



eccenca
mastering complexity

Edgard Marx | Principal Data Scientist

LEVERAGING AGENTIC DATA FLOWS WITH LARGE LANGUAGE MODELS



Outlook

Introduction

- Day-to-Day tasks

Motivation

- Diving into Large language Models (LLMS)
- A shared Vision

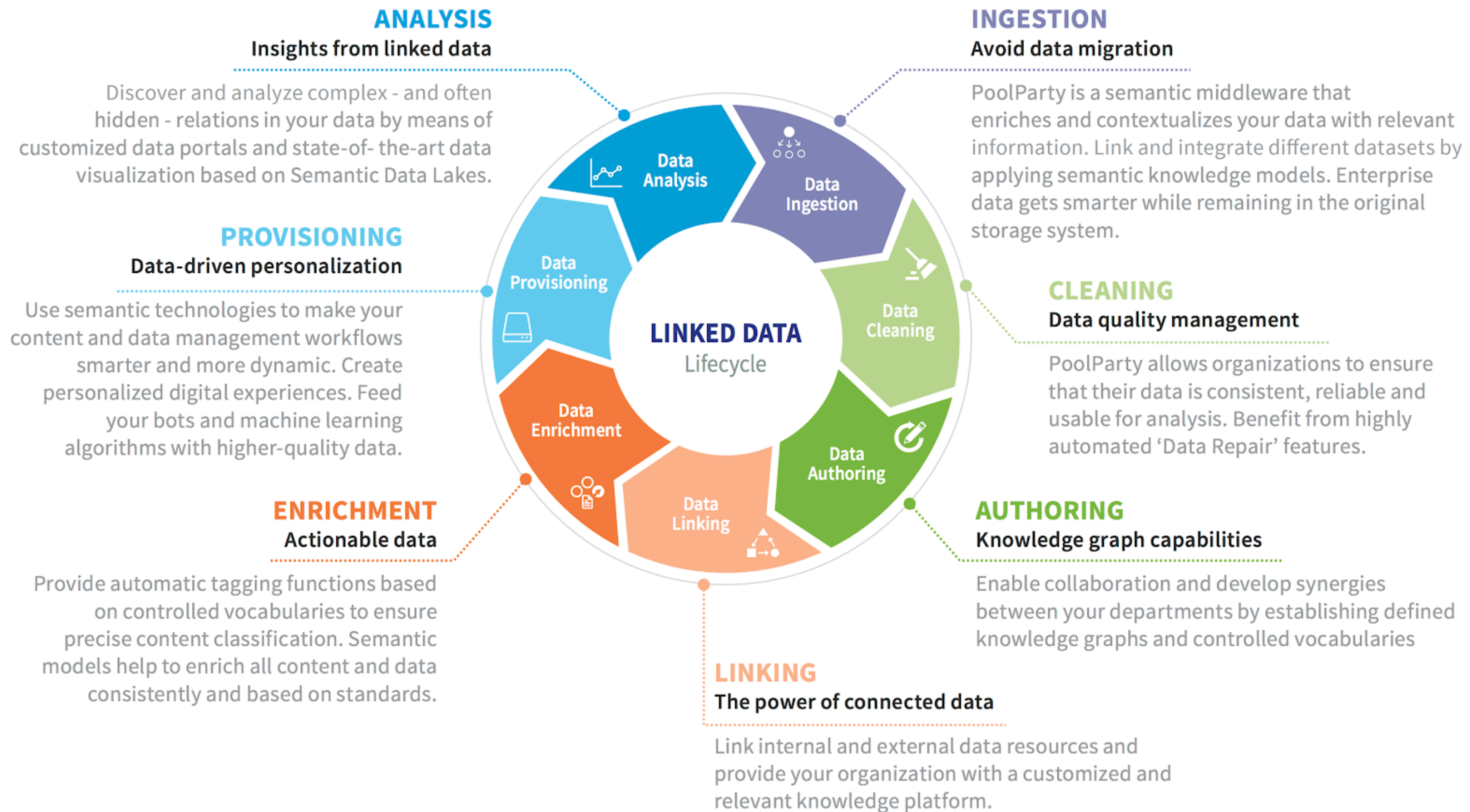
Showcase

- Practical Application

Conclusion

- Challenges
- Future Directions
- Questions

The (linked) data life cycle



Day-to-Day tasks



“Data practitioners spend 80% of their valuable time finding, cleaning, and organizing the data”

