# HybridRAG: A Fusion of Graph and Vector Retrieval - Mitesh Patel, NVIDIA

8,818 views  Jul 22, 2025  AIEWF 2025 Complete Playlist

Interpreting complex information from unstructured text data poses significant challenges to Large Language Models (LLM), with difficulties often arising from specialized terminology and the multifaceted relationships between entities in document architectures. Conventional Retrieval Augmented Generation (RAG) methods face limitations in capturing these nuanced interactions, leading to suboptimal performance. In our talk, we introduce a novel approach integrating Knowledge Graph-based RAG (GraphRAG) with VectorRAG, designed to refine question-answering (Q&A) systems for more effective information extraction from complex texts. Our approach employs a dual retrieval strategy that harnesses both knowledge graphs and vector databases, enabling the generation of precise and contextually appropriate answers, thereby setting a new standard for LLMs in processing sophisticated data.

About Mitesh Patel

Mitesh Patel is a developer advocate manager at NVIDIA. His team is responsible for creating workflows to showcase how developers can harness GPU acceleration in their workflows using tools and frameworks popular in the developer community. Before NVIDIA, he was a senior research scientist at Fuji Xerox Palo Alto Laboratory Inc. (a research subsidiary of Fuji Xerox), where he worked on developing indoor localization technologies for applications such as asset tracking in hospitals and delivery cart tracking in manufacturing facilities. Mitesh received his Ph.D. in Robotics from the Center of Autonomous Systems (CAS) at the University of Technology Sydney, Australia in 2014.
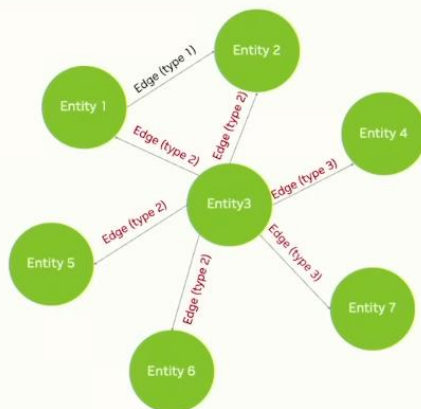


**⊚ NVIDIA.**

# HybridRAG: A fusion of Graph and Vector Retrieval to Enhance Data Interpretation

Mitesh Patel, Sr. Manager – Developer Advocate



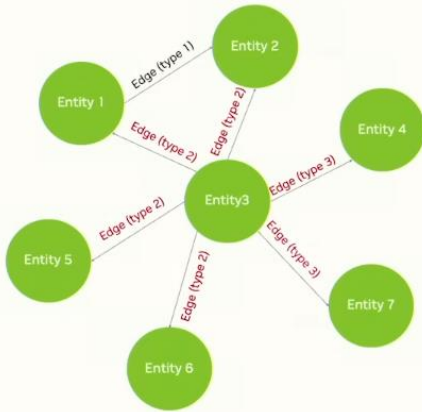## What is Knowledge Graph ?
### A refresher

- Network that represent relationships between different entities

- Entities can be objects, places, people, concepts or events

- The edges represents the relationship between entities

- "Triplets" : **[Entity 1 – Relationship – Entity 2]**

```
[['Exxon Mobil', 'COMP', 'Cut', 'Spending on oil and gas exploration', 'ACTIVITY']]
```

