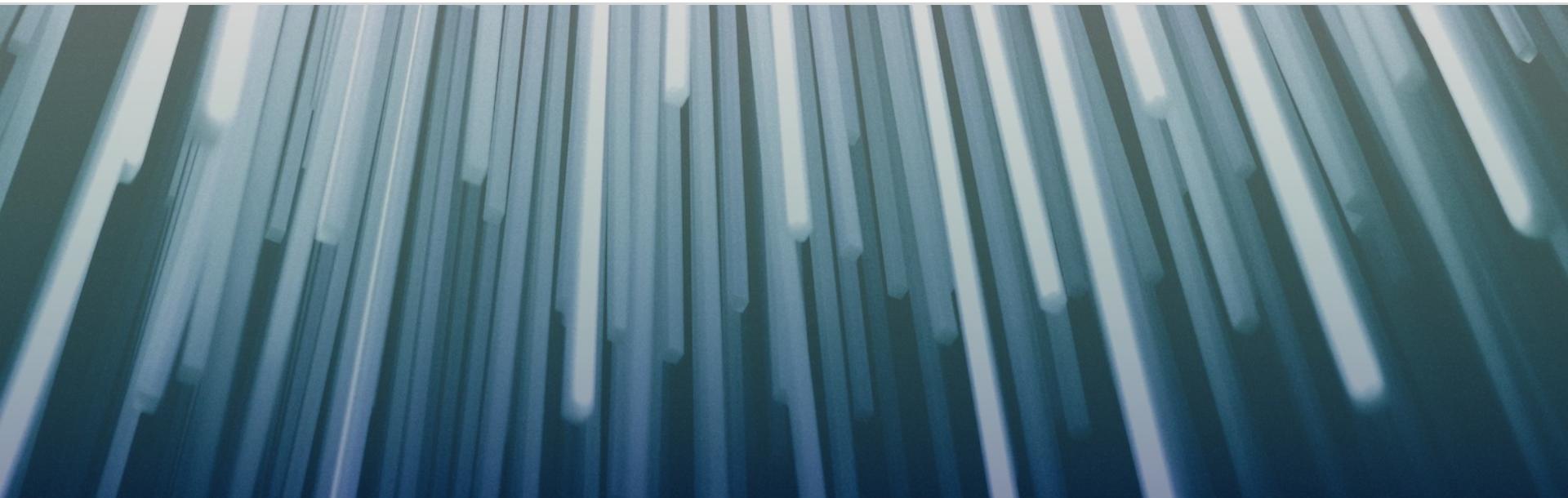


March 19, 2024

Architecting Enterprise AI success with RAG and LLMs

Lessons from the First 12 Months of Building Generative AI Solutions





AGENDA

- 01** Generative AI - A Phase Transition in Computing
- 02** Enterprise Application Patterns
- 03** Retrieval Augmented Generation
- 04** Case Studies
- 05** Lessons Learnt Architecting LLM Deployments

01

Busy executive's guide to LLMs

How did we get here?

GENERATIVE AI - A PHASE TRANSITION IN COMPUTING

- Jensen Huang, NVIDIA



VALUE CREATION WITH NEW APPLICATIONS

AI-Employee Startups are all the rage in Venture Capital

Emulating the role (copilot)

