in terms of the other impacts estimated by our tool.



Models	Hours	Expected energy (kWh)	Estimated energy (kWh)	Expected electricity cost (\$)	Estimated electricity cost (\$)
1	120	41.7	93	5	11
24	2880	983	2230	118	268
4789	239942	82250	186000	9870	22320

We can see that we still obtain values approximately twice as high as the ones presented. This fact can be mostly explained by the difference between using a PUE of 1.58 and a dynamic ratio of 3.1

### 6.2.3 integrating Life cycle to previous analyses

```
estimated impacts: {'gwp': {'embodied': 2800.0, 'direct': 68800.0, 'total': 72000.0, 'unit': 'kg CO2e'},
to put impacts in perspective: {'relative_SNBC': {'value': 36.0, 'unit': 'Emissions of CO2e per kWh'},
Direct energy consumption: 186000.0 kWh, translates to a cost of 22320.000000 $
(22320.0, 186000.0)
```

We can see that the full impacts estimated for performing the whole model search, hyperparameter tuning and training represents the annual impacts of 36 persons if we place ourselves in a scenario where we would respect the "Stratégie Nationale Bas Carbone" for France by 2050. If we place ourselves in the framework of the Planetary boundaries, where if we want to stay sustainable, societies must not overpass the planetary boundaries. The whole process accounts for the maximal annual impacts of 73 persons in terms of Green House Gas emissions and the annual impacts of 25 persons in terms of ressource depletion.

Of course, if computations were to run in a country with a less carbon intensive electricity mix, green warming potential would be lower. Still, the impacts on ressources depletion are very important, and, in this estimation, we do not take into account any (1 GB) memory on the server that runs the experiments.

If we were to add memory, for instance 512 GB of memory, we would obtain the following estimation

```
estimated impacts: {'gwp': {'embodied': 12000.0, 'direct': 121000.0, 'total': 130000.0, 'unit': 'kg CO2e'},
to put impacts in perspective: {'relative_SNBC': {'value': 67.0, 'unit': 'Emissions of CO2e per kWh'},
Direct energy consumption: 327000.0 kWh, translates to a cost of 39240.000000 $
(39240.0, 327000.0)
```

with expected impacts as high as the maximal anual ones of 140 persons in terms of GWP and 34 persons in terms of Ressources depletion.

As a title of comparison, if we were to make the same estimates but running in France, we would obtain the following (with a carbon intensity of 98gCO<sub>2</sub> e/kWh)

```
estimated impacts: {'gwp': {'embodied': 12000.0, 'direct': 32100.0, 'total': 45000.0,
to put impacts in perspective: {'relative_SNBC': {'value': 22.0, 'unit': 'Emissions of
Direct energy consumption: 327000.0 kWh, translates to a cost of 39240.000000 $
```

It would still represent the maximal annual emissions of 45 persons in terms of GWP and the maximal impacts of 33 persons in terms of ressources depletion

## 7 comparing manufacture impacts with Dell LCAs

First define a helper functions to print pie charts with the repartition of impacts by components :

### 7.1 Dell R740

```
GWP: {'manufacture': 2400.0, 'use': 1170.0, 'unit': 'kgCO2eq'}
PE: {'manufacture': 31000.0, 'use': 39700.0, 'unit': 'MJ'}
ADP: {'manufacture': 0.19, 'use': 0.000198, 'unit': 'kgSbeq'}
RAM gwp: 540.0 kgCO2eq
SSD gwp: 24.0 kgCO2eq
Other component GWP impacts: 302 kgCO2eq
```

```
impact_dict
```

```
{'MAINBOARD': 111.69999999999999,
 'OTHER': 302.0,
 'RAM': 540.0,
 'SSD': 24.0,
 'SSD_3.84GB': 1440.0}
```

```
wdf sdf
```

```
#l = [(n, v) for (n,v) in impact_dict]
tabulate.tabulate(impact_dict.items(), tablefmt='orgtbl')
```

MAINBOARD	111.7	
RAM	540	
SSD	24	
OTHER	302	
SSD_3.84GB	1440	

Component	ACV DELL - GWP (kgCO2eq)	MLCA - GWP (kgCO2eq)	Boavizta - GWP (kgCO2eq)
CPU	47	45.6	
RAM	533	540	
SSD	64	24	
OTHER	266	368	
TOTAL	910	970	

We can see that we obtain more or less the same results even after the modifications to the way CPU impacts are computed or some bugfixes.

