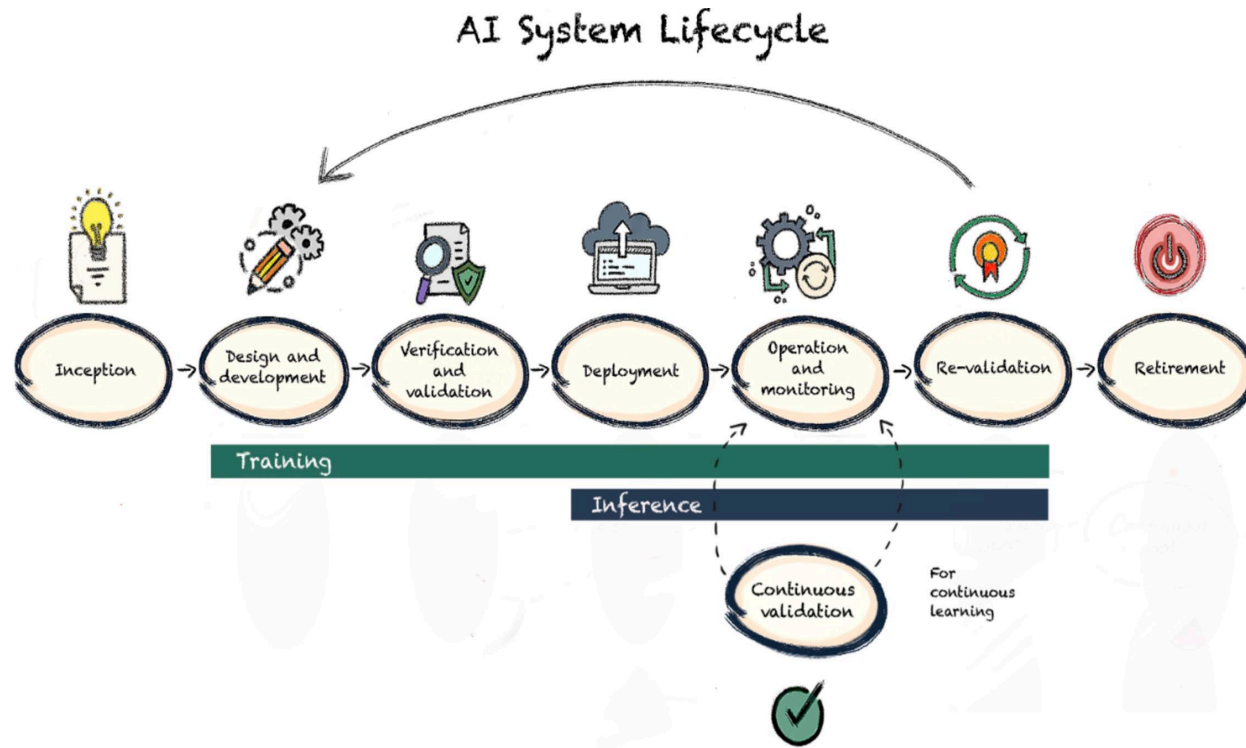


Smarter, Smaller, Stronger:

Resource-Efficient Generative AI &

the Future of Digital Transformation

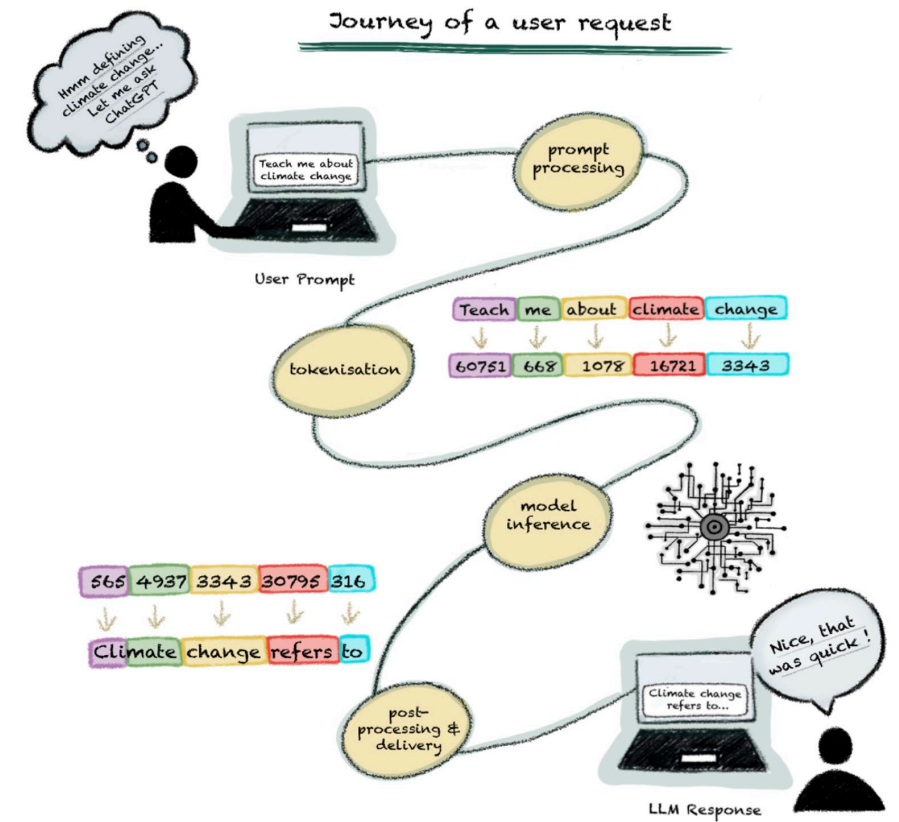


Authored by: Leona Verdadero^{1*}, Ivana Drobnjak^{2*}, Hristijan Bosilkovski^{2*}, Zekun Wu²³, Emma Fischer, and María Pérez-Ortiz².