Accelerating operations

finpilot
Financial Analyst

V
Product Manager
Copilot

TextQL
Data Scientist

NormAi
Compliance

Cognition

AI Software
Developer

Magic

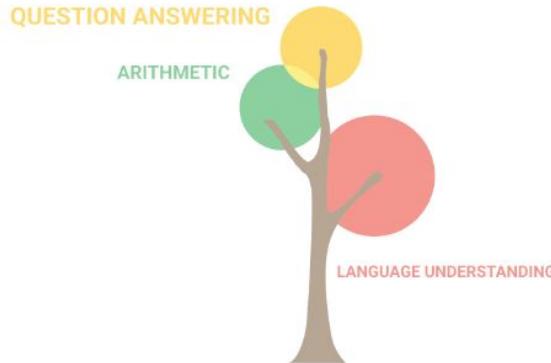
Copilot for Software
Developer

SIERRA

Customer Support

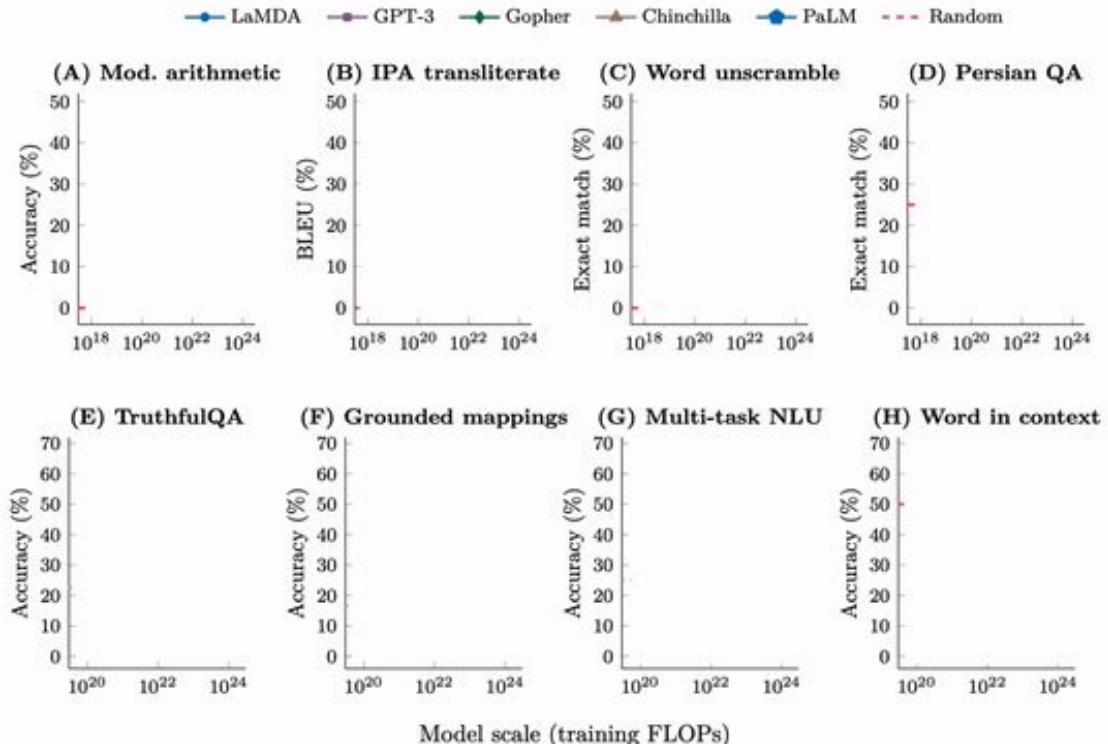
Funding Raised (in Millions)

EMERGENT CAPABILITIES + IN-CONTEXT LEARNING IS ... MARVELOUS



8 billion parameters

EMERGENT CAPABILITIES + IN-CONTEXT LEARNING IS ... MARVELOUS



02

Enterprise Applications

How do Enterprises want to use Generative AI?

