

# LINMA2710 - Scientific Computing

## Power Consumption

*P.-A. Absil and B. Legat*

☐ Full Width Mode    ☐ Present Mode

### **Table of Contents**

**Energy consumption**

**Carbon intensity**

**Power consumption of computing**

**Reducing power consumption**

# Energy consumption

## Primary energy consumption

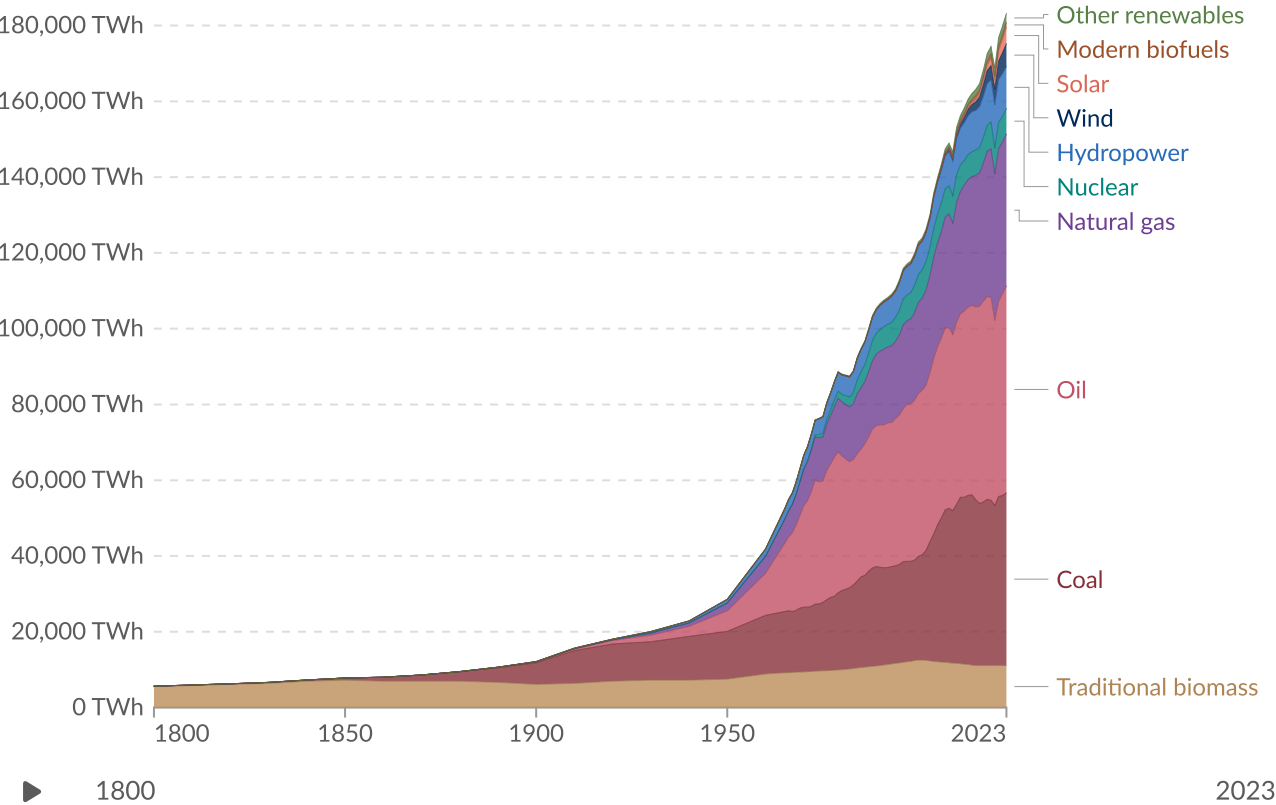
Our World  
in Data

### Global primary energy consumption by source

Primary energy is based on the substitution method and measured in terawatt-hours.

Table Chart

Settings



Data source: Energy Institute - Statistical Review of World Energy (2024); Smil (2017) - [Learn more about this data](#)

Note: In the absence of more recent data, traditional biomass is assumed constant since 2015.

