



## Main Goal

Evaluating the impact of quantization on operational-level LLM metrics.

### RQ1

Quanti:  
benchmarking  
system for LLMs

### Metrics

GPU[%], Memory  
[GB], Energy [kWh]

### RQ2

Accuracy-efficiency  
tradeoffs

### Metrics

MMLU accuracy  
[%], Energy [kWh],  
Latency [ms]

### RQ3

Architecture-  
specific tradeoffs

### Metrics

MMLU accuracy  
[%], Energy [kWh]

### RQ4

CO2-aware  
scheduling impact

### Metrics

Emissions [gCO<sub>2</sub>],  
Energy [kWh]