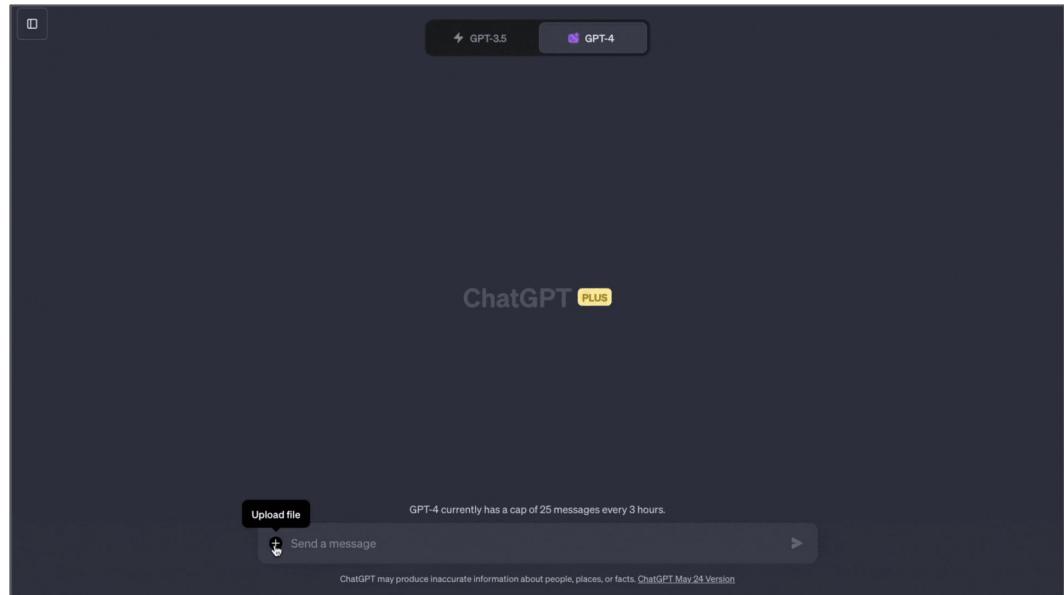
We conducted 300
workshops

This is what we have learnt

We conducted 300
workshops

This is what we have learnt

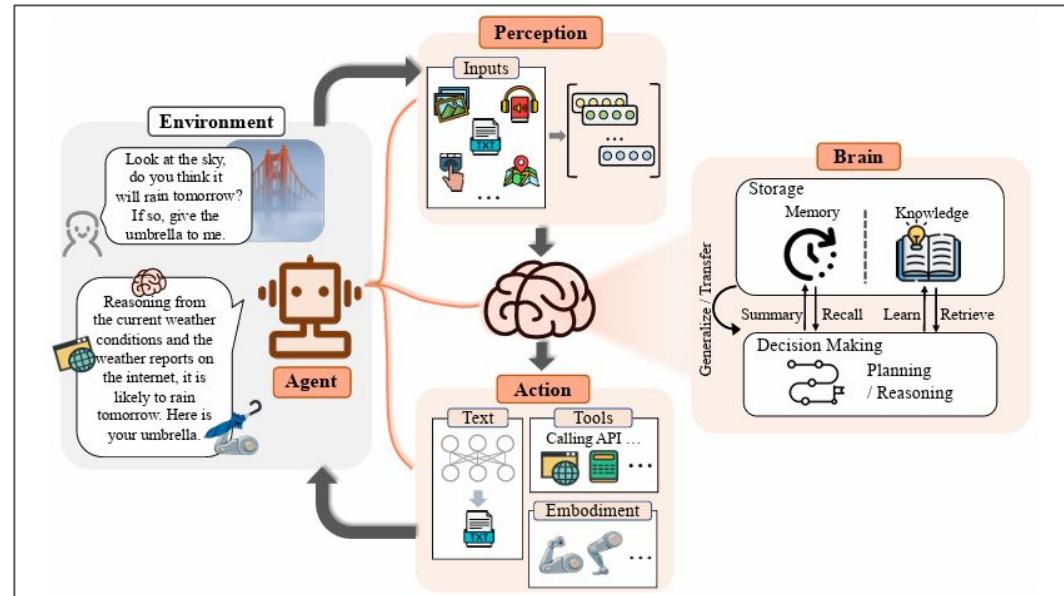
- 1 Help me talk to my data



We conducted 300 workshops

This is what we have learnt

- 1 Help me talk to my data
- 2 Help me Automate and Plan Tasks



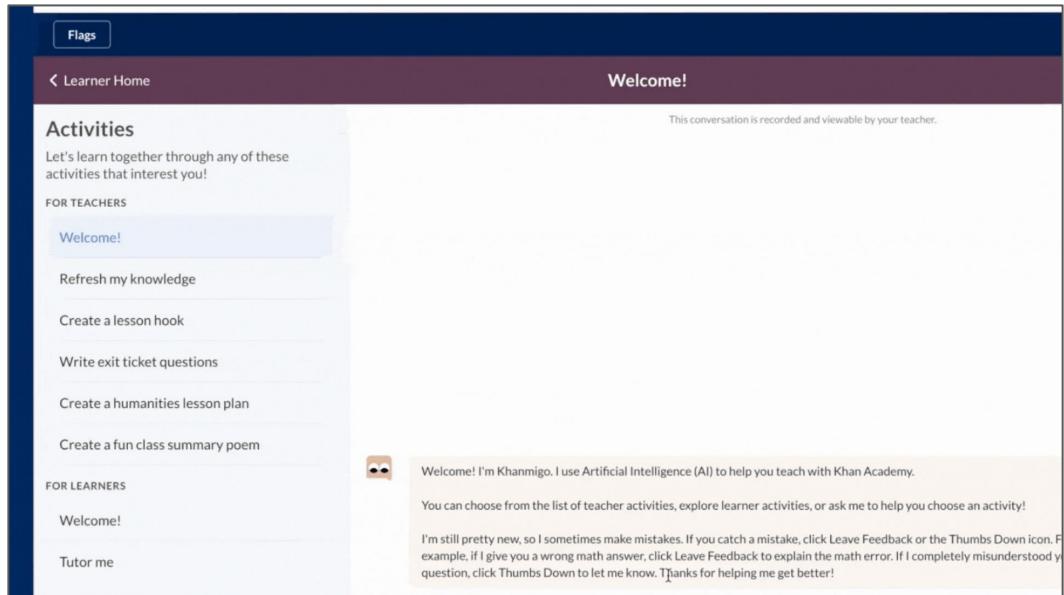
We conducted 300 workshops

This is what we have learnt

1 Help me talk to my data

2 Help me Automate and Plan Tasks

3 Elevate Conversations



The screenshot shows the Khan Academy Learner Home page. At the top right, there's a "Welcome!" message with a note: "This conversation is recorded and viewable by your teacher." Below it, under "Activities", there's a section for teachers with options like "Welcome!", "Refresh my knowledge", "Create a lesson hook", "Write exit ticket questions", and "Create a humanities lesson plan". There's also a "Create a fun class summary poem" option. Under "FOR LEARNERS", there are "Welcome!" and "Tutor me" options. On the right side, there's a sidebar with a cartoon character icon and a message from "Khanmigo": "Welcome! I'm Khanmigo. I use Artificial Intelligence (AI) to help you teach with Khan Academy. You can choose from the list of teacher activities, explore learner activities, or ask me to help you choose an activity! I'm still pretty new, so I sometimes make mistakes. If you catch