# What is Unique about Knowledge Graph
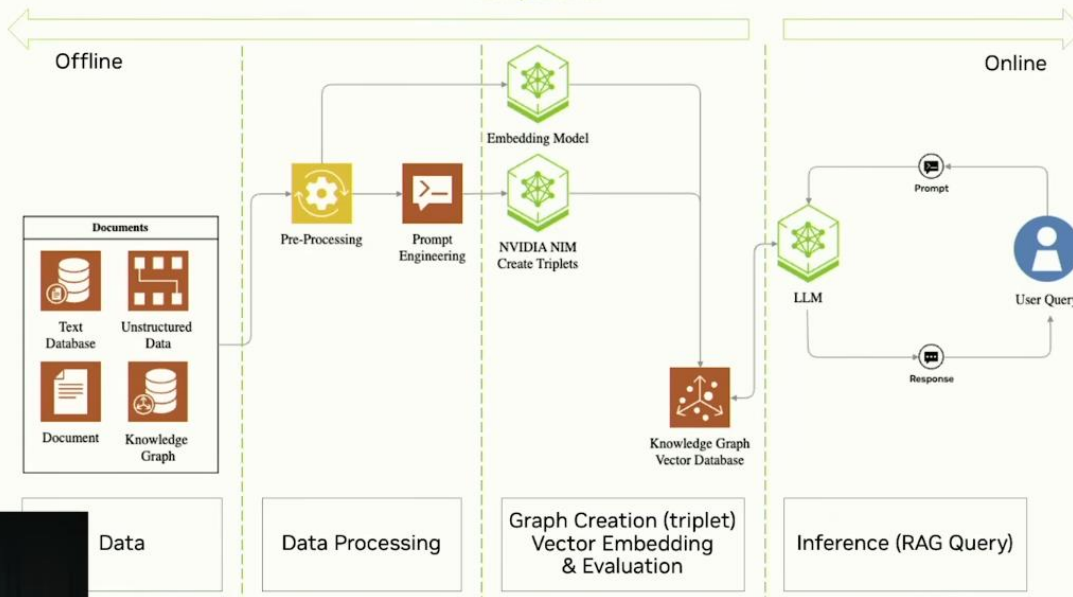## Where do they Excel



- Capture information about entities for a given domain or task

- Entity connection gives a comprehensive view of the knowledge base

- Knowledge Graph organize data from multiple sources

```
[['Exxon Mobil', 'COMP', 'Cut', 'Spending on oil and gas exploration', 'ACTIVITY']]
```
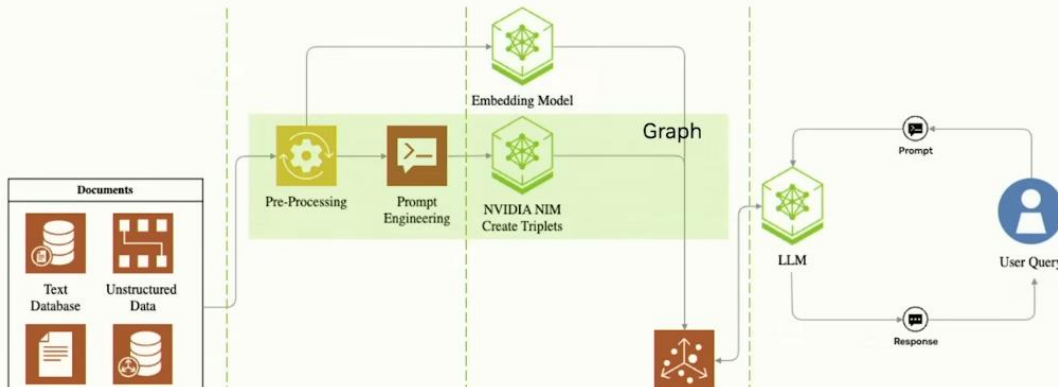
# How to Create a GraphRAG (Hybrid) System
## Components



| Data | Data Processing | Graph Creation (triplet) Vector Embedding & Evaluation | Inference (RAG Query) |