Challenges in Data Tasks

- Handling Missing Data
- Removing Duplicate Entries
- Correcting Data Types
- Standardizing Data Formats
- Dealing with Outliers
- Data Cleaning Tools and Techniques
- Addressing Inconsistent Data



Clonflicting Visions

“In fact, we’re not even close to matching the understanding of the physical world of any animal, cat or dog.”

Lecun

Sam Altman: OpenAI's New Model Passes AGI Threshold

BY PYMNTS | JANUARY 6, 2025

[f](#) [X](#) [in](#) [e](#) [v](#)



Meta Chief AI Scientist Slams Quest for Human-Level Intelligence

BY PYMNTS | JANUARY 8, 2025

[f](#) [X](#) [in](#) [e](#) [v](#)



TECH

Elon Musk Says Tesla Vehicles Will Drive Themselves in Two Years

BY KIRSTEN KOROSK

December 21, 2015 at 2:00 PM EST

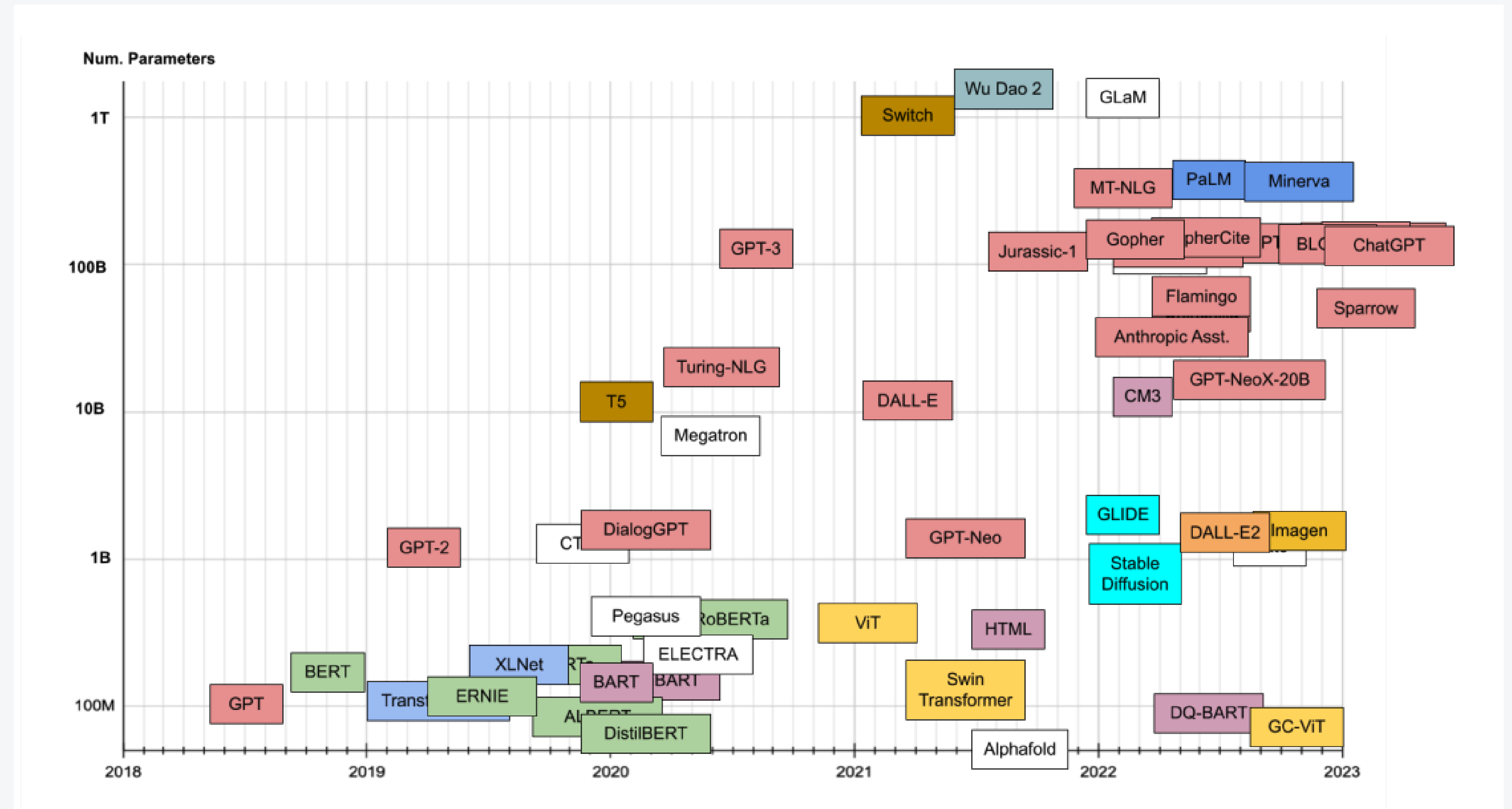


PHOTOGRAPH BY PATRICK FALLON — REUTERS

LLM Evolution

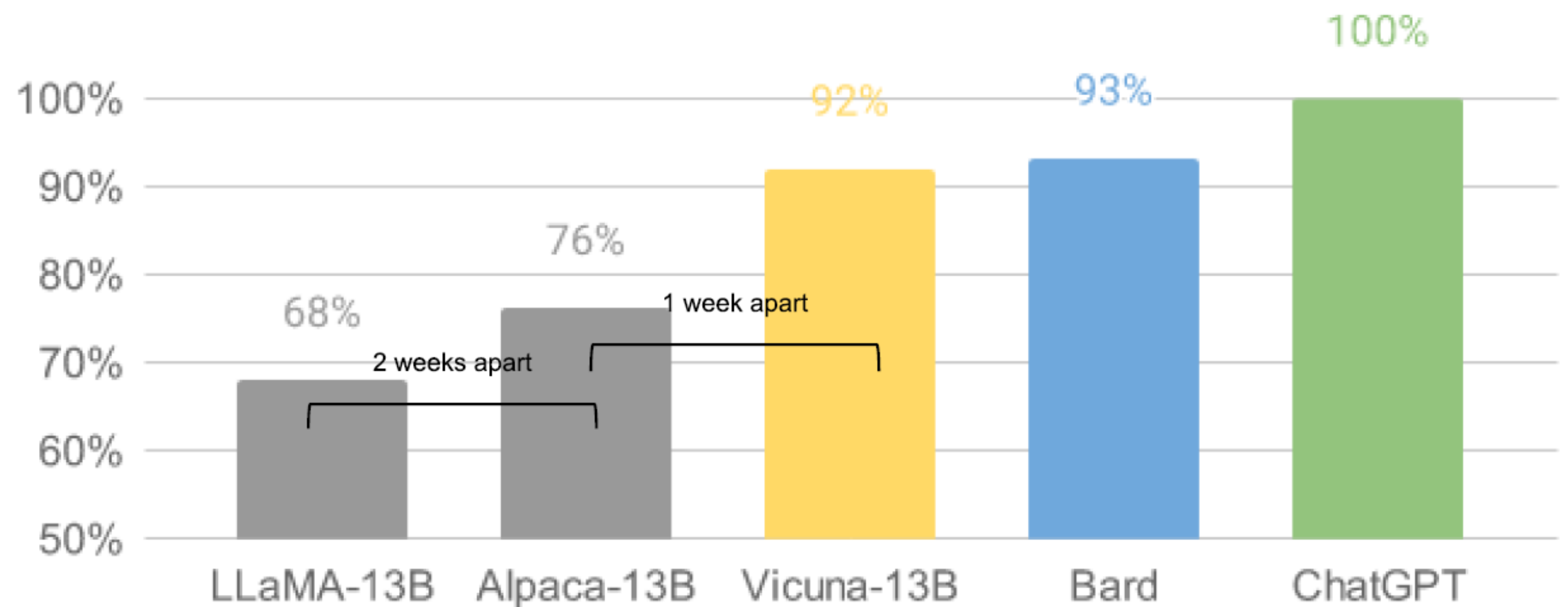
Open-Source LLMs

- MosaicML
- Falcon
- Open-Assistant
- Llama
- Qwen
- Mistral
- Phi
- DeepSeek



Fine-tuning

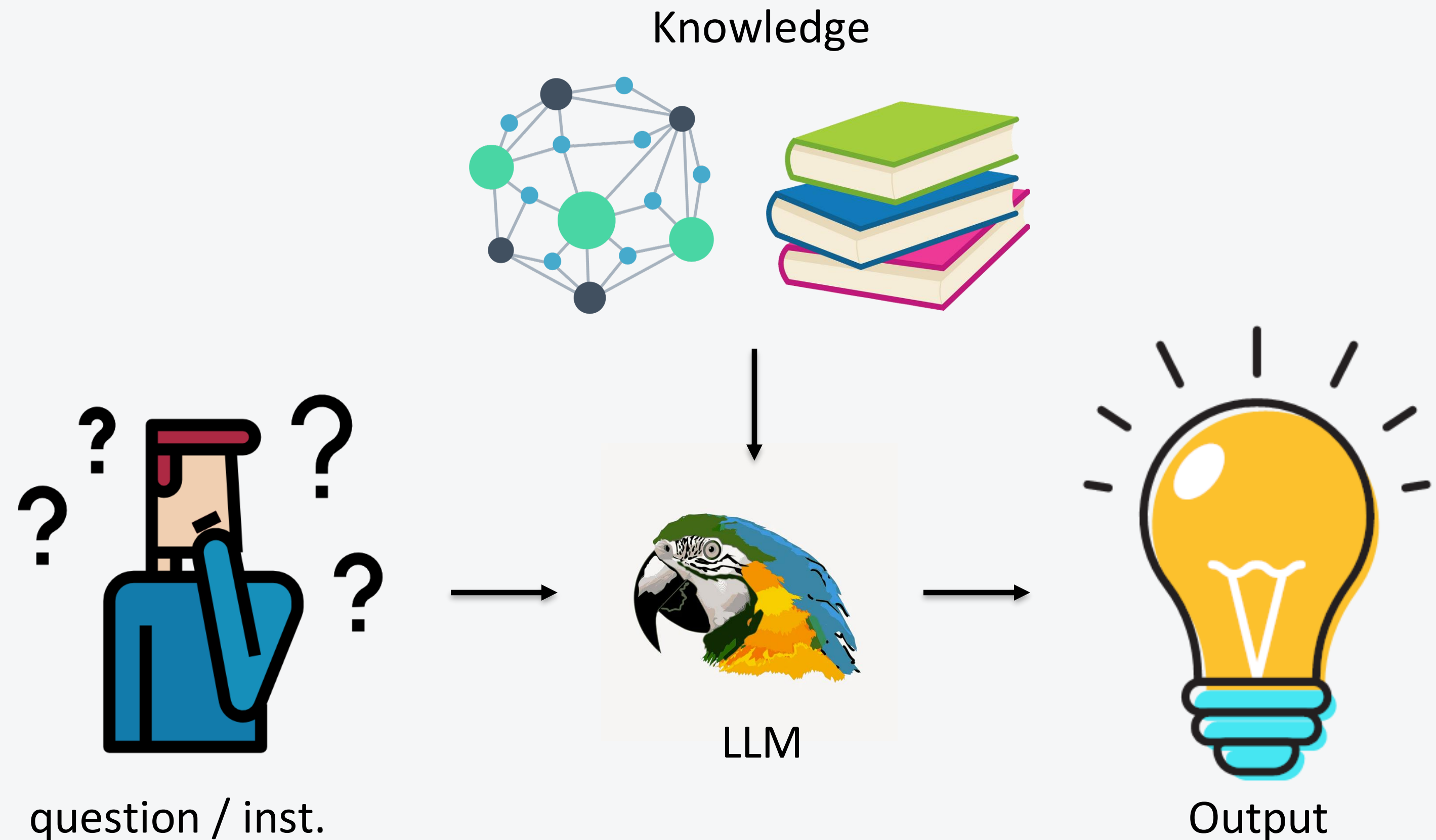
- **Deep Task Specialization**
 - **Adapts LLMs** to master niche domains (e.g., legal, medical, or brand-specific language).
 - **How?** By training on small, high-quality datasets → reshapes model weights for **domain expertise**.
→ *Results in unmatched accuracy for targeted use cases.*
- **Resource-Intensive but Powerful**
 - **Requires:**
 - Curated task-specific data (100s–1000s of examples).
 - Computational power (GPUs/TPUs) & technical expertise.
 - **Trade-off:** Higher upfront effort than RAG, but delivers **self-contained, optimized models** needing no live retrieval.
 - **Key Takeaway:** Fine-tuning creates purpose-built AI experts but demands significant investment.



