

# OPEA-based Retrieval Augmented Generation (RAG) on Intel® Gaudi with OpenShift AI

**Red Hat Summit Connect 2024 Denmark**

*Copenhagen, 29 October 2024*



# Codrin Bucur

Principal AI Specialist Solution Architect,  
EMEA, Red Hat





Over **25** Years of Collaboration



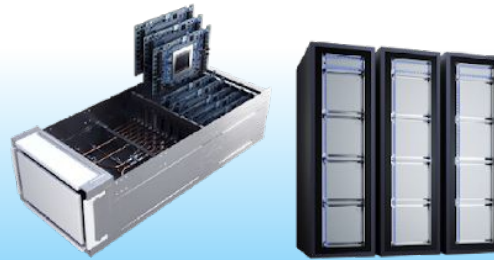
# Bringing AI Everywhere

## Intel's AI Strategy



AI PC Node  
AI Developer Productivity & Light  
Inference

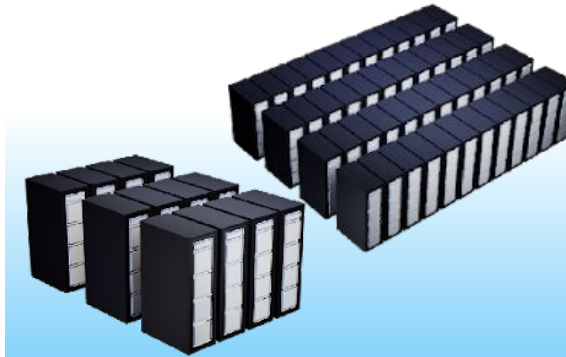
AI PC  
Broadest AI SW Ecosystem



Node  
Fine-tuning,  
Inference

Cluster  
Light Training, Tuning, Peak  
Inference

ENTERPRISE AI & EDGE AI  
Open Standard, "Ready to Use"



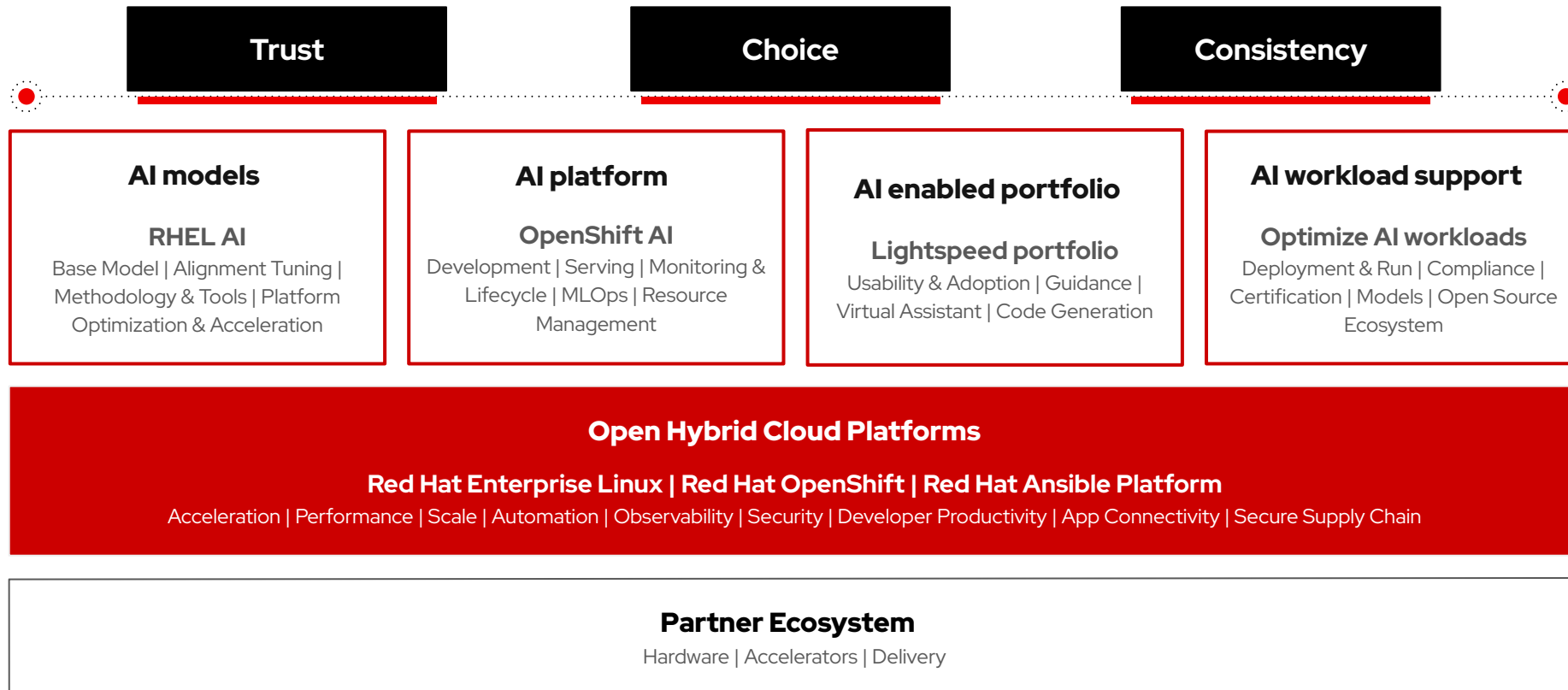
Super Cluster  
Training, Tuning, Peak  
Inference