a mistake, click Leave Feedback or the Thumbs Down icon. For example, if I give you a wrong math answer, click Leave Feedback to explain the math error. If I completely misunderstand your question, click Thumbs Down to let me know. Thanks for helping me get better!"

We conducted 300
workshops

This is what we have learnt

- 1 Help me talk to my data
 - 2 Help me Automate and Plan Tasks
 - 3 Elevating Conversations
 - 4 Help me get productive and ship code quicker

The screenshot shows a Jupyter Notebook interface with three code cells:

- utils.py**:

```
1 import torch
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.metrics import accuracy_score
6
7 def load_data():
8     iris = datasets.load_iris()
9     X = iris.data
10    y = iris.target
11
12    # Standardize the data
13    scaler = StandardScaler()
14    X = scaler.fit_transform(X)
15
16    X_train, X_test, y_train,
17    y_test = train_test_split(X, y, test_size=0.3,
18                             random_state=42)
19
20    # Convert numpy data to PyTorch tensors
21    X_train = torch.tensor(X_train, dtype=torch.
22                          float32)
```
- main.py**:

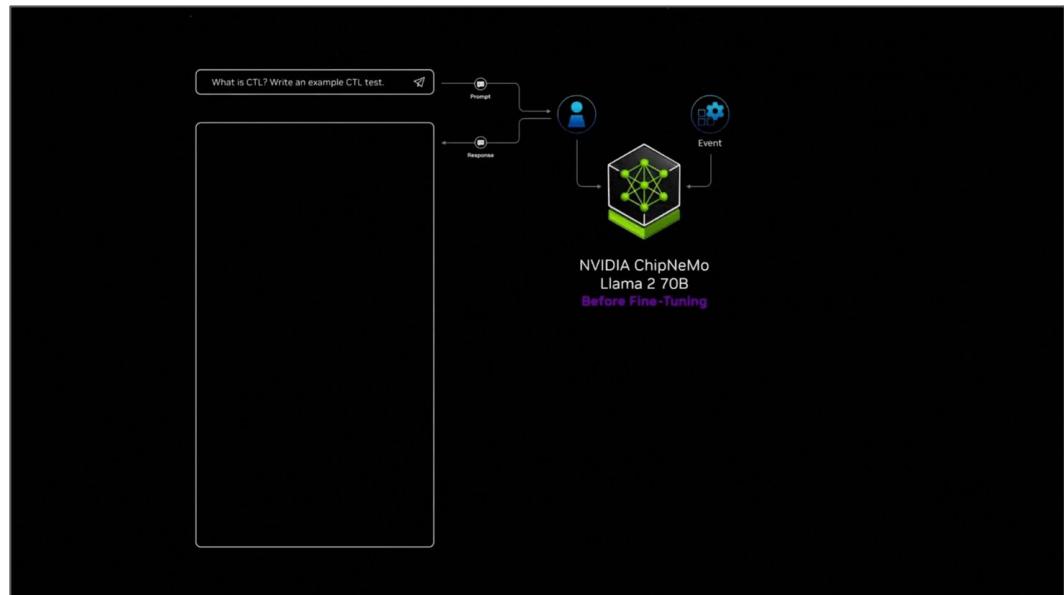
```
1 from utils import load_data,
2         evaluate_predictions
3 from model import IrisClassifier as Classifier
4
5 def main():
6     # Model training and evaluation
7
8     if __name__ == '__main__':
9         main()
```
- model.py**:

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from torch.utils.data import
5     DataLoader, TensorDataset
6
7 class IrisClassifier(nn.Module):
8     def __init__(self):
9         super(IrisClassifier, self).
10            __init__()
11         self.fc = nn.Sequential(
12             nn.Linear(4, 16),
13             nn.ReLU(),
14             nn.Linear(16, 3)
15         )
16
17     def forward(self, x):
18         return self.fc(x)
19
20     def train_model(self, X_train,
21                    y_train, epochs, lr, batch_size):
22         criterion = nn.
23             CrossEntropyLoss()
24         optimizer = optim.Adam(self.
25             parameters(), lr=lr)
26
27         # Create DataLoader for
28         # batches
29         dataset = TensorDataset
30         (X_train, y_train)
```

We conducted 300 workshops

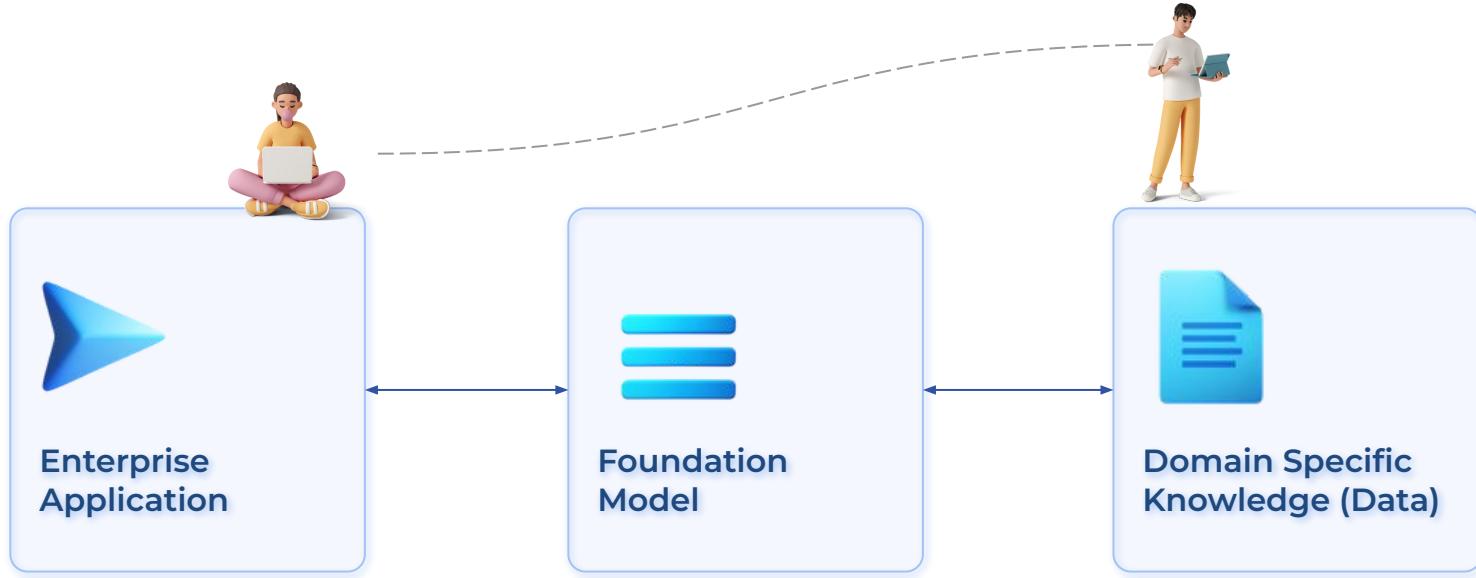
This is what we have learnt

- 1 Help me talk to my data
- 2 Help me Automate and Plan Tasks
- 3 Elevating Conversations
- 4 Help me get productive and ship code quicker
- 5 Help me catalyze growth with my Data Moat



AUGMENTING LLMS WITH RETRIEVAL: THE OPEN BOOK EXAM

Enhancing Task-Specific Performance



Enhanced Context Understanding

Tailored Responses

Reduced Irrelevance

03

CASE STUDIES: Early success with boosting knowledge worker productivity

RAG SUCCESS: CASE STUDIES



Dialog Agent for Field Technicians

Intelligent RAG-based agent to assist Telecom Field Operators in configuration, installation, and troubleshooting network and customer devices and minimize supervisor involvement