[width=.9]results/./results/17-05-2303 – 07<sub>R</sub>740

## 7.2 Dell R6515, R7515, R7525

```
total_energy = (213.92 + 595.11 + 347.16 + 207.39)*4
CI_europe = (3450 - 1343) / total_energy
CI_US = (4280 - 1343) / total_energy
print(total_energy, CI_europe, CI_US)
```

```
print(0.152160000000000002 * 4 * 365 * 24)
```

```
5454.32 0.386299300371082 0.5384722568532833
5331.68640000000005
```

```
with open("boaviztapi/data/devices/server/R6515.json", 'r') as m:
    R6515 = json.load(m)
```

```
R6515["usage"]["workload"] = {
    "100": {
        "time": 2.4/24,
        "power": 1.0
    },
    "50": {
```

```

        "time": 8.4/24,
        "power": 184.1/244.2
    },
    "10": {
        "time": 7.2/24,
        "power": 132.1/244.2
    },
    "idle": {
        "time": 6/24,
        "power": 94.7/244.2
    }
}

```

```

print("Europe Scenario")
R6515["usage"]["usage_location"] = "EEE"
R6515["usage"]["gwp_factor"] = 0.386299300371082
out = run_experiment_server(R6515, "R6515_Europe", directory='../results')
print_impacts_server(out)

```

```

print("US Scenario")
R6515["usage"]["usage_location"] = "USA"
R6515["usage"]["gwp_factor"] = 0.5384722568532833
out = run_experiment_server(R6515, "R6515_USA", directory='../results')
print_impacts_server(out)

```

```

R6525 = R6515
R6525["configuration"]['cpu']['units'] = 2
R6525["configuration"]['ram'][0]['units'] = 16

```

```

print("R6525 US Scenario")
out = run_experiment_server(R6525, "R6525_USA", directory='../results')
print_impacts_server(out)

```

```

Europe Scenario
GWP: {'manufacture': 1200.0, 'use': 2060.0, 'unit': 'kgCO2eq'}
PE: {'manufacture': 15000.0, 'use': 68600.0, 'unit': 'MJ'}
ADP: {'manufacture': 0.13, 'use': 0.000342, 'unit': 'kgSbeq'}
US Scenario

```

```

GWP: {'manufacture': 1200.0, 'use': 2870.0, 'unit': 'kgCO2eq'}
PE: {'manufacture': 15000.0, 'use': 60600.0, 'unit': 'MJ'}
ADP: {'manufacture': 0.13, 'use': 0.000526, 'unit': 'kgSbeq'}
R6525 US Scenario
GWP: {'manufacture': 1600.0, 'use': 2870.0, 'unit': 'kgCO2eq'}
PE: {'manufacture': 21000.0, 'use': 60600.0, 'unit': 'MJ'}
ADP: {'manufacture': 0.17, 'use': 0.000526, 'unit': 'kgSbeq'}

```

In the case of Dell R6515, the manufacturing has a contribution of 1,343 kg CO<sub>2</sub>e, approximately 39% to the total of the life cycle impact in the light-medium use scenario

For the manufacture of the R6515, we obtain an estimate of 1200 kgCO<sub>2</sub> e when the expected results stand at 1343 kgCO<sub>2</sub> e.

For the R6525, we obtain an estimate of 1600 kgCO<sub>2</sub> e when the expected result stands at 1709 kgCO<sub>2</sub> e.

```

import matplotlib.pyplot as plt
labels = 'Frogs', 'Hogs', 'Dogs', 'Logs'
sizes = [15, 30, 45, 10]

```

```

fig, ax = plt.subplots()
ax.pie(sizes, labels=labels)

```

```

([<matplotlib.patches.Wedgeobjectat0x7fec40692d90>, <matplotlib.
patches.Wedgeobjectat0x7fec406faeb0>, <matplotlib.patches.Wedgeobjectat0x7fec406fa430>
, <matplotlib.patches.Wedgeobjectat0x7fec406fa760>], [Text(0.9801071672559598,
0.4993895680663527, 'Frogs'), Text(-0.33991877217145816, 1.046162142464278,
'Hogs'), Text(-0.49938947630209474, -0.9801072140121813, 'Dogs'), Text(1.
0461621822461364, -0.3399186497354948, 'Logs')])

```

## 8 replicating the Bloom estimates from [?]

### 8.1 Gathering information about the setup

To replicate their experiments, we first need to gather some information on the time duration and hardware setup for the training phase.

We can see in the paper that the training phase lasted for 118 days, 5 hours and 41 mins for a total of 1,082,990 GPU hours. (table 1)

in section 4.1, we can read that training used on average 48 computing nodes with 8 GPUs each.

Combining the real time and these information about the setup, we obtain an estimate of the number of GPU hours of 1,089,670.4 hours this gives us a pretty close figure to the real GPU time.