Mega Cluster  
Large Scale Training  
& Inference

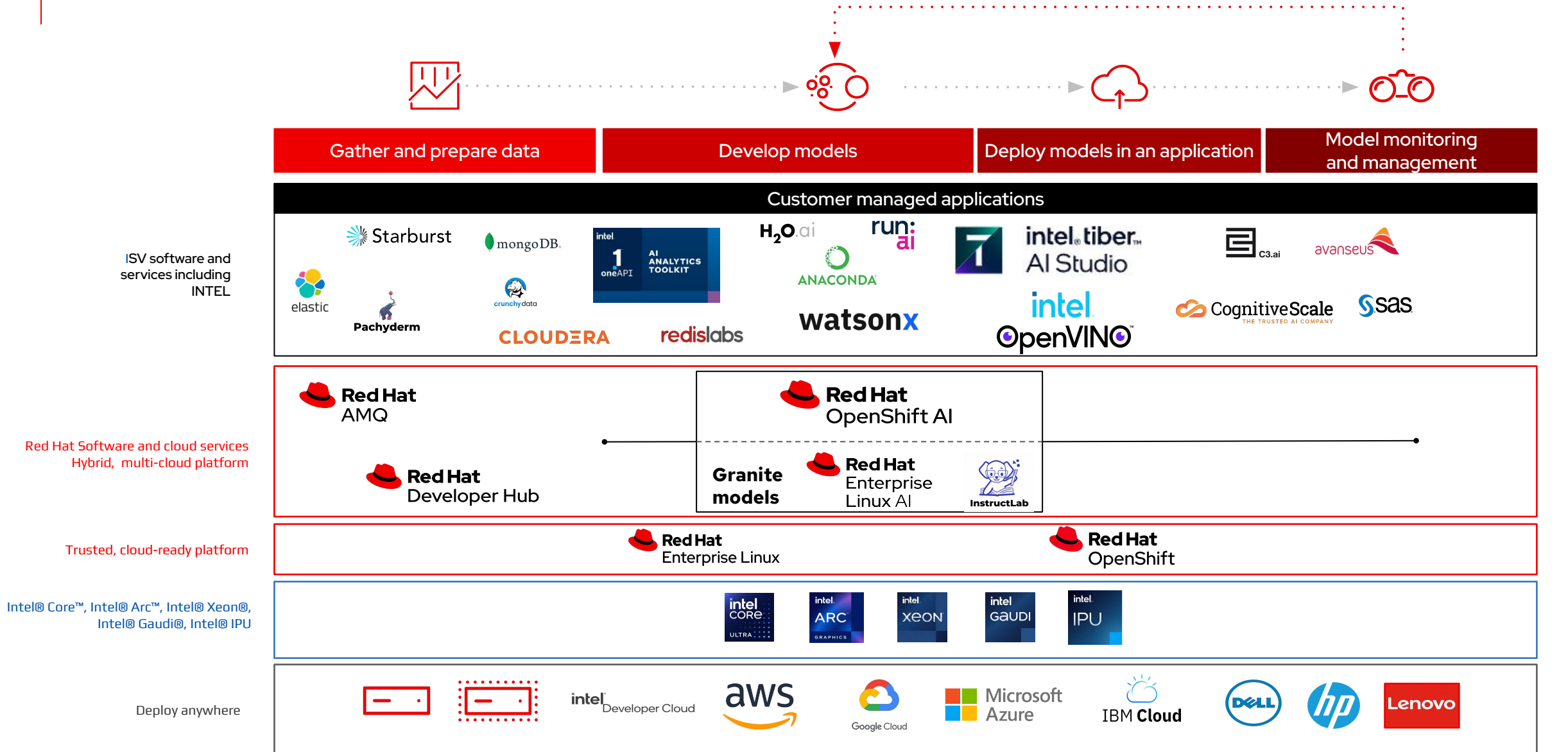
DATA CENTER AI  
AI Open, Scalable Systems & Reference Arch