90% Fewer Supervisor Calls, quicker turnaround time and cost reduction

LLM Agents for Banking IT Operations Management

LLM-powered Q&A and insights from SQL logs, and IT product marketplace using knowledge graphs, to monitor and manage IT Ops based on defined IT access rules and protocols



5X Improvement in IT Incident Prioritization and Resolution*



Clinical Trial Protocol Matching

RAG-powered system to automate matching process of trial-eligible patients using EHR data, based on alignment with trial protocols and providing eligibility rationale



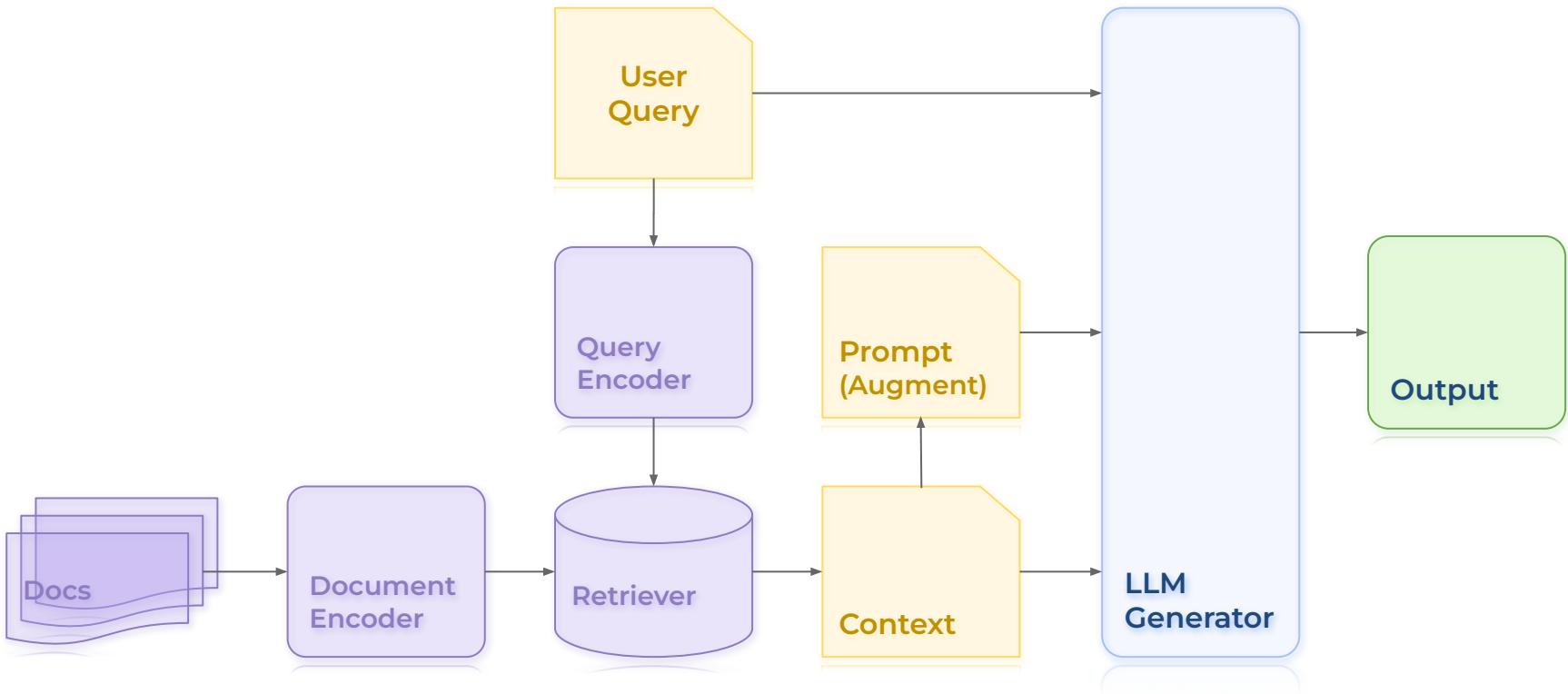
90% Reduction in Patient Matching Turnaround Time (TAT)



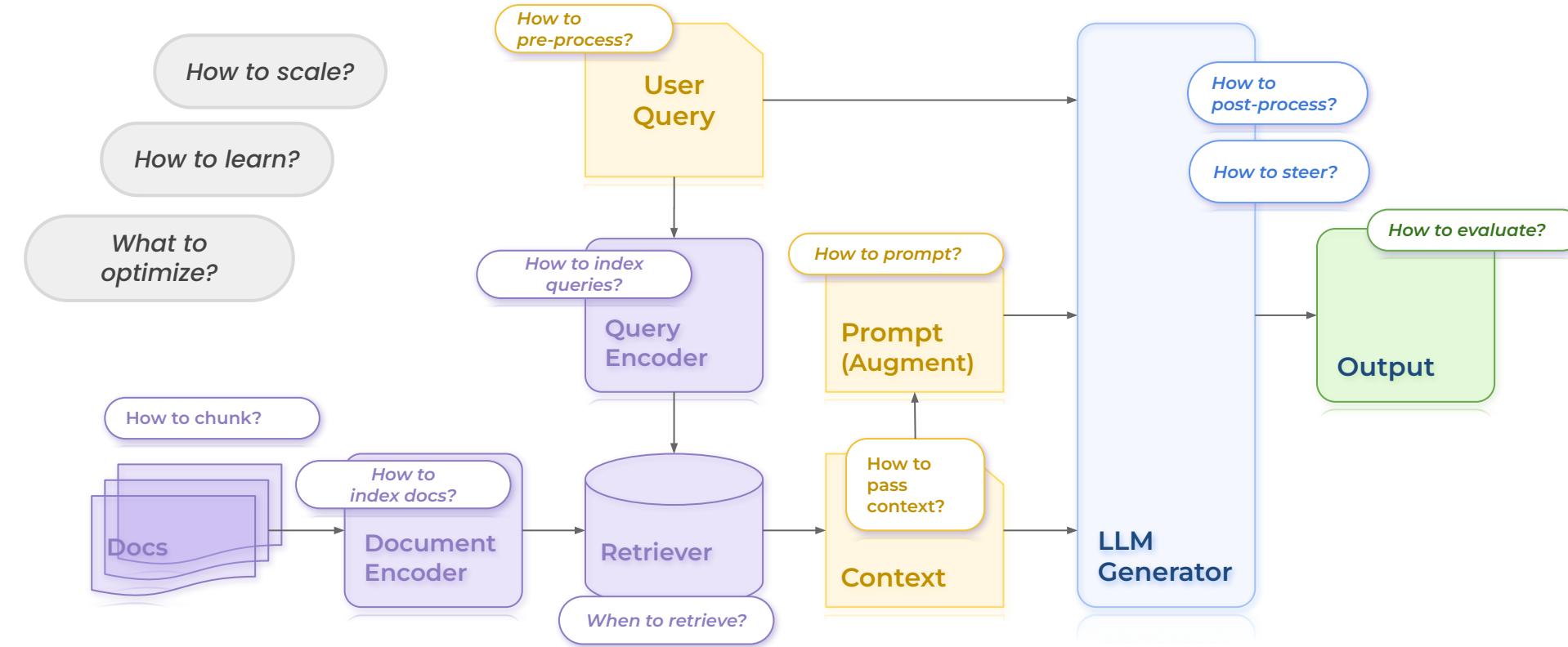
04

No Free Lunch: The Hidden Technical Debt of Building RAG Applications

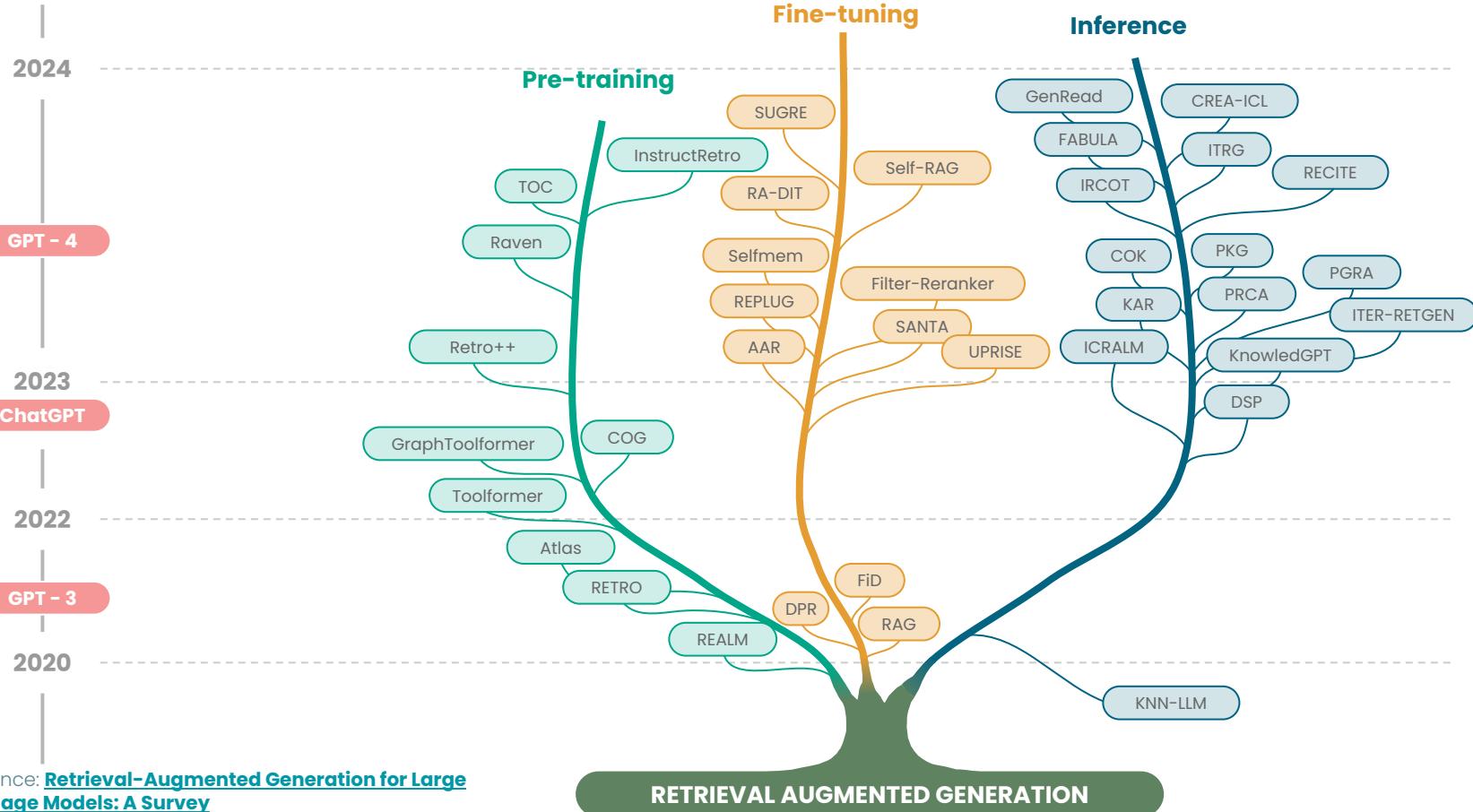