# How to Create a GraphRAG (Hybrid) System
## Components

## Creating the Triplets

| Text | Triplets |
|---|---|
| LONDON (Reuters) - Despite the strongest start for oil prices in four years, the world's top oil companies are hesitating to accelerate the search for new resources as a determination to retain capital discipline trumps the hope of making bonanza discoveries. | [ [ "Exxon Mobil", "COMP", "Cut", "Spending on oil and gas exploration", "ACTIVITY" ], [ "Royal Dutch Shell", "COMP", "Cut", "Spending on oil and gas exploration", "ACTIVITY" ], |
| Exxon Mobil, Royal Dutch Shell, Total and their peers are set to cut spending on oil and gas exploration for a fifth year in a row in 2018, according to consultancy Wood Mackenzie (WoodMac), despite a growing urgency to replenish reserves after years of reining back investment. | [ "Total", "COMP", "Cut", "Spending on oil and gas exploration", "ACTIVITY" ], [ "World's top oil companies", "ORG", "Hesitate", "Accelerate the search for new resources", "ACTIVITY" ], |
| (For graphic 'Global spending on oil and gas exploration' click reut.rs/2CjAONv ) | [ "Consultancy Wood Mackenzie", "ORG", "Estimate", "Global investment in exploration", "ECON_INDICATOR" ], |
| Global investment in exploration, vital to increase output and offset the natural decline of existing fields, will reach $37 billion in 2018, down 7 percent from a year earlier and over 60 percent below the 2014 peak, according to WoodMac. | [ "Global investment in exploration", "ECON_INDICATOR", "Reach", "$37 billion", "VALUE" ], [ "Global investment in exploration", "ECON_INDICATOR", "Decrease", "7 percent"," PERCENTAGE"] |

Triplets expose the relationship/information between 2 entities, this information is helpful. We can use LLMs to extract this information for us to save it in our KB in triplet format.
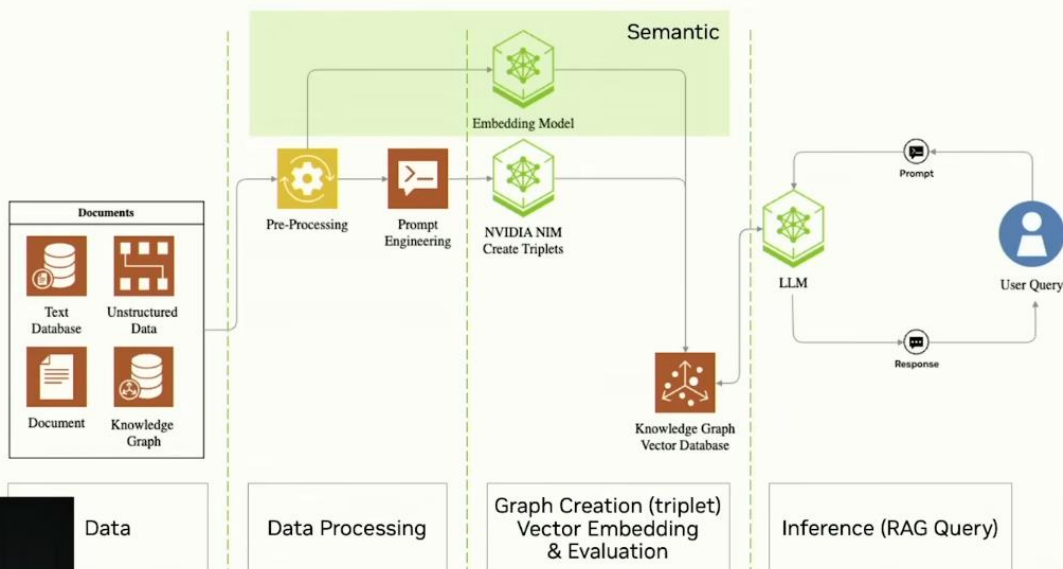
## Creating the Triplets
### Prompts is all that Matters

```
def extract_triples(text, llm):
    prompt = ChatPromptTemplate.from_messages(
    [("system", """Note that the entities should not be generic, numerical, or temporal (like dates or percentage
- ORG: Organizations other than government or regulatory bodies
- ORG/GOV: Government bodies (e.g., "United States Government")
- ORG/REG: Regulatory bodies (e.g., "Food and Drug Administration")
- PERSON: Individuals (e.g., "Marie Curie")
- GPE: Geopolitical entities such as countries, cities, etc. (e.g., "Germany")
- INSTITUTION: Academic or research institutions (e.g., "Harvard University")
- PRODUCT: Products or services (e.g., "CRISPR technology")
- EVENT: Specific and Material Events (e.g., "Nobel Prize", "COVID-19 pandemic")
- FIELD: Academic fields or disciplines (e.g., "Quantum Physics")
- METRIC: Research metrics or indicators (e.g., "Impact Factor"), numerical values like "10%" is not a METRIC;
- TOOL: Research tools or methods (e.g., "Gene Sequencing", "Surveys")
- CONCEPT: Abstract ideas or notions or themes (e.g., "Quantum Entanglement", "Climate Change")
```