# Red Hat's AI Strategy



# Intel Enterprise AI with Red Hat® OpenShift® AI



# OPEA – Open Platform for Enterprise AI

# OPEA – Open Platform for Enterprise AI

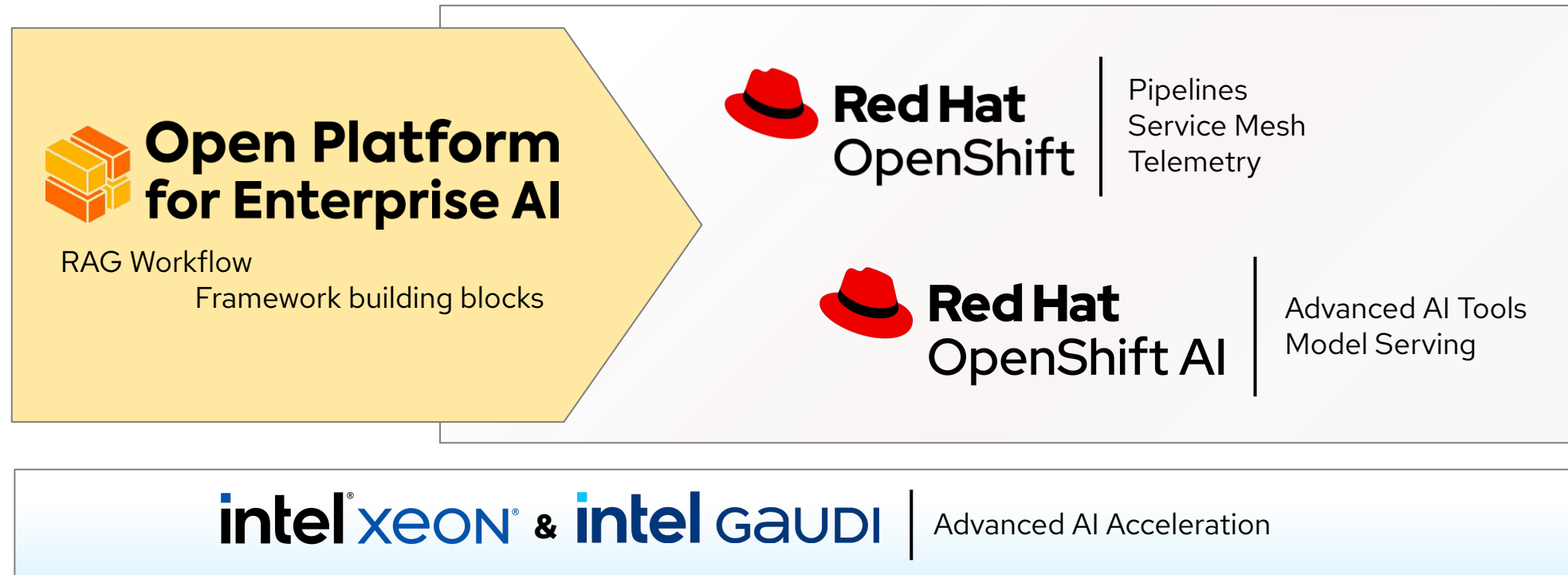
By The Linux Foundation

- ▶ Ecosystem orchestration framework for GenAI
- ▶ OPEA.dev
- ▶ GitHub: <https://github.com/ozea-project>
- ▶ Contributors:



# OPEA with OpenShift AI

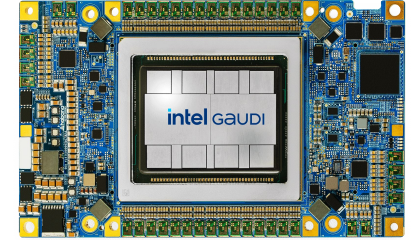
OpenShift AI makes OPEA more enterprise ready



# Intel Gaudi AI Accelerators

# Introducing the Intel® Gaudi® 3 Accelerator

Breaking benchmarks, not budgets



## Competitive Gen AI Performance over H100

- Projected **50% faster time to train**<sup>1</sup>
- Projected **50% faster inferencing**<sup>2</sup>
- Projected **40% better power efficiency**<sup>3</sup>



## Freedom to Scale without Lock-in

- Open standard ethernet networking vs proprietary InfiniBand
- 24x200 GbE ports of industry-standard RoCE on every Gaudi®<sup>3</sup>
- 33% more I/O peak throughput vs H100 for massive scale-up within the server<sup>4</sup>



## Open Development on GenAI platforms

- Integrated open-source PyTorch framework with optimized model library on Hugging Face
- Migrate models on open software from H100 with as few as 3 lines of code

<sup>1</sup> NV H100 comparison based on : <https://developer.nvidia.com/deep-learning-performance-training-inference/training>, Mar 28th 2024 -> "Large Language Model" tab.

<sup>2</sup> Source: NV H100 comparison based on <https://nvidia.github.io/TensorRT-LLM/performance.html#h100-gpus-fp8> , Mar 28th, 2024. Reported numbers are per GPU.

