

Building with Open Source Al

A Crash Course: World Summit Al

Cedric Clyburn Senior Developer Advocate Red Hat AI BU @cedricclyburn Legare Kerrison Developer Advocate Red Hat Al BU @legarekerrison





cloud computing and automation

artificial intelligence and machine learning

Jan 2023: The Problem

For 2023 and the start of 2024, closed dramatically outpaced open.



The Power of Open

There has been an explosion of capability from open-source over the last 2 years.



The Power of Open

Open models are deployment targets **today**. And the trend is not slowing down.



Advantages of Open Source Models

Open-source models play an important role in the **Enterprise Al landscape**.



Cost

- Self managed infrastructure
- 1B-405B size match task difficulty to model



Customization

 Improve accuracy and costs with task specific tuning



Control

- Model lifecycle (no changes to the model in place)
- Resources (no rate limits / API downtime)



Security

 Complete data privacy (no 3rd party APIs)



Open source is great driver for Innovation

Open source is about more than developing software. It's how we built Red Hat.

And it's completely revolutionized the AI ecosystem too.

			8 8		5 3	- 8 -	8.3	S.,	S. 1	2.1	- 8 - 8	5 - S	1. A	1 8	- (ð	 8 - 3 1	- 8	. 8	e) - 3	5 (5)	- 21	8	8 - S	÷.	8. 8	- 5	5	* *	1	- 86 - 8	8	3. A	8	e 1	e.	- 21	6									
			s - 3	5.2						e						8 - B	×.,				1					20		2. 2				e) - 1	2							2. 3						2
*	34 - 3		÷				14 - A	-			¥		÷		14	9 - 14	×	2		- (*		14				÷.	4			85 - 33			14				14	s - 14	1.04	$\mathbf{x} = \mathbf{x}$		14			-	
										2						1										8															- 8					
+		 +					· ·									÷ +										+				•								e: 14							× •	2
			÷					4		a (1				. 1		а ц	÷												1			-						a (4)		a - 1					4 - 1 a 1	4
																										20																				
+1			a) (a																				e - 14							+1 3						8						- 14				
		 																																												2
			a. 54					-									÷.			. (s.		14			a. 1	•	4			a								a (4)		a. 4		14				3
																																									- 61			-		
																										•																				
																																									- 42		1.1.4	41		4

Open Source Principles

There are various organizations that define open source, but generally...





Transparency

Open source ensures that code and processes are visible and accessible to everyone, fostering trust and accountability.

Community Collaboration

Development thrives through collective input, enabling diverse contributors to innovate and improve projects together.



No Vendor Lock-in

Users have the freedom to choose, modify, and migrate solutions without being tied to a single provider or proprietary system.











these two
personas get
most of the
attention

Data Scientists "Developing Models"

Developers "Writing Code"

Operationalizing AI is one of the biggest challenges





The reality of enterprise IT environments





The reality of enterprise IT environments







Poorly designed systems lead to failed ML projects

Lack of focus on end-to-end system builds technical debt



Technical debt is a barrier to production



Skills: The new AI stack and key OS projects



The 2024 MAD (Machine learning, Artificial Intelligence & Data) overwhelming Landscape

INFRASTRUCTURE	ANALYTICS	MACHINE LEARNING & ARTIFICIAL INTELLIGENCE	APPLICATIONS - ENTERPRISE					
STORAGE MP DB: DUTALASS/ Composition	BPLATFORMS FLOGER & TO RESIDENT FLOGER & T	DATA SCIENCE MOTESONS DATA SCIENCE PLATA SCIENCE SCIENCE PLATA SCIENCE PLATA SCIENCE PLATA SCIENCE PLATA SCIENCE PLATA SCIENCE PLATA SCIENCE	SALES MARKETING GOODINGSAL					
INSCRIMENTATION New Concentration INSCRIMENTATION INSCRIMENTATION INSCRIMENTATION INSCRIMENTATION <tr< td=""><td>DATA ANALYST PLATFORMS Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft Microsoft alteryx Microsoft Microsoft alteryx Microsoft Microsoft Microsoft Microsoft Microsoft Microsoft Avertyn Microsoft Av</td><td></td><td></td></tr<>	DATA ANALYST PLATFORMS Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft alteryx Microsoft Microsoft alteryx Microsoft Microsoft alteryx Microsoft Microsoft Microsoft Microsoft Microsoft Microsoft Avertyn Microsoft Av							
Image: States A point of the states A poin	Loo AndALYTCS Shine the sector of							
OPEN SOURCE INFRASTRUCTURE								