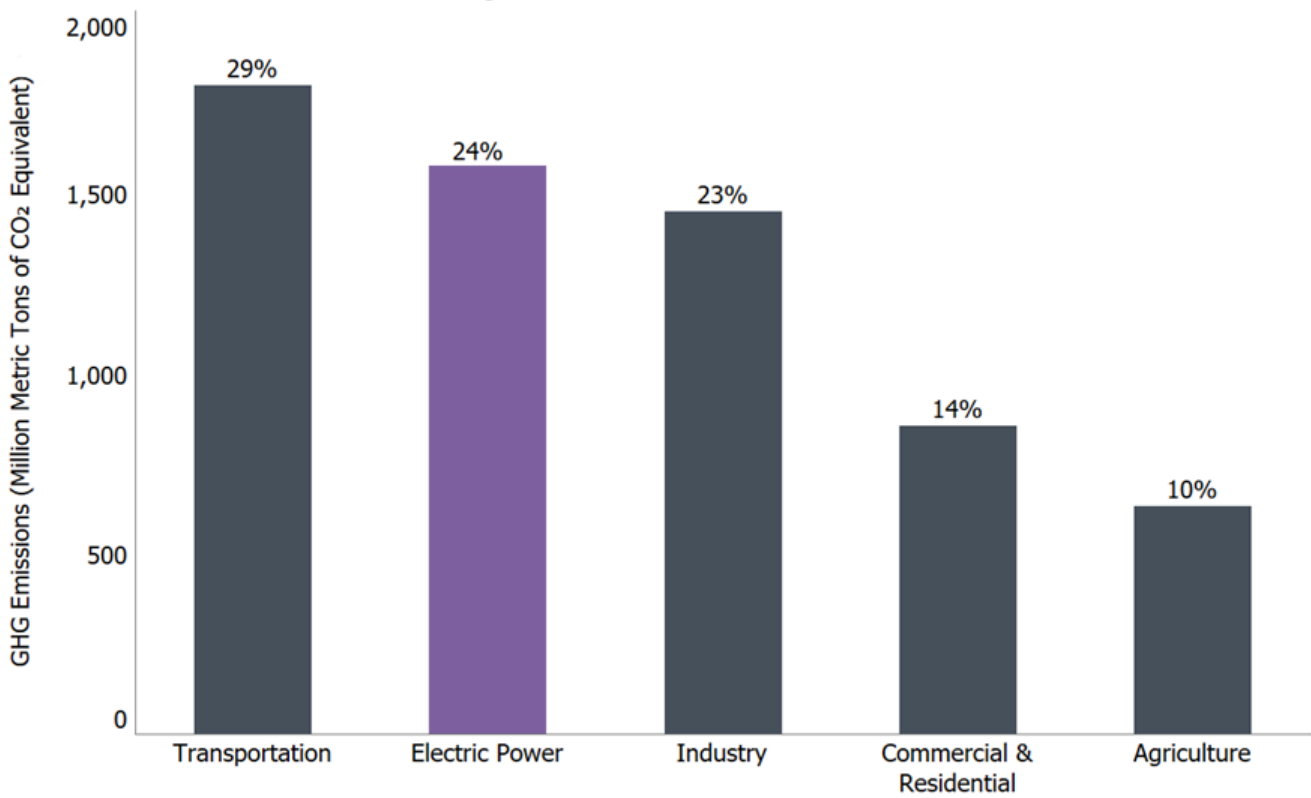
OurWorldinData.org/energy | CC BY

Download Share Explore the data →

# The emisson from the power sector not negligible

Total U.S. Greenhouse Gas Emissions by Economic Sector in 2022

Greenhouse Gas Emissions by Economic Sector



# Electricity generation per country

## Electricity generation, 2024

Total electricity generated in each country or region, measured in terawatt-hours.