## How to Create a GraphRAG (Hybrid) System
### Components

# Create Semantic Vector Database

- Document chunking size

- Chunking overlap size

- Embedding Model

This is where you convert your documents into chunks and put in a vector DB

# Accelerate : Knowledge Graph Retrieval Strategies
## Knowledge Graphs can Increase RAG Accuracy

- Search beyond a single hop/node (Multi-hops):
  - Helps in retrieve all relevant nodes
  - Provides far greater context between entities and relations

- Accelerate Search on GPUs through cuGraph
  - Accelerate multi-hop & multi-step search
  - BFS/DFS strategies accelerated

You can query your data using different strategies like depth or breadth approaches

**Evaluate Knowledge Graph**
Is Your KG good enough

- RAGAS: evaluation framework for RAG systems
  - Faithfulness
  - Answer Relevancy
  - Context Precision
  - Context Recall
- Nemotron-340b-reward: A reward model that scores the text output
  - Helpfulness
  - Correctness
  - Coherence
  - Complexity
  - Verbosity

RAGAS is a pip install library that allows you to bring your own model



**The Last Mile**
Improve Performance & Accuracy

The way you create your Knowledge Graph (KG) can help improve your system performance and results

# The Last Mile
## Improve Performance & Accuracy

### Apostrophe

```
[['HP', 'COMP', 'Recall', 'Laptop Batteries', 'PRODUCT'],
['Laptop Batteries', 'PRODUCT', 'Operate_In', 'All HP Laptops',
'GPE'], ['Affected Laptops', 'GPE', 'Has', 'Burn Hazard Risk',
'CONCEPT'], ['Eligible Batteries', 'PRODUCT', 'Replace',
'Free', 'FIN_INSTRUMENT'], ['Consumers', 'PERSON', 'Check',
'HP's Website', 'ORG'], ['HP', 'COMP', 'Push', 'Battery Safety
Mode Update', 'PRODUCT'], ['Affected Computers', 'GPE',
'Activate', 'Battery Safety Mode', 'PRODUCT']]
```

**Add to prompt**
Do not output apostrophe removes apostrophes for example
convert Apple's to Apple , Korea's to Korea etc