¹UNESCO

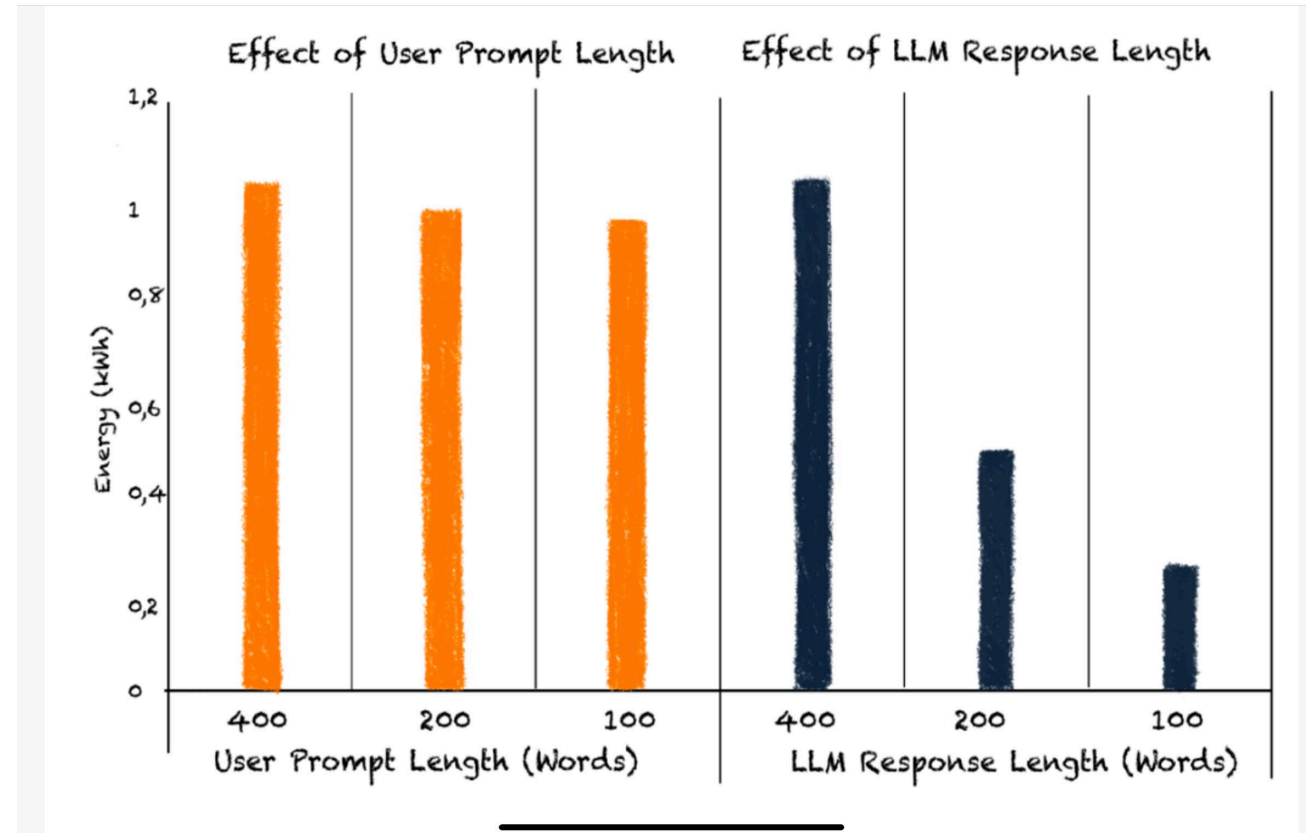
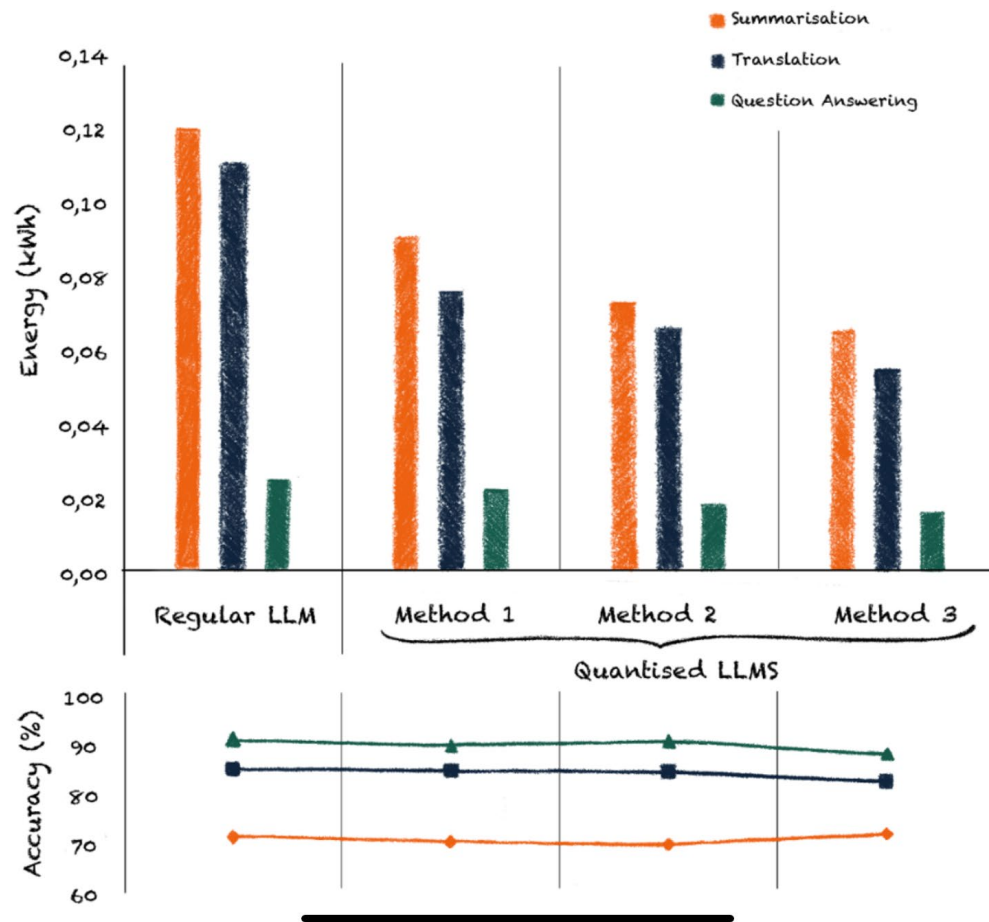
²University College London

³Holistic AI



General models: 75% energy reduction

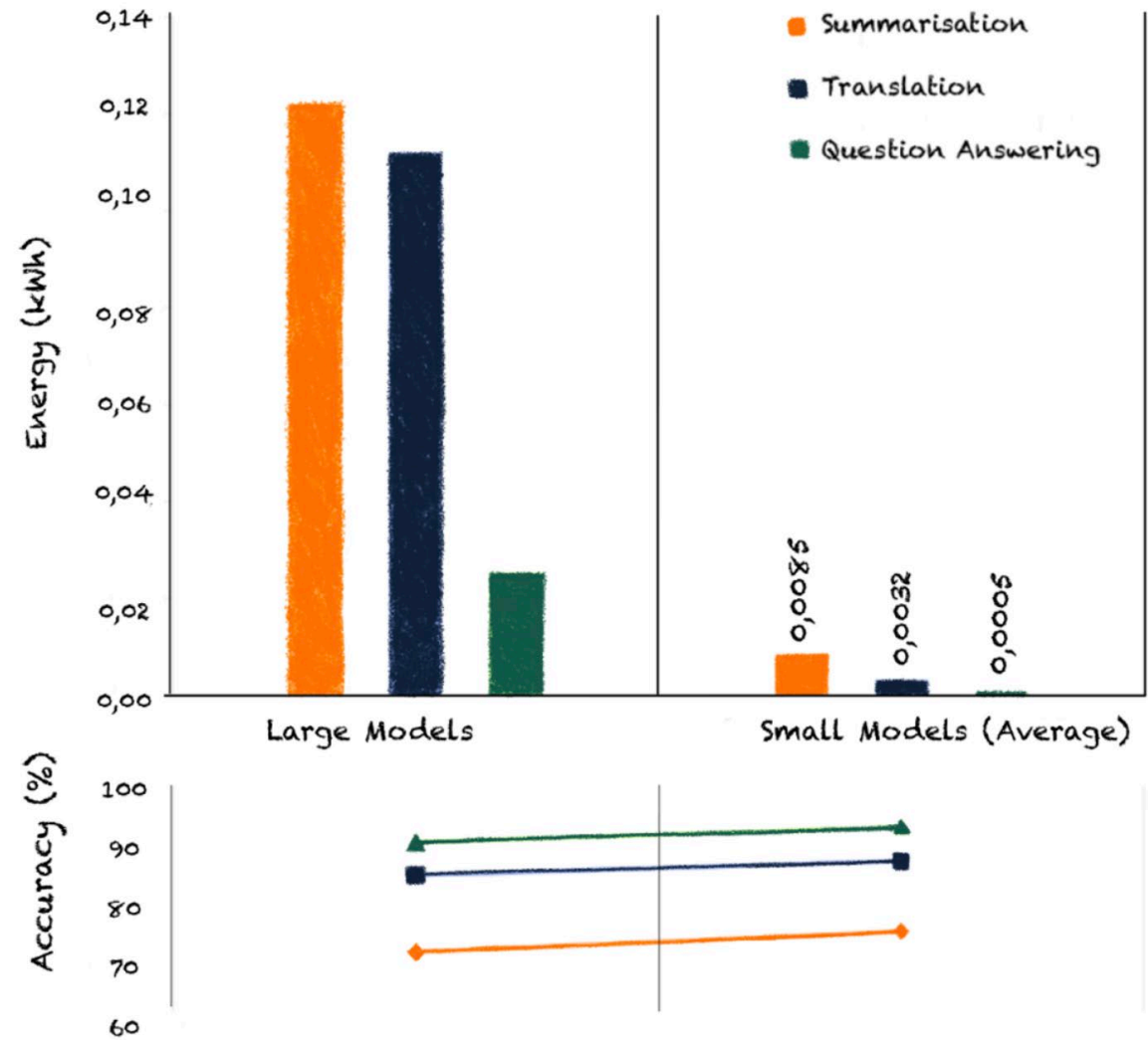
Effects of Quantisation



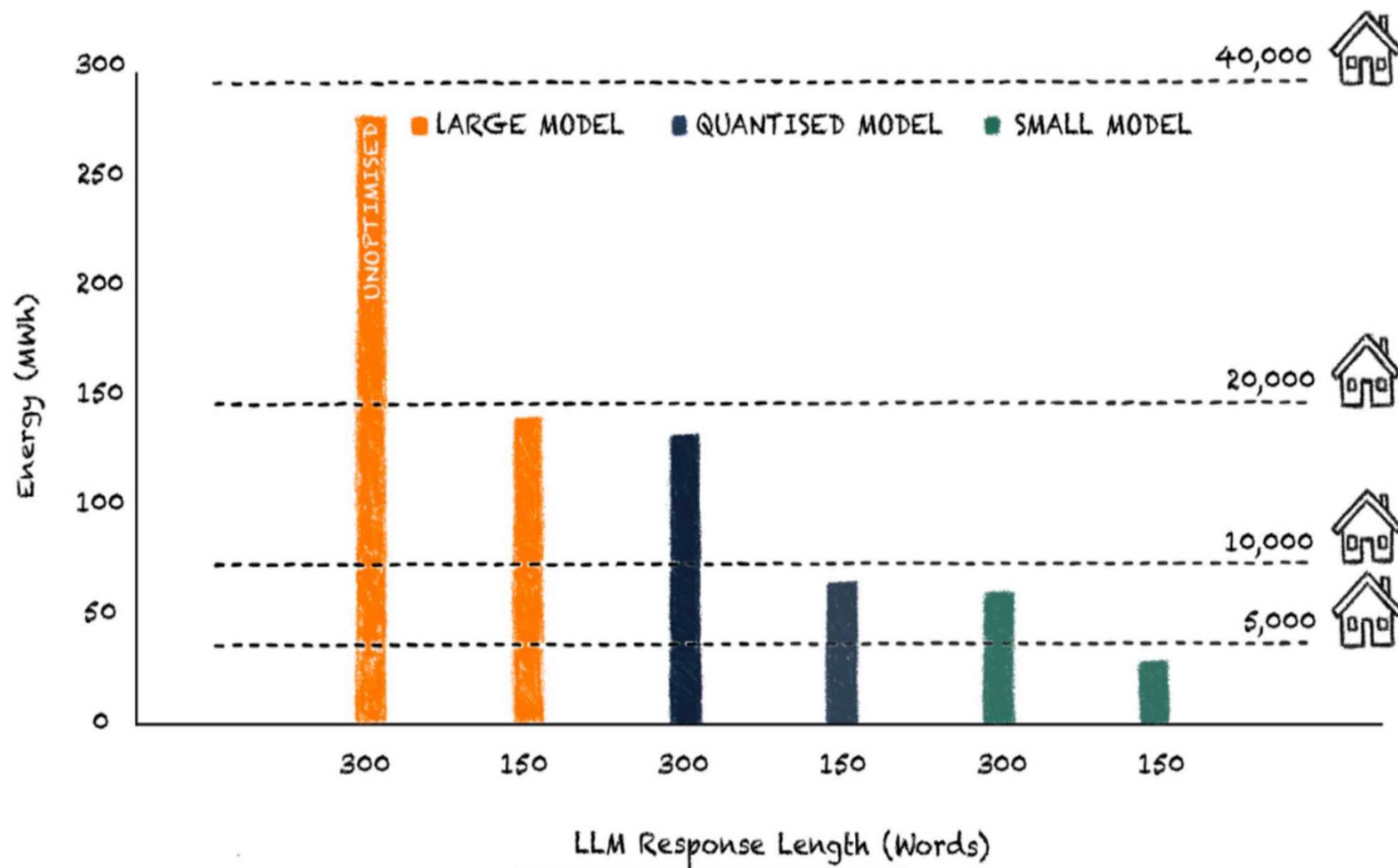
For specialised tasks: Small models

over 90% energy reduction

Effects of Model Size



DAILY IMPACT OF OPTIMISATION ON ENERGY SAVING



The future

- Optimising large models (e.g. quantisation etc.)
- Using small models where appropriate
- More efficient large model architectures (work in progress):
 - Mixture of Experts
 - Multi Agent
 - Sparse and conditional computation
 - Retrieval-augmented generation
 - Neurosymbolic and brain inspired architectures
- More work needs to be done – designing with efficiency in mind



Smarter, Smaller, Stronger:

Resource-Efficient Generative AI &

the Future of Digital Transformation

Smarter, Smaller, Stronger: Resource-Efficient Generative AI & the Future of Digital Transformation