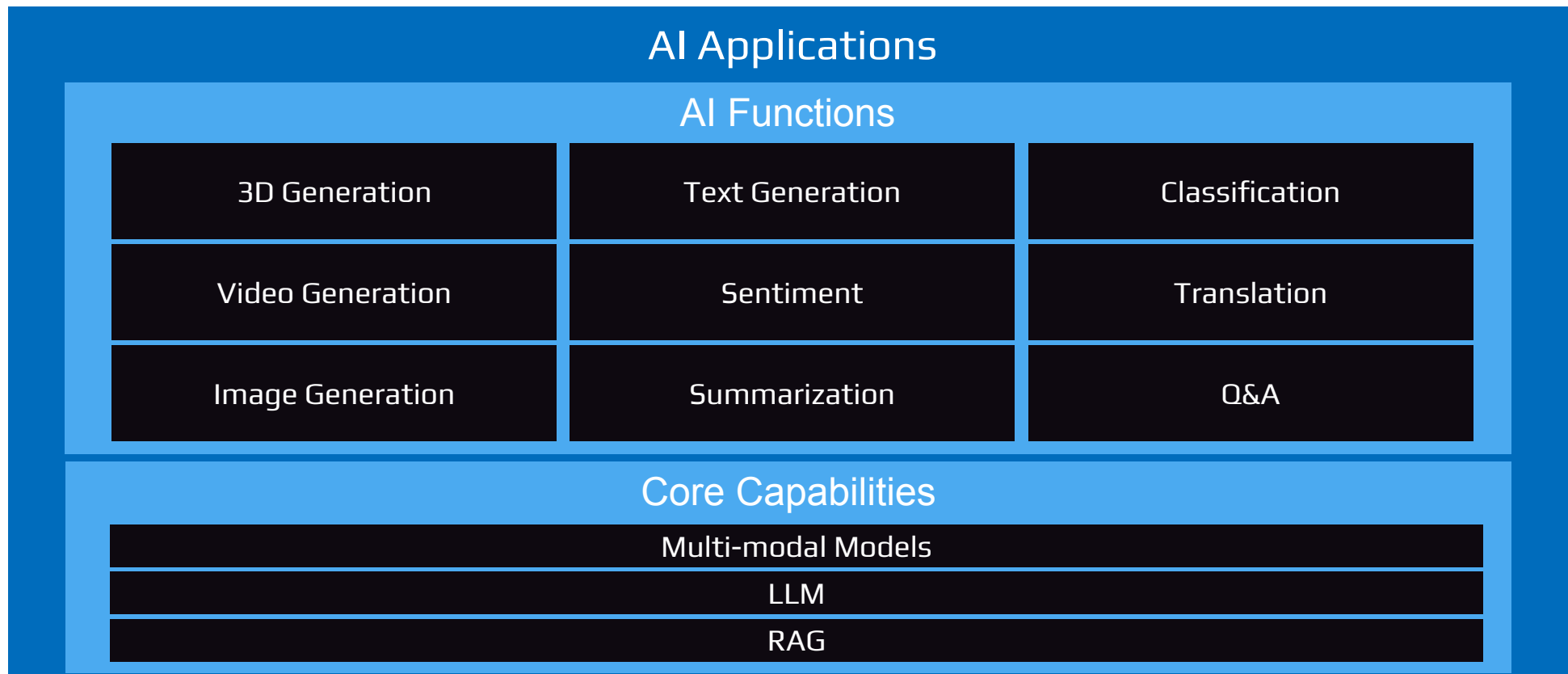
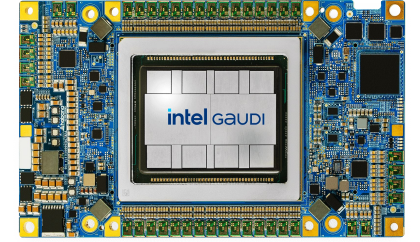
<sup>3</sup> Source: NV comparison based on <https://nvidia.github.io/TensorRT-LLM/performance.html#h100-gpus-fp8> , Mar 28th, 2024. Reported numbers are per GPU.

<sup>1-3</sup> Vs Intel® Gaudi® 3 projections for LLAMA2-7B, LLAMA2-70B & Falcon 180B Power efficiency for both Nvidia and Gaudi3 based on internal estimates. Results may vary.