It is written in the paper that training took place on the Jean Zay supercomputer, using HPE's Apollo 6500 Gen10 Plus. We can read on their website that it uses AMD EPYC 7000 Series CPUs. Combining this information with informations about the Jean Zay supercomputer on IDRIS's website, we can see that only the **gpu<sub>p5</sub>** partition uses such CPUs. We can conclude that for each of the 48 used nodes, the server configuration is :

- 2 CPUs : AMD Milan EPYC 7543
- 512 Go of Memory
- 8 NVIDIA A100 SXM4 80Go

## 8.2 comparing the server footprint with the PCF sheet.

In section 4.1, it is stated that they use values provided in the HPE ProLiant DL345 Gen10 Plus PCF, the closest server with information provided. In this PCF sheet, we can read that servers are of type rack and that the estimated Carbon Footprint is of 2503.2 kg CO<sub>2</sub> e. If we try our tool with the server configuration used for training, we obtain :

```
GWP: {'manufacture': 2300.0, 'use': 1170.0, 'unit': 'kgCO2eq'}
PE: {'manufacture': 29000.0, 'use': 39700.0, 'unit': 'MJ'}
ADP: {'manufacture': 0.17, 'use': 0.000198, 'unit': 'kgSbeq'}
RAM impact GWP: {'value': 1800.0, 'unit': 'kgCO2eq'}
```

we can see manufacture impacts of 2300 kg CO<sub>2</sub> e. This impact is close to the 2500 kgCO<sub>2</sub> e provided on the PCF sheet and is mainly impacted by the quantity of memory used, as it accounts for 1800 kg CO<sub>2</sub> e.

## 8.3 comparing the GPU footprint with the chosen value

In section 4.1, it is stated that a value of 150 kg CO<sub>2</sub> e is chosen. Taking a look at the source, there is no real justification given for that value. Given that in [?] a small GPUs manufacture is estimated at emitting around 30 kg CO<sub>2</sub> e, we can hypothesize that GPU manufacture impacts would be in the order of 50 to 150 kg CO<sub>2</sub> e.

For the specific model used, the "NVIDIA A100 SMX4 80GB", we can see a manufacture impact of 310 kgCO<sub>2</sub> e. this impact is mainly influenced

by the quantity of memory on the GPU with 290 kg CO<sub>2</sub> e. These are preliminary results since the base value for gpu impacts is not properly set yet.

## 8.4 Estimating the total impacts

with all of the previous information, we can run the estimation

we can see in the results (full result in results/datetime bloom.json) that we obtain close figures to those in the paper. with embodied impacts of 7T CO<sub>2</sub> e for the servers and 7.6T for the GPUs to compare with the 7.6T for the servers and 3.6 T for the GPUs in the paper. Most of the difference is due to estimated impacts of 300 kgCO<sub>2</sub> e for one GPU while it was estimated to 125 kgCO<sub>2</sub> e in the paper.

For the dynamic consumption, we obtain an estimate of 26.8T CO<sub>2</sub> e, mainly due to the GPUs (accountable for 25T, the only difference with the figure obtained in the paper being the slightly off conversion from real time to GPU hours) while the memory, not accounted for in the paper brings another 1.5T CO<sub>2</sub> e.

The only thing that differs greatly is the value for the idle consumption. (not so surprising since figures differ quite a lot).

## 9 Conclusions

After these experiments trying to evaluate the validity of our tool, we can draw some conclusions, firstly about the challenges of replicating results and then about the validity of our tool.

### 9.1 about the replication of results

Overall, reproducing results from different papers proved way harder than expected. Indeed, Unless a real effort is made by authors to allow replication of their results, it is most of the time really difficult to find enough information to run estimates and reproduce their results. This is also particularly true for results produced using a measurement tool, indeed, If the hardware on which those results were produced isn't explicited, it is impossible to reproduce the experiments and check the quality of the results presented. We were only able to conduct experiments for all of these papers because we were able to contact the authors and they were able to give us some insight about the hardware configuration of their experiments.

Even when we had enough information to run our estimates precisely enough to hopefully match the expected results, we faced multiple times important errors and inconsistencies in the data presented in different tables. This was for example the case with the results presented in [?] and in [?]. This was also the case to a lesser extent in [?] where a notable effort for reproducibility was realised by the authors but there were still some problems and hypotheses that needed to be made in order to reproduce the results. After pointing out the problems with the data presented in [?], the authors conducted new experiments to resolve the problems with their data and we were able to reproduce these new results.

## 9.2 about the validity of the tool

Running new experiments often required us to gather some information about a CPU not present in our database. This was not needed for GPUs. It seems like there is much more diversity in CPUs used than in GPU used. However, it was relatively easy to find all the information we needed when encountering a new CPU and when running estimations about GPU intensive tasks such as training NLP models, the CPU usage is often set close to 0. Moreover, CPU manufacturing does not play a huge part in the manufacturing impacts of a server in terms of GWP, it does however play an important part of the impacts in terms of mineral resource usage (ADP)

We were unfortunately not able to find experiments to demonstrate the validity of other indicators than the Global Warming Potential.

Still, we can see that overall, we were able to reproduce results for the dynamic consumption and for the embodied impacts. These experiments also demonstrate the usability of our tool in diverse scenarios.

## 10 TODO :noexport:

faire des jolis camemberts ?

faire un truc sur la variation du dynamic ratio ? variation de la durée de vie ? baisser le dynamic ratio mais baisser la durée de vie du matériel. (scenario à la patterson)

/!\ résultats de variabilité qui exploitent ADP dans strubell, à modifier quand changement dynamic ratio sera fait.

On en est à Dinarelli dans la relecture.