**clean the output**
.replace("'s", "s")

### Really long output

```
[[['I', 'PERSON', 'Visit', 'Puerto Rico', 'GPE'],
['Hurricane Maria', 'EVENT', 'Relate_To', 'Puerto Rico',
'GPE'], ['I', 'PERSON', 'Feel', 'Guilty', 'CONCEPT'],
['Everybody', 'PERSON', 'Exhausted', 'CONCEPT'],
…..
…..
['U.S. Army Corps of Engineers', 'ORG', 'Classified',
'Total_Impact_On', 'Puerto Rico', 'GPE'], ['FEMA', 'ORG',
'Classified', 'Impacts', 'ECON_INDICATOR', 'Puerto Rico',
''',……… // Really long Output - Incomplete output
```
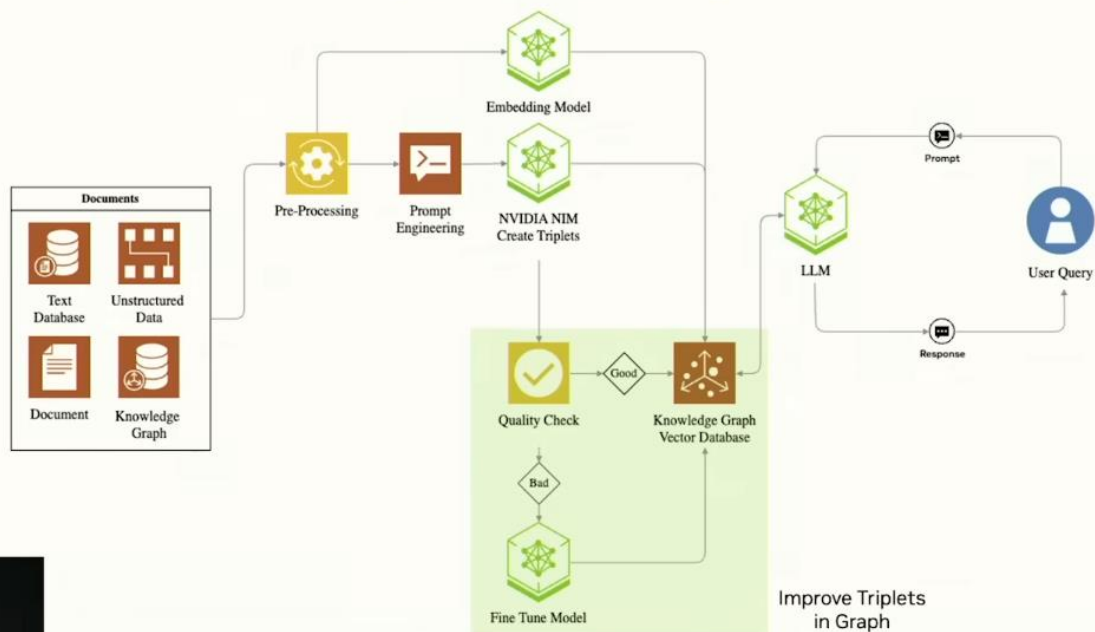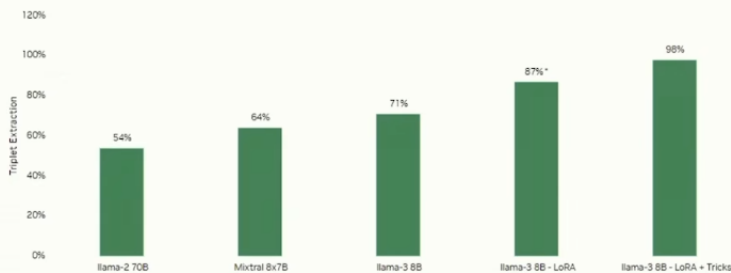
**Implement Correction**
In case of missing entries at the end of list
truncate it for example
INPUT : `[['Polaris Industries Inc.', 'COMP',
'Announce', 'fourth quarter and full-year 2017
financial results', 'ECON_INDICATOR'], ['Polaris
Industries Inc.', 'COMP', 'Hold', 'webcast and
conference call', 'EVENT'], ['webcast and conference
call', 'EVENT', 'Hosted_By', 'Scott Wine',
'PERSON`
OUTPUT : [['Polaris Industries Inc.', 'COMP',
'Announce', 'fourth quarter and full-year 2017
financial results', 'ECON_INDICATOR'], ['Polaris
Industries Inc.', 'COMP', 'Hold', 'webcast and
conference call', 'EVENT']]

# The Last Mile
## Improve Performance & Accuracy



**Improvement in Accuracy over 100 document**

- llama-2 70B: 54%
- Mixtral 8x7B: 64%
- llama-3 8B: 71%
- llama-3 8B - LoRA: 87%*
- llama-3 8B - LoRA + Tricks: 98%

# The Last Mile
## Improve Performance & Accuracy



If your graph gets really big with lots of nodes, you might start having latency and network issues. Tweak your logic and measure the network performance