<sup>4</sup> 900 GB/s NVLink connectivity on H100 vs. 1200 GB/s on Gaudi 3

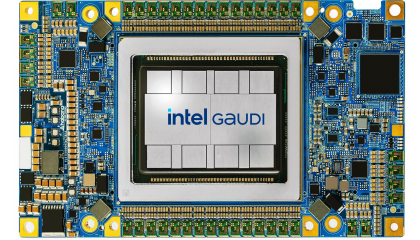
# Intel Gaudi AI Accelerators

Broad Application Support with Focus on Multi-Modal, LLM and RAG



# Intel® Gaudi® 3 AI Accelerator

Launch Partners

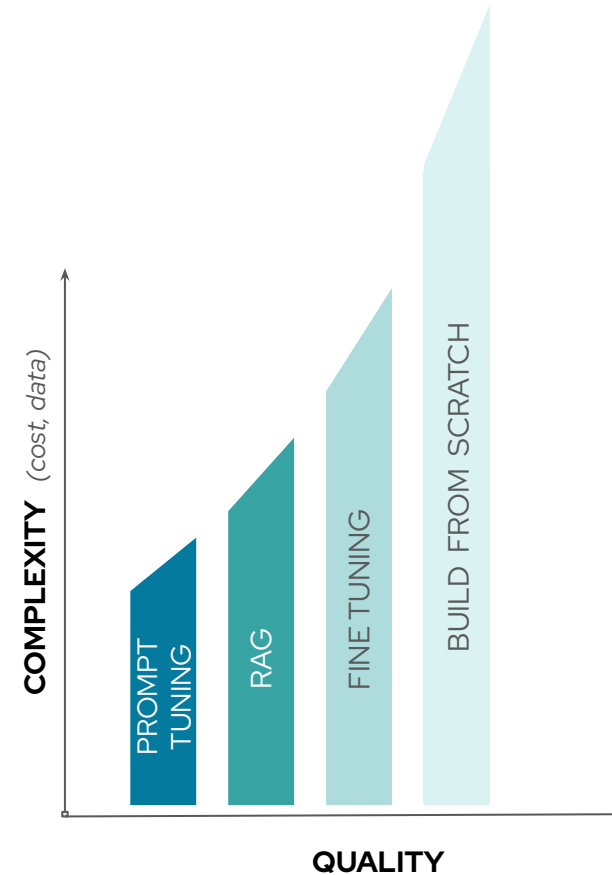


# Retrieval Augmented Generation (RAG) Explained

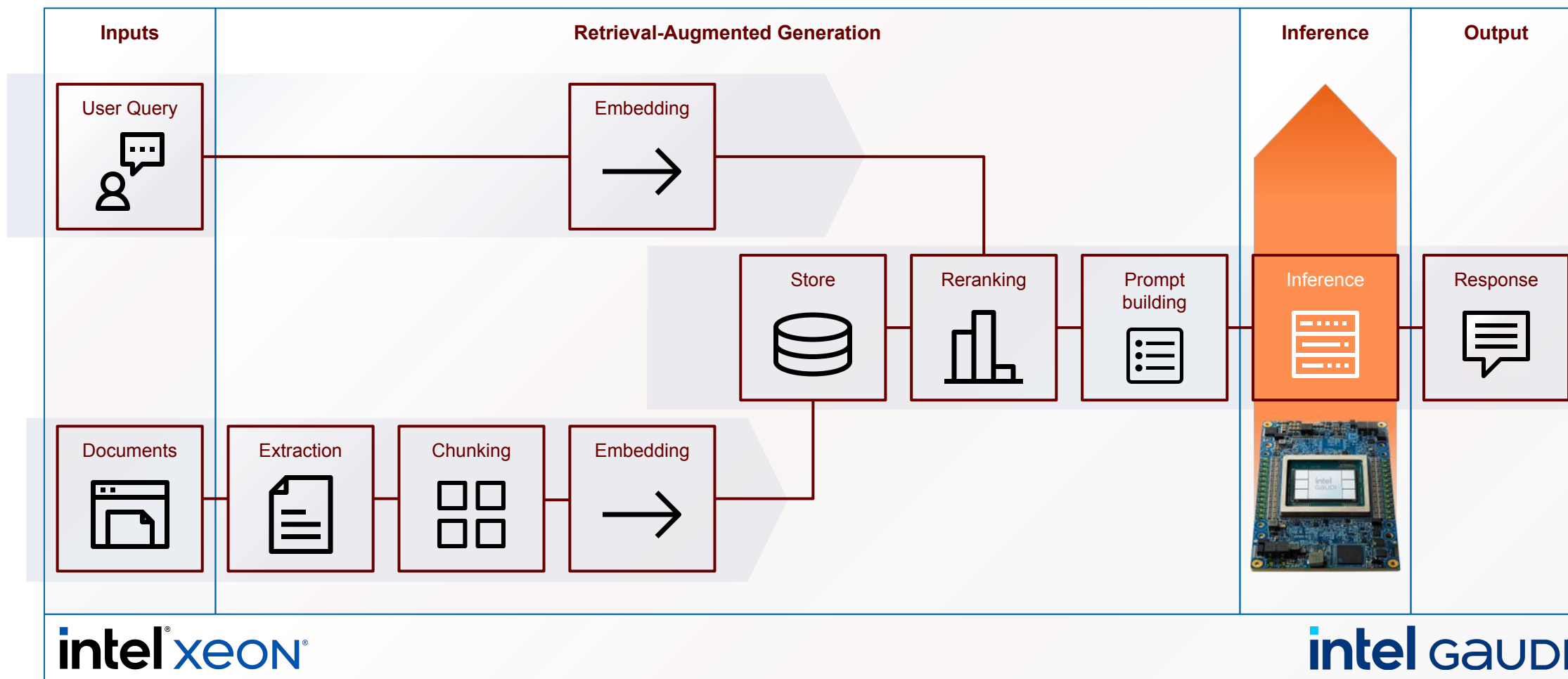
# The balancing act of using foundation models

Foundation models will still need more work to be useful

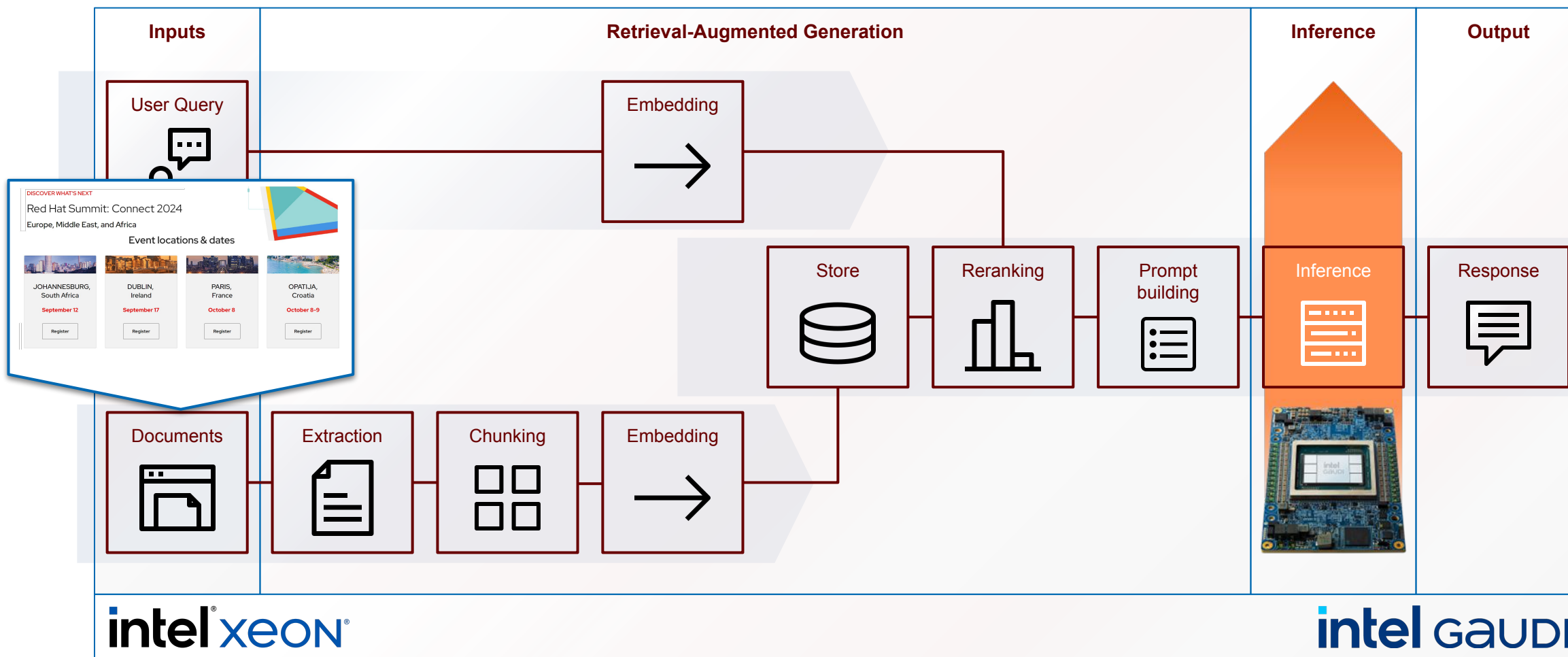
- ▶ Prompt tuning
- ▶ Retrieval-Augmented Generation (RAG)
- ▶ Fine tuning foundation models
- ▶ Training a Foundation Model from scratch