	DATA PARAMENDESS POINT SCAP WARENESS POINT SCAP W	OLAP CONCENTRATION DIFFACTOR	AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRERAL OPJORT PROMISSION AFRAMEWORKS TOOLS & UNOPERAL MIRE AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRE AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UNOPERAL MIRE AFRAMEWORKS TOOLS & UNOPERAL MIRERAL AFRAMEWORKS TOOLS & UN	ALMORELS Commany Walking Memory and Control and All All All All All All All All All Al
--	--	------------------------------	--	---

		DATA SOURCES	& APIs			DATA & AI CONSULTING	
DATA MARKETPLACES	PRANCIAL & MARKET DATA Biomberg Consistentions D SOM JOILS Plant Biomberg 4	AR/SPACE/SEA Splite PRELIDE: S (Porel WINDWARD Gravestate armsterrates Motion Differe Gravestate Structure Motion Encoded for Management Structure Structure Motion	PEOPLE / ENTITIES	LOCATION INTELLIGENCE POISSONNE © mapbas () surgestion @ surgestion @ dataplar unacast. Piace@ @ent c.k.r.@ A Rador © Parmar & Quebla @ Southernine USP OTVTTLAL GOLIGEO "I/V veraset © @	ESO WINKER SUSTAINALYTICS MSCI TRUVALUE LABS: Replice Gsgbook ISSESGS CLAPITY AI Winterhold Arcadia nervista OWLEG WINKO WINTER	Consummiliarie BCG Delotte IEFE Consultar Cons	KPMG Lighthouse

FIRSTMARK

EARLY STAGE VENTURE CAPITAL

The 2024 MAD (Machine learning, Artificial Intelligence & Data) overwhelming Landscape



teractive version: MAD.firstmarkcap.com Comments? Email MAD2024@firstmarkcap.com

Scope For Today

 Training a Foundation Model from Scratch No, that would be expensive
 Working off existing open source base models
 Integrating your data to the model

- Building applications with Al
- Platforms, serving, and operationalizing AI



Session Slides <u>red.ht/open-source-ai</u>







Line Ideation & Prototyping	Building & Refining	/ /
Find LLMs	Connect to data source	Model serving
Try prompts	Exception handling	Endpoints
Experiment with your data	Limited fine tuning	Monitoring
Benchmarking	Retrieval-Augmented Generation (RAG)	Integrate with apps
	Chaining	
	Evaluate flows	
How do I evaluate models and pick the best one for my use case?		How do I deploy my application with LLMs, scale, etc?
	As a developer, how do I build	







Let's begin with the model!

Determining your use case What problem are we solving? **Types of Models** Instruct vs Base vs Embed...





So, which model should you select?

Well... it depends!

- It depends on the use case that you want to tackle.
- DeepSeek models excel in reasoning tasks and complex problem-solving.
- Granite SLM models perform well in various
 NLP tasks and multimodal applications.
- Mistral and Llama are particularly strong in summarization and sentiment analysis.



Types of Models

Instruct Models

These are generative models (like ChatGPT) that have been fine-tuned to follow **natural language instructions**

Use Cases

- Summarize this Article
- Explain Quantum Physics in simple terms

Embedded Models

These models **convert text into vectors** (high-dimensional numerical representations) that capture semantic meaning.

Use Cases

Text Classification by Semantic Similarity

Example output

▶ $[0.023, -0.448, \dots, 1.205] \leftarrow (a vector representing that sentence's meaning)$

You then **compare vectors** (using cosine similarity, etc.) to see how similar two pieces of text are.

Let's begin with the model!

Model Nomenclature

Understanding the naming conventions

Model Transparency

Is the training data biased & understanding openness

Benchmarks/Evaluations

Understanding "vibes" and model evaluations

Also! There's a naming convention

Kind of like how our apps are compiled for various architectures!



Open Source provides transparency behind models

Understanding model architecture & training data is critical!



Check the Model Openness Framework for understanding model components & licensing

Integrating Your Data To The Model

How can we specialize a "generic" language model to our unique data?

RAG, Fine-Tuning, etc Customizing responses **Processing Enterprise Documents** Not all data is Al-ready!

Add Your Data And Preferences <

Make it your own



First, consider your data!

Our data isn't always in one format, it's text, image, audio, etc.



Customization of Models



RAG Retrieval Augmented Generation

Enhance Gen Al model-generated text by retrieving relevant information from external sources, improving accuracy and depth of model's responses.

Fine tuning



Adjust a pre-trained model on specific tasks or data, improving its performance and accuracy for specialized applications without full retraining.

What is RAG?

A method that retrieves facts from an external knowledge base and causes the LLM to generate answers **based on accurate information**. The **original LLM is not modified**.



vector store (Corporate/personal data)

What is RAG?

A method that retrieves facts from an external knowledge base and causes the LLM to generate answers **based on accurate information**. The **original LLM is not modified**.



What is Fine-tuning?