Table

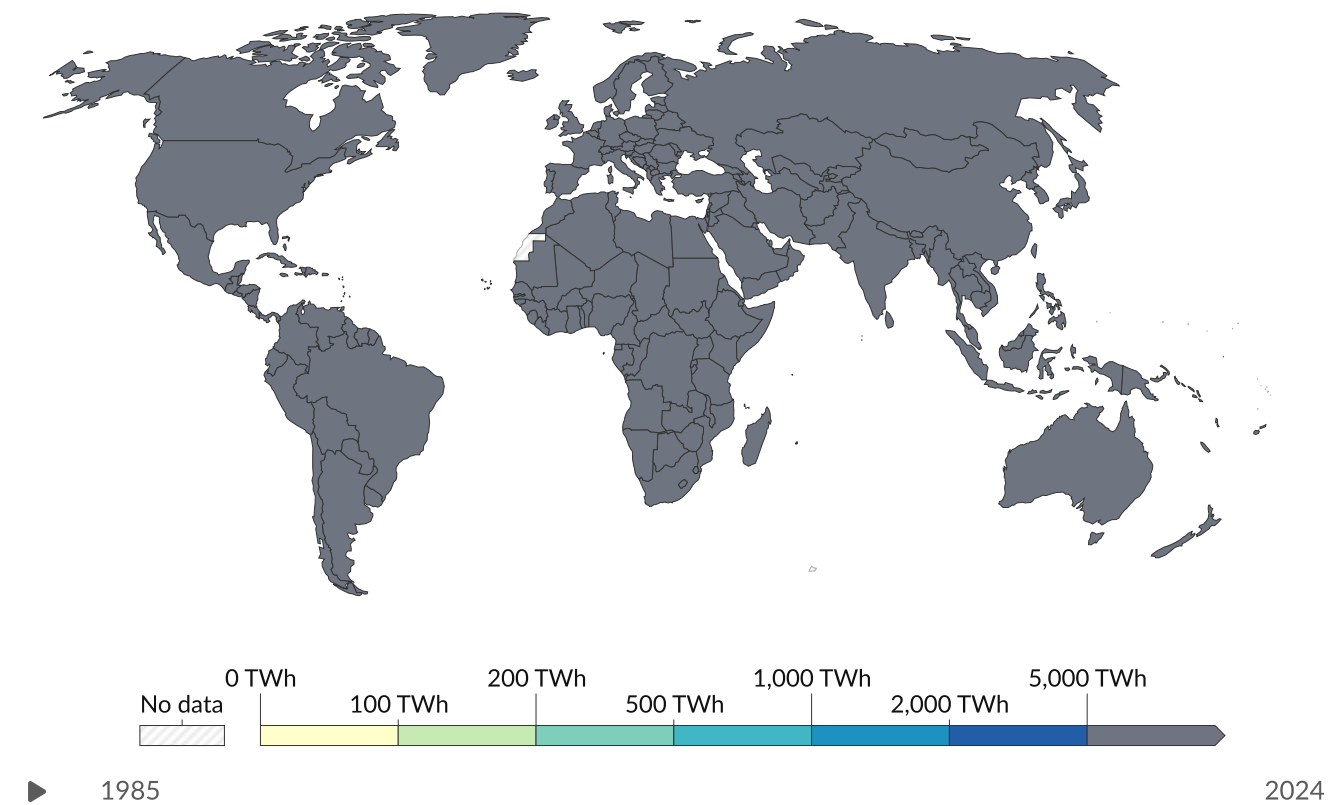
Map

Chart

Zoom to...

2D3D

Select



Data source: Ember (2025); Energy Institute - Statistical Review of World Energy (2024) – [Learn more about this data](#)  
OurWorldinData.org/energy | CC BY

[Explore the data](#) →

# Electricity generation per capita

---



# Carbon intensity

---

## Intensity per source

---

Energy Source	Carbon Intensity (kg/MWh)
Coal	995
Petroleum	816
Natural Gas	743
Solar	48
Geothermal	38
Nuclear	29
Hydroelectricity	26
Wind	26

Source

# Share of production by sources

---

# Share production by group

---



# Carbon Intensity per country



# Power consumption of computing

---

## RAM consumption

---

- 2-5 W per slot when the computer is on
- Drops to around 1 W when in sleep mode
- This is the reason your computer still consumes while sleeping.
- To turn the computer off, you need to lose all data in the RAM.
- Hibernation consists of copying the RAM to disk so that the RAM can be completely shut-off and copy it back when you turn it on.