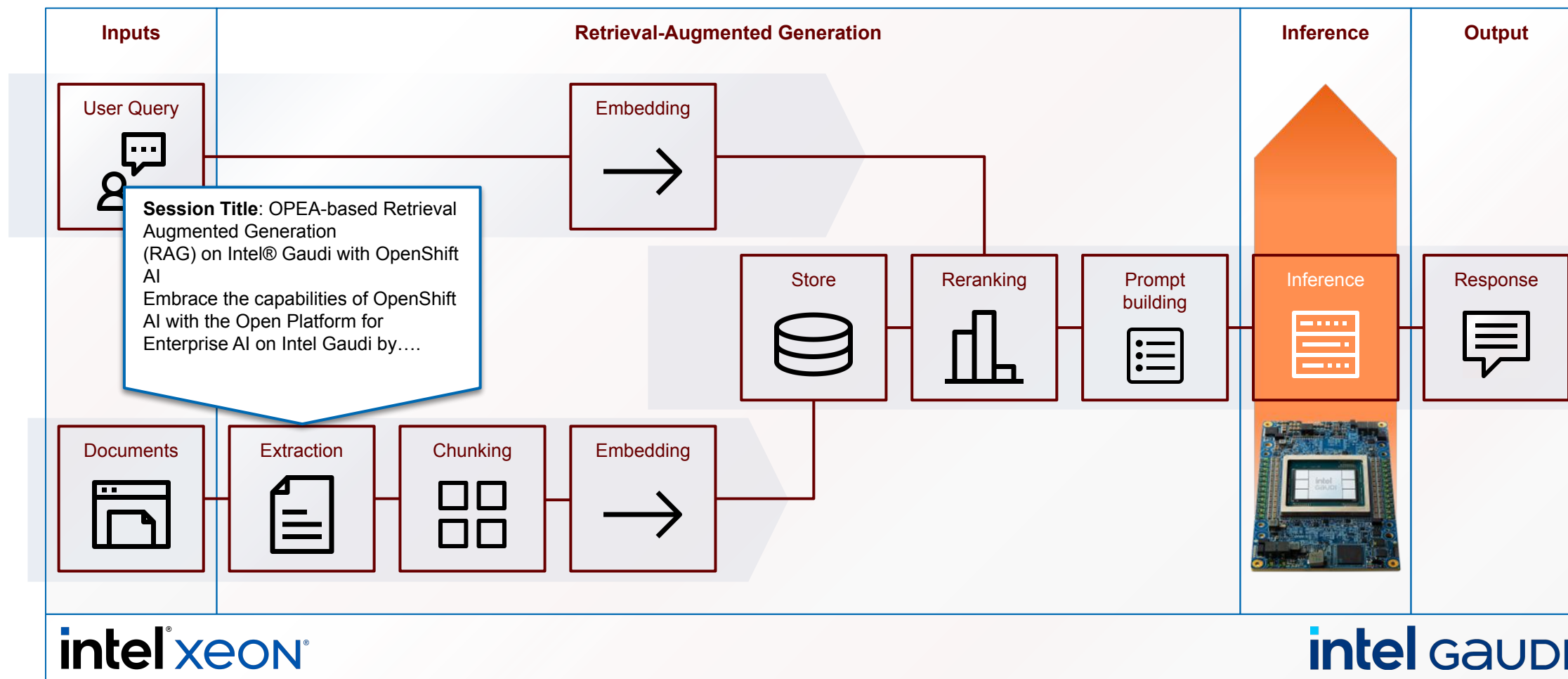
# Retrieval Augmented Generation (RAG)