Fine-tune the original LLM with your own data by retraining part or the entire LLM with a different data set.



What is Fine-tuning?

Fine-tune the original LLM with your own data by retraining part or the entire LLM with a different data set.



InstructLab vs. Alternative Model Alignment



InstructLab provides more accessible fine tuning & complements RAG

InstructLab powers simple & cost-effective LLM refinements

C	Foundational skills	-	Math
Taxonomy	Compositional skills		Writing Extraction
		A	Technical manual
	Knowledge		
		4	Textbook/history





Taxonomy-Driven Data Curation

Folder structure with Q&A pairs for topics to teach a model

Large-Scale Synthetic Data Generation

Generate additional training data to expand dataset automatically

Multi-Phase Alignment Tuning

Full parameter, phased alignment tuning for custom knowledge and skills Lab QR Code

39



https://catalog.demo.redhat.com/workshop/cf2e6x

Introducing the Docling Project docling

- Document Parsing: Extracts context, tables, graphs, and valuable data from PDF, Docx, etc
- AI-Ready: Direct integrations into LlamaIndex, LangChain, etc. for RAG and Dataset Gen.
- Performative: Fastest among all open-source parsers on CPU (+ 200k pages/day on GPU)
- For Developers: Released as a CLI for simple usage or a library for integrating into applications.





Building applications that use AI models

How can developers use, build, and test AI capabilities?

Serving Models Locally

How to not depend on 3rd party LLM API's

Building Apps with Al Capabilities Think RAG, Agentic, etc

Testing Model Output

How do we test non-deterministic output?



Introducing: Ramalama

To make AI boring by using containers

- Al in Containers: Run models with Podman/Docker with no config needed.
- Registry Agnostic: Freedom to pull models from Hugging Face, Ollama, or OCI registries.
- GPU Optimized: Auto-detect & accelerate performance.
- Flexible: Supports llama.cpp, vLLM, whisper.cpp & more.





Introducing: Podman AI Lab

For developers looking to build AI features

- For App Builders: Choose from various recipes like RAG, Agentic, Summarizers
- Curated Models: Easily access Apache 2.0 open-source options.
- Container Native: Easy app integration and movement from local to production.
- Interactive Playgrounds: Test & optimize models with your custom prompts and data.





Run Model in Production



Introducing vLLM vLLM docs



VLLM

For LLM inference serving in production environments

- Research-Based: UC Berkeley project to improve model speeds and GPU consumption
- Standardized: Works with Hugging Face & OpenAl API.
- Versatile: Supports NVIDIA, AMD, Intel, TPUs & more.
- Scalable: Manages multiple requests efficiently, ex. with Kubernetes as an LLM runtime





How can you save on GPU resources?

Quantization! It's a way to compress models, think like a .raw to .jpg

- Quantization: A technique to compress
 LLMs by reducing numerical precision.
- Converts high-precision weights (FP32) into lower-bit formats (FP16, INT8, INT4).
- Reduces memory footprint, making models easier to deploy.







How can you save on GPU resources?

Quantization! It's a way to compress models, think like a .raw to .jpg

- The Benefit? Run LLMs on "any" device, not just your local machine but IoT & Edge too
- Results in faster and lighter models that still maintain reasonable accuracy
 - Testing with Llama 3.1, for W4A16-INT resulted in 2.4x performance speedup and 3.5x model size compression
- Works on GPUs & CPUs!

Source:









& there's a open repository of Quantized Models

Check it out on Hugging Face & save resources on LLM serving!

Broad Collection



Comprehensive Validation

Benchmark	Meta-Llama-3.1- 70B-Instruct	Meta-Llama-3.1-70B-Instruct- FP8(this model)	Recovery
MMLU (5-shot)	83.83	83.73	99.88%
MMLU-cot (0-shot)	86.01	85.44	99.34%
ARC Challenge (0-shot)	93.26	92.92	99.64%
GSM-8K-cot (8-shot, strict-match)	94.92	94.54	99.60%
Hellaswag (10-shot)	86.75	86.64	99.87%
Winogrande (5-shot)	85.32	85.95	100.7%
TruthfulQA (0-shot, mc2)	60.68	60.84	100.2%
Average	84.40	84.29	99.88%

Extensive Selection

Formats W4/8A16 W8A8-INT8 W8A8-FP8 2:4 sparse	Algorithms GPTQ / AWQ SmoothQuant SparseGPT RTN
<u>Hardw</u>	<u>vare</u>
NVIDIA . GPUs AMDAMD Instinct	Google TPUs intel. CPUs



49

Cut GPU costs in half ready-to-deploy inference-optimized checkpoints



But what does an AI platform need to be successful?

Hint: It stems from open source technology!







So, we curate these projects into a trusted AI platform!





Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.