## CPU consumption : Thermal Design Power

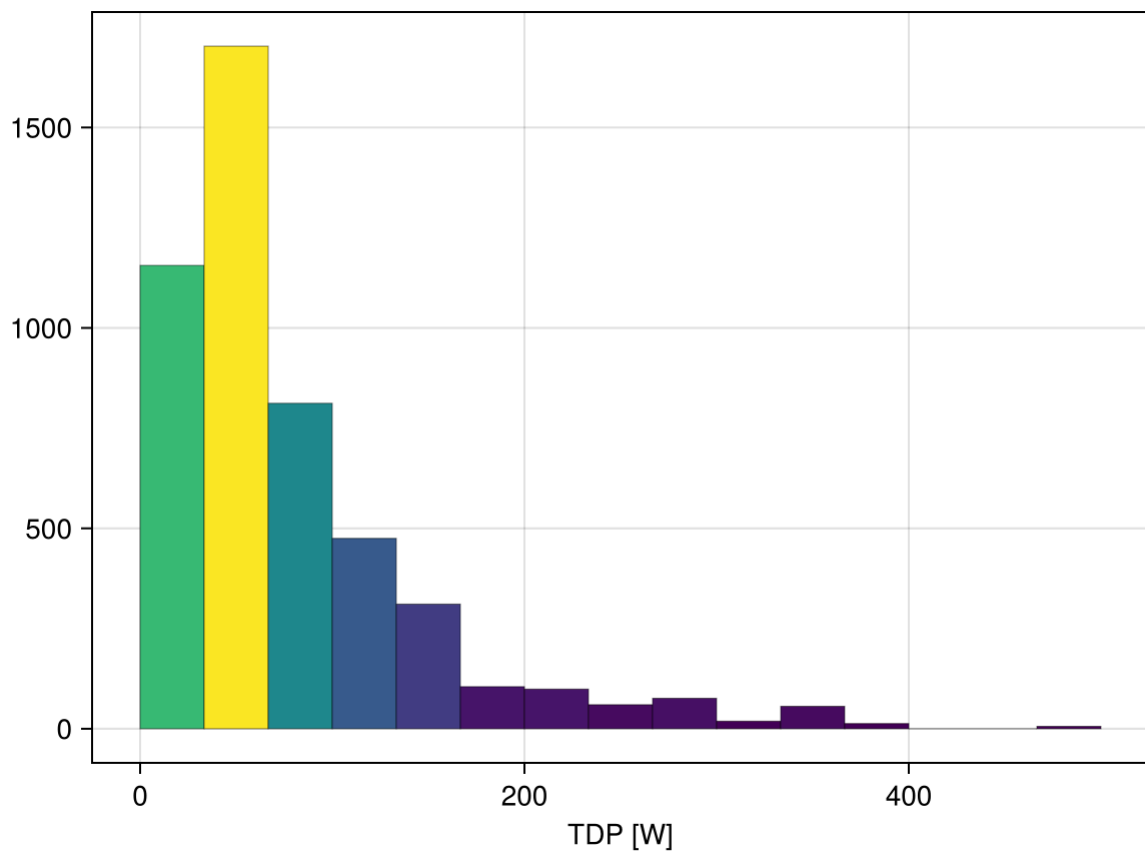
---

Thermal Design Power (TDP) : Maximum amount of heat a CPU is designed to generate (important for cooling).

