

I've Got 99 Problems, but LLMs Ain't One!

Cedric Clyburn · Senior Developer Advocate, Red Hat



Why private and local AI?

83% of organizations are moving workloads back to private cloud or on-prem for privacy and cost ([Barclays CIO Survey 2024](#)).

I've Got 99 Problems, but LLMs Ain't One!

Cedric Clyburn · Senior Developer Advocate, Red Hat



Why private and local AI?

83% of organizations are moving workloads back to private cloud or on-prem for privacy and cost ([Barclays CIO Survey 2024](#)).

Control & Privacy

Keep data on your infra, control updates and versions

Cost Predictability

Avoid unpredictable egress/inference bills

Customization

Tune models, add guardrails, ship your cadence

Portability

Run on Linux, Kubernetes, or bare metal

About me

Cedric Clyburn

Senior Developer

Advocate, Red Hat



Organizer, [KCD New York](#)

Open source +
community builder ·
[@cedricclyburn](#)

Technical Educator

- [RAG vs Fine-Tuning](#)
- [Ollama](#)
- [vLLM](#)

Let's get started!

OpenSource Model Explosion

Foundation models' new in OSS

Llama 3.x, Qwen, DeepSeek, Mistral, Gemma

Private & sovereign AI options

Instruct

base

MoE

Modalities

Text, vision, audio, diffusion,
multi-modal

VLM (Llava, Idefics)

Diffusion

Why hosted API ecosystem

Always enough

Tokenizers, safetensors, GGUF, serving
runtimes

How to deploy: local → production

llama.cpp

vLLM

TGI

GGUF

safetensors

Tooling: llama.cpp, vLLM, TGI

Cut GPU cost by ~50%: quantization

Options to use LLMs today

Private and Sovereign AI

Owning your stack from weights → runtime → platform.

Infra

Linux, GPUs, Kubernetes, observability

Serving

llama.cpp (CPU/GGUF) & wrappers,
vLLM (GPU), TGI

Governance

Security, access, model registries,
evaluations

Transformers

llama.cpp

vLLM

What if we cut costs in half?

Quantization shrinks weights and KV cache → fewer/lower VRAM GPUs, faster tokens/sec.
OpenAI-compatible server

What it is

Open-source LLM serving (UC Berkeley).

Approaches

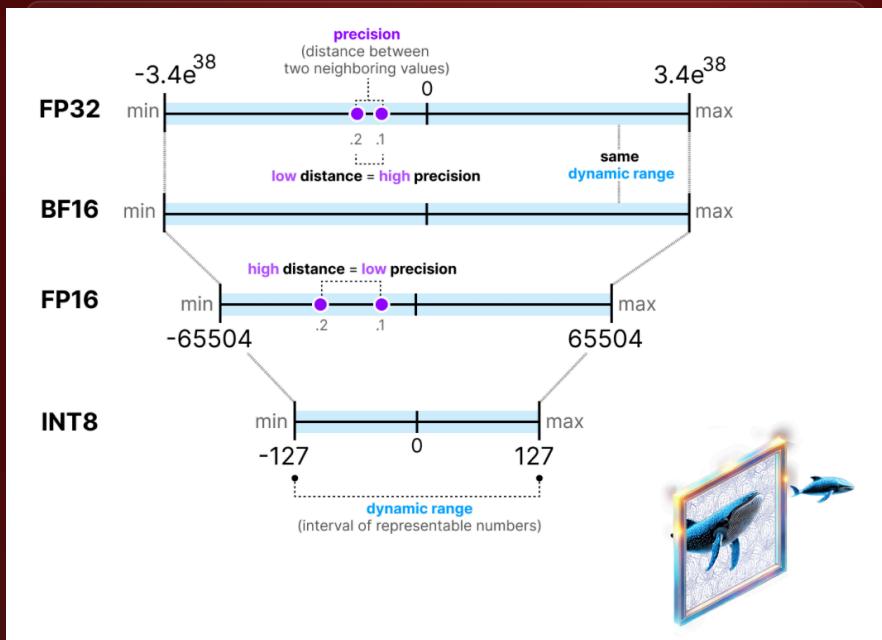
High throughput, low-latency GPU inference with
AVG/GPTQ/GQA-aware INT8/INT4, FP8
Paged Attention and efficient KV cache!

Strengths

- **Typical wins**
Excellent throughput with batching/continuous batching
- **2x VRAM reduction, similar quality** with right
OpenAI compatible server, easy clients
- **calibration**
Work with HF models; supports tensor/quant formats

Trade-offs

- **Trade-offs**
Different GPUs and VRAM; plan for scheduling
- **Slight complexity/batching, memory, layer-wise parallel**
sensitivity matters



Quantization demo

Local → Production

Local

- llama.cpp binary, GGUF weights
- vLLM single-GPU serve for quick iteration
- Experiment in notebooks/Transformers before serving

Production

- Kubernetes + GPU scheduling (time-slicing/MIG)
- Autoscaling, observability, circuit-breakers, canaries

Local → Production

Local Takeaways

Docker container, GGUF weights

- vLLM single-GPU serve for quick iteration
- Experiment in notebooks/Transformers before serving

Production

- Kubernetes + GPU scheduling (time-slicing/MIG)
- Autoscaling, observability, circuit-breakers, canaries

SRE checklist: health probes, request budgets, token limits, logging redaction, evals, red-team prompts, rollout policy.

Thanks!

Resources

Takeaways

- Quantization at scale (Red Hat article)
- LLM Compressor (quantization library)
- vLLM (OpenAI-compatible server)
- Open models are thriving
- Red Hat AI models on Hugging Face
- across sizes and modalities
- ramalama (containers / local LLM tooling)
- llama.cpp (GGUF inference)

Quantization can halve GPU needs with minimal quality loss

Follow

@cedricclyburn · YouTube, X/Twitter, GitHub, LinkedIn



Private + sovereign AI
control, privacy, and
governance

You might still have 99
LinkedIn
problems – LLMs ain't
one



Slides

Ai4