THE EASY: BUILDING A RAG APPLICATION, GREAT TOOLING!



THE COMPLEX: SYSTEM DESIGN FOR RAG APPLICATIONS



THE UNKNOWN: EVER EVOLVING FRONTIER OF RAG



06

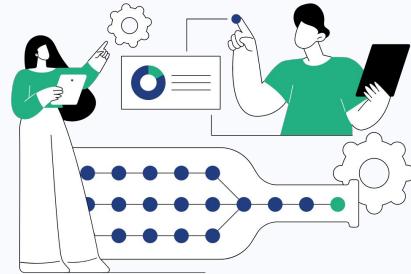
Conclusion

LESSONS LEARNT FROM PRACTICAL DEPLOYMENTS

PRACTICAL GENERATIVE AI SYSTEMS ARE COMPOUND AI SYSTEMS



MAP DOMAIN REQUIREMENTS TO SYSTEM CONSTRAINTS EARLY



BUILD WITH NVIDIA AI ENTERPRISE SOFTWARE



“

RAGS are a Full Stack Accelerated Computing Challenge

..... Each Prompt is an Inference Request

A Call to Multiple Components

Each Component Needs Acceleration

”





2024
Americas Service Delivery
Partner of the Year

3X
Americas Service Delivery
Partner of the year

Join us at:

- **Booth #1513**
AI CoE Pavillion
- **Booth #G129**
Generative AI Pavillion



*Scan to
accelerate
your AI
journey*

Thank you