*GPT-4 grades LLM outputs. Source: <https://vicuna.lmsys.org/>

Retrieval Augmented Generation (RAG)

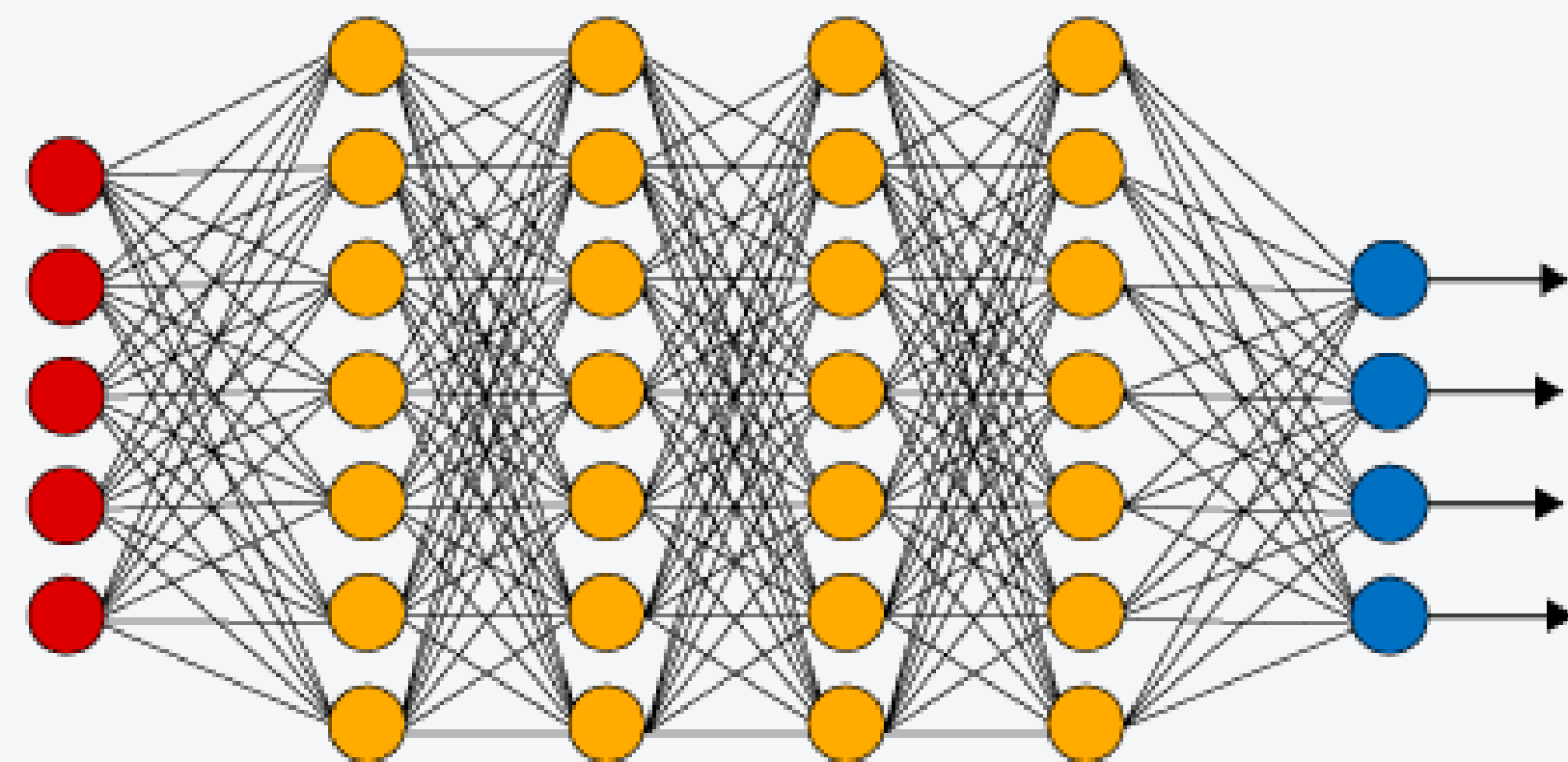
- **Retrieval-First Architecture**
 - **Step 1:** Retrieve relevant data from trusted sources (e.g., databases, documents).
 - **Step 2:** Augment LLM prompts with this context.
→ *Grounds responses in real-time, domain-specific knowledge.*
- **Mitigates Hallucinations & Outdated Info**
 - Uses live external data instead of relying solely on pre-trained model knowledge.
 - Reduces factual errors and keeps outputs current/accurate.
- **Cost-Efficient & Adaptable**
 - No retraining needed for new data—update the knowledge base instantly.
 - Works with any LLM (e.g., GPT, LLaMA) for flexible deployment.



