

# Measure energy use of Apertus

Understanding the energy behaviour of Apertus

Agustín Herrerapicazo Luis Barros Stefan Aeschbacher agustin.it@proton.me luisantoniio1998@gmail.com stefan@aeschbacher.ch

#### Goals

#### The goal of the project was to

- measure the energy consumption of inference on Apertus
- understand the influence of prompt and response on the energy consumption
- Bonus: compare to llama

## Measuring Infrastructure

Our measuring infrastructure consisted of

- Mac Studio with remote access.
- powermetrics tool
- Begasoft BrandBot





### **Experiments**

**Experiment 1:** Prompt/Response length

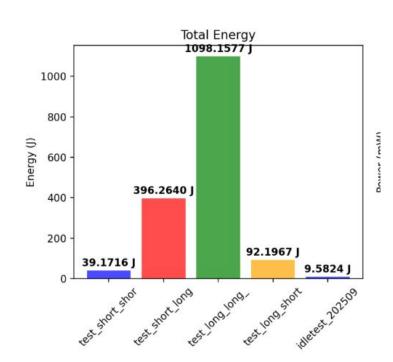
We tested different prompt and response lengths.

Ex: Short-Long: Explain photosynthesis in detail.

**Experiment 2:** Different subjects

**Experiment 4:** Different languages

**Experiment 3:** Compare with Llama



## Learnings

#### We failed successfully on:

- getting a stable setup X
- excluding overhead 🔀
- proper energy calculations X
- getting 2nd system to run X
- do physical measurements X
- change measurement method X
- automate everything X
- doing proper statistics X

#### We succeeded successfully on:

- measure idle consumption
- measure different prompts
- 🕨 🛮 Integrated a new team member 🔽



having fun!

#### Results



- Measuring is hard!
- Prompt size:
  - o long answers dominate the energy use
  - long prompts have some influence
- Different types of prompts: Basic knowledge, Chemistry, Geography, Math
  - o does not seem to influence the energy consumption (other than the length of the response)
- Different languages:
  - o apart from the length, the language seems to have some influence (portugese: +20%)
- Different models: inconclusive but similar

https://decs.google.com/decument/d/1.u.hNlvOV.v.TVgd.tafgQUI\_2EuEVNIZenNlcO\N/c2IIghE4/cdi+2tah\_++0