Power of CPU is approx. proportional to its utilization/load until below 10%

$$\text{power} = \text{TDP} \times \max(0.1, \text{load})$$

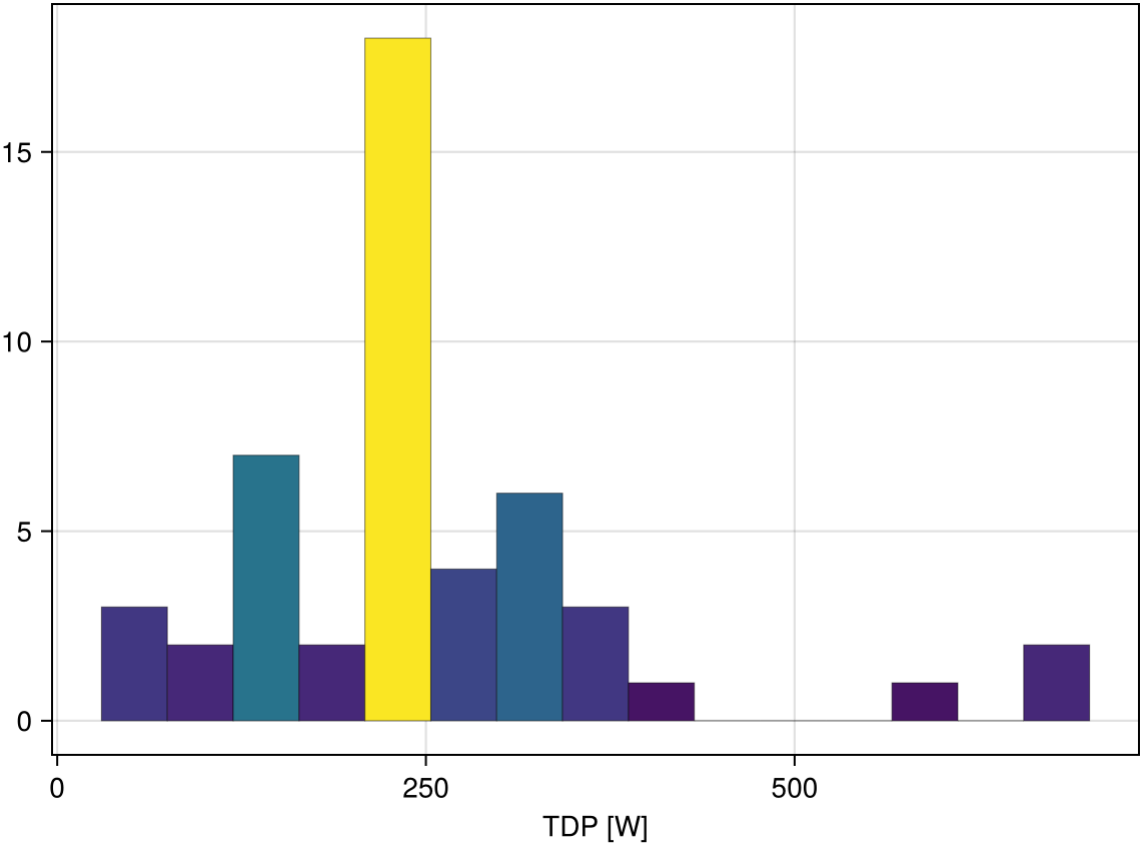
Source



	Name	TDP
<b>1</b>	"AMD EPYC 9755"	500.0
<b>2</b>	"AMD EPYC 9965"	500.0
<b>3</b>	"Intel Xeon 6960P"	500.0
<b>4</b>	"Intel Xeon 6972P"	500.0
<b>5</b>	"Intel Xeon 6979P"	500.0
<b>6</b>	"Intel Xeon 6980P"	500.0
<b>7</b>	"AMD EPYC 9475F"	400.0
<b>8</b>	"AMD EPYC 9565"	400.0
<b>9</b>	"AMD EPYC 9575F"	400.0
<b>10</b>	"AMD EPYC 9655"	400.0
⋮ more		
<b>4891</b>	"Intel Quark Microcontroller D1000"	0.0

# Power consumption of GPUs

.....