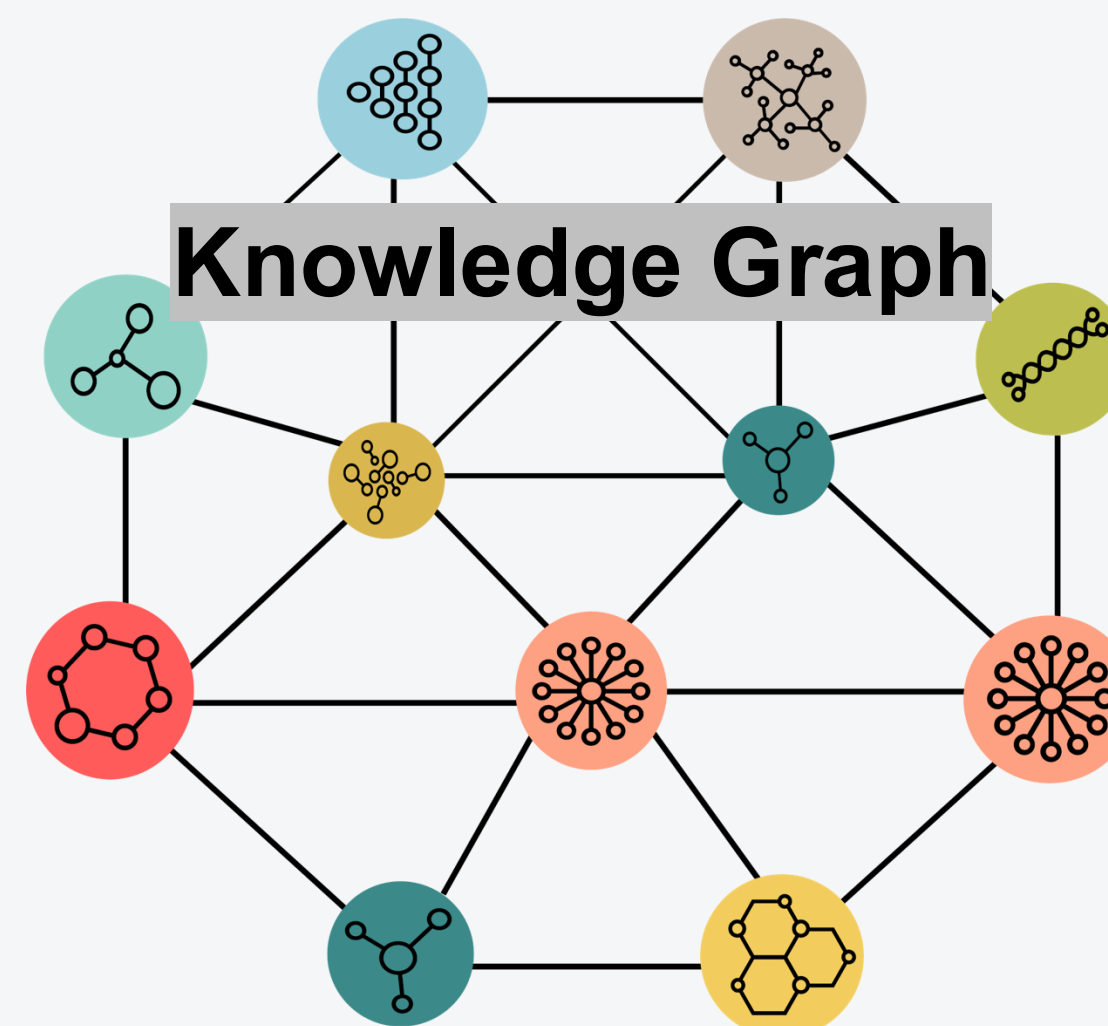
Best of both Worlds - Neuro-Symbolic AI

Integration of Knowledge Graphs with Machine Learning Models

Machine Intelligence



Background Knowledge



Human Intelligence



Knowledge authoring,
curation, validation


KG-LLM Integration Opportunities

Knowledge Graphs

- Represent a base of validated, trustworthy knowledge
- Help organize and integrate enterprise knowledge from various sources
- Provide input for enterprise and domain-specific training and fine-tuning of LLMs

Large Language Models

- Help curating knowledge in the KG by suggesting and recommending
- Create mappings, queries etc. for the KG
- Can become a frontend for human interaction with the KG



Practical Application

Navigating Trade-offs in LLM Customization

Fine-tuning

- **Resource Intensity**
 - Demands heavy computation (GPU/TPU), time, and technical expertise
 - Requires large *curated datasets* → costly/data-scarce domains struggle
- **Knowledge Rigidity**
 - Model "freezes" post-training → fails with new trends/updates
 - Retraining needed for adaptation → operational bottlenecks

RAG

- **Retrieval Dependency**
 - Output quality = $f(\text{knowledge base quality} + \text{retrieval accuracy})$
 - Hallucinations persist if retrieval fails or sources are outdated
- **Latency & Complexity**
 - Real-time retrieval → slower responses than pure LLMs
 - Synchronizing vector DBs/APIs adds architectural overhead
 - Size of the Context

Any Questions? Get in Touch!



eccenca GmbH

Hainstraße 8
D-04109 Leipzig
Germany

Try for free our Community
Edition Sandbox!

<https://eccenca.my>



+49 341 2650 8028



info@eccenca.com
<https://eccenca.com>