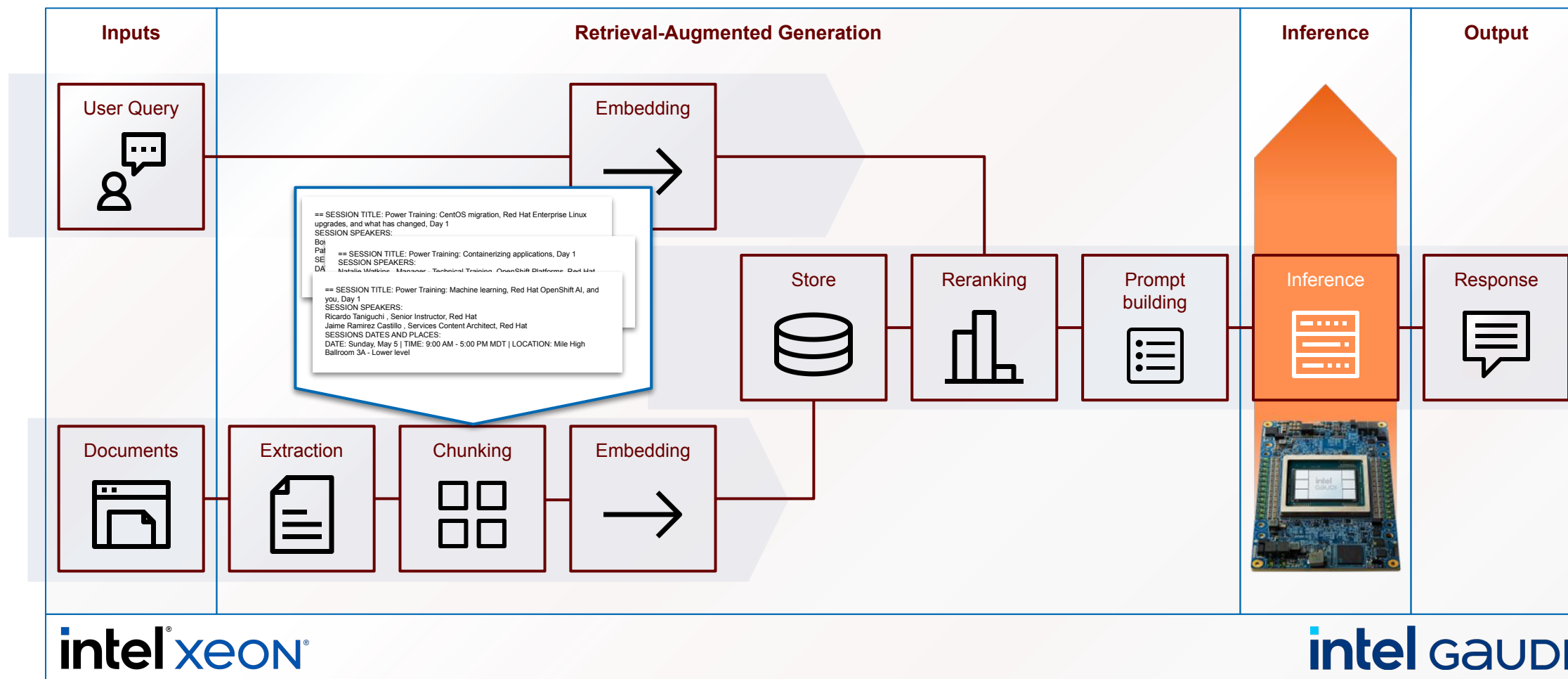
# Retrieval Augmented Generation (RAG)



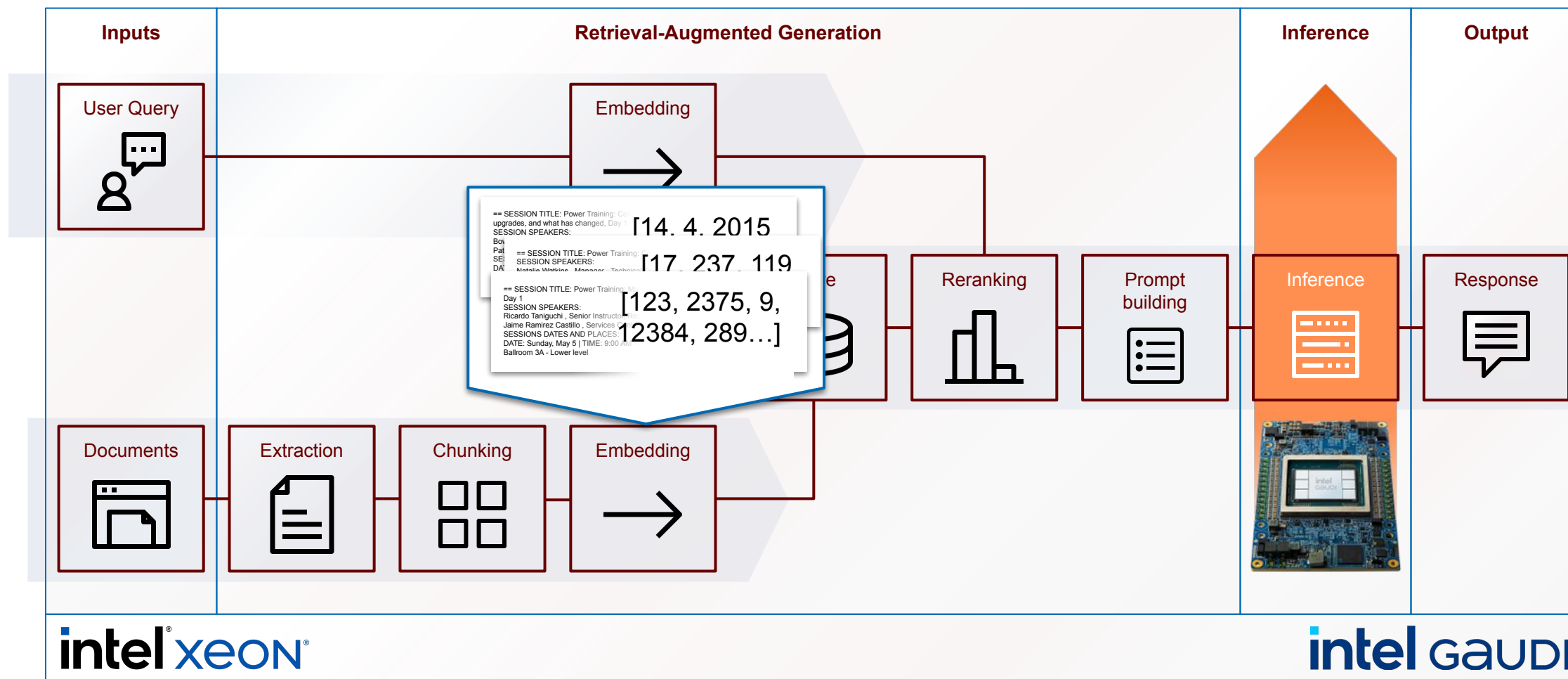
# Retrieval Augmented Generation (RAG)