# The Last Mile
## Improve Performance & Accuracy

The full end-to-end accelerated NetworkX experience is now enabled by setting the NX_CUGRAPH_AUTOCONFIG environment variable to True.

```
%env NX_CURGAPH_AUTOCONFIG=True

import pandas as pd
import networkx as nx

url = "https://data.rapids.ai/cugraph/datasets/cit-Patents.csv"
df = pd.read_csv(url, sep=" ", names=["src", "dst"], dtype="int32")
G = nx.from_pandas_edgelist(df, source="src", target="dst")

%time result = nx.betweenness_centrality(G, k=10)
```
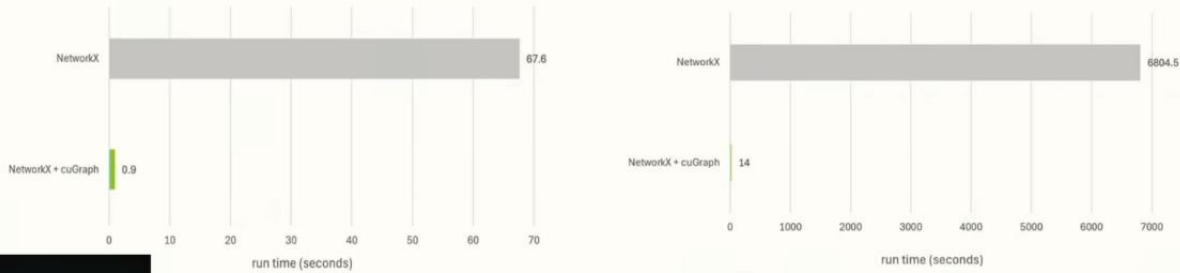
| | run time (seconds) | | run time (seconds) |
|---|---|---|---|
| NetworkX | 67.6 | NetworkX | 6804.5 |
| NetworkX + cuGraph | 0.9 | NetworkX + cuGraph | 14 |

*...m run on a citation graph of U.S. patents (4M nodes, 16M edges) is 70x faster than NetworkX on CPU;*
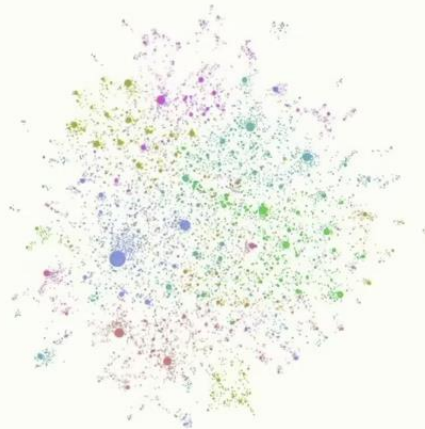
*Figure 2. The betweenness centrality algorithm run on the Live Journal social network (5M nodes, 69M edges) is 485x faster than NetworkX on CPU for number of samples (k) set to 100.; SW: NetworkX 3.4.1, cuGraph/nx-cugraph 24.10; GPU: NVIDIA A100 80GB; CPU: Intel Xeon w9-3495X (56 cores) 250GB*

---

# Should I Use Graph or Semantic or Hybrid
## What to use ?

- Your data
  - Structured data
  - Semantic/Unstructured data
  - Can graph be created

- Application/Usecase
  - Complex relationship
  - Semantic relationship
  - Latency

---

# Developer Tools and Resources
## Accelerate innovation and growth

Learn more: developer.nvidia.com

## Individuals

| Software | Training | GPU Sandbox |
|---|---|---|
| 100s of APIs, models, SDKs, microservices, and early access to NVIDIA tech | Hands-on self-paced courses, instructor-led workshops, and certifications | Approval basis, multi-GPU and multi-node |
| **Learning** | **Community** | **Ecosystem** |
| Tutorials, self-paced courses, blogs, documentation, code samples | Dedicated developer forums, meetups, hackathons | GTC, NVIDIA Partner Network |

## Organizations

| Startups | Venture Capital | Higher Education |
|---|---|---|
| Cloud credits, engineering resources, technology discounts, exposure to VCs | Deal flow and portfolio support for Venture Capital firms | Teaching kits, training, curriculum co-development, grants |
| **ISVs and SIs** | **Research** | **Enterprises** |
| Engineering guidance, discounts, marketing opportunities | Grant programs, collaboration opportunities | Tailored developer training, skills certification, technical support |