	name	type	tdp_watts	TFLOPS32	TFLOPS16	GFLOPS32/W	GFLC
1	"Tesla H100-PCIE-80GB"	"gpu"	700	missing	missing	missing	miss
2	"Tesla H200-SXM"	"gpu"	700	missing	missing	missing	miss
3	"Tesla H200-PCIE"	"gpu"	600	missing	missing	missing	miss
4	"A100 SXM4 80 GB"	"gpu"	400	missing	312.0	missing	miss
5	"RTX 3080 TI"	"gpu"	350	34.1	34.1	NaN	NaN
6	"RTX 3090"	"gpu"	350	35.58	35.58	NaN	NaN
7	"Tesla H100-PCIE"	"gpu"	350	missing	missing	missing	miss
8	"RTX 3080"	"gpu"	320	29.77	29.77	NaN	NaN
9	"Tesla K80"	"gpu"	300	4.113	missing	13.71	miss
10	"Tesla V100-PCIE-16GB"	"gpu"	300	14.13	28.26	4.71	94.2
	: more						
49	"AGX Xavier"	"gpu"	30	16.0	32.0	533.33	1066

# Power consumption of a cluster

xAI Colossus will be the largest cluster being built (not counting mega-clusters made of several ones connected by optic fibers). Made of 100k H100 GPUs and 100k H200 GPUs so a total of 150 MW. This will need a whole new Gas turbine just to power it:

	Source	Power	xAI need
1	"Nuclear plant"	1000.0 MW	0.15
2	"Gas turbine"	150.0 MW	1.0
3	"Wind turbine"	3.0 MW	50.0
4	"1000 PV panels"	0.32 MW	468.8

Source

# Reducing power consumption

---

## Break down

---

The power consumption of a chip is the sum of two sources:

- *Static power* : primarily due to leakage currents, which become more important as the transistor size decreases.
- *Dynamic power* : Switching power given by  $CV^2Af$  where
  - $C$  : Capacitance being switched
  - $V$  : Voltage
  - $A$  : *Activity factor*, i.e., number of switches of transistors per clock cycle.
  - $f$  : Clock frequency

The higher the voltage is, the higher are the leakage currents hence the power consumption but the voltage cannot be lowered without lowering the frequency hence the two are often done together → DVFS. Example of application in [YCC23].

## Dynamic voltage and frequency scaling (DVFS)

---

The higher the voltage is, the higher are the leakage currents hence the power consumption but the voltage cannot be lowered without lowering the frequency hence the two are often done together → DVFS. Example of application in [YCC23].

► **If the clock frequency is decreased, does the time performance always get worse ?**

[YCC23] J. You, J.-W. Chung and M. Chowdhury. Zeus: Understanding and Optimizing GPU Energy Consumption of DNN Training. In: 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23) (2023); pp. 119–139.



# Gating

---

When a core is idle, it may first be

- *clock-gated* : part of the circuit stops switching. Corresponds to states C1-C3 in intel CPUs.

If it continues being idle, it may then

- Reduces the voltage as the clock frequency is now zero. This will reduce the leakage currents. Corresponds to state C4 in intel CPUs.

If it continues being idle, it may even be

- *power-gated* : turn off circuit blocks. Which eliminates the leakage currents. It has a fixed power cost but is worth it if it is unused for a long time [WYCC11]. Corresponds to states C6-C10 in intel CPUs.

You can inspect the states of your cores on your laptop with

- Intel Power Gadget (Windows/macOS)
- powertop (Linux)
- HWInfo or Throttlestop (Windows)

[WYCC11] P.-H. Wang, C.-L. Yang, Y.-M. Chen and Y.-J. Cheng. Power Gating Strategies on GPUs. ACM Trans. Archit. Code Optim. **8**, 13:1–13:25 (2011).

# Reducing the power consumption of your code

---

- **DVFS and gating are automatically handled. So what actions can be taken in the design of your code or GPU kernels to allow these optimizations ?**



```
Activating project at '~/work/LINMA2710/LINMA2710/Lectures'
```



**biblio =**

►CitationBibliography("/home/runner/work/LINMA2710/LINMA2710/Lectures/references.bib", Alp

① Loading bibliography from `/home/runner/work/LINMA2710/LINMA2710/Lectures/references.bib`...

① Loading completed.