0. Abstract: what is being done, why a
what was achieved

— ̄ — ̄ — ̄ — ̄

1. Intro

bcg a motivation → prop

our contributions
scope
paper overview

2. Theoritical bcg
→ LLM & LLM context, sampler, etc
→ M1 SOC, GPU, NPU, Shared
address space
→ zero copy reads (servers data streams,)
→ quantization
→ data leakage

3 Related work
→ ollama & llama.file
→ relevant papers from proposal

→ ghotz rev eng

→ how coreml model uses NPU

4 methodology / design

↳ vec dump & mmap

↳ NPU hack
↳ RAG pipeline

5 Implementation & architecture

( UI, static lib, etc)
↳ context pool & model reuse

6 Experimentation
bertscore, models, quant, etc
Results & analysis

7 Conclusion & future work

↳ limitations

↳ NPU backend for llama