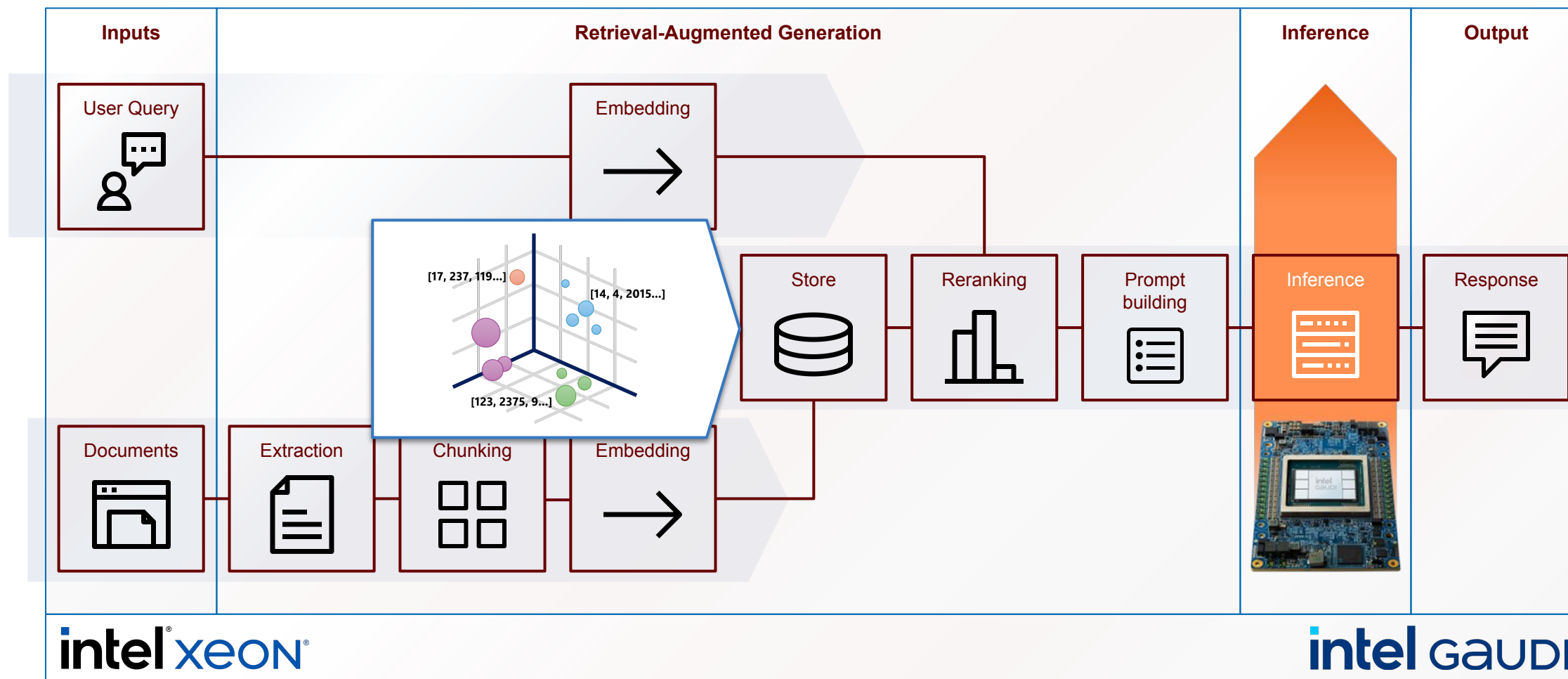
# Retrieval Augmented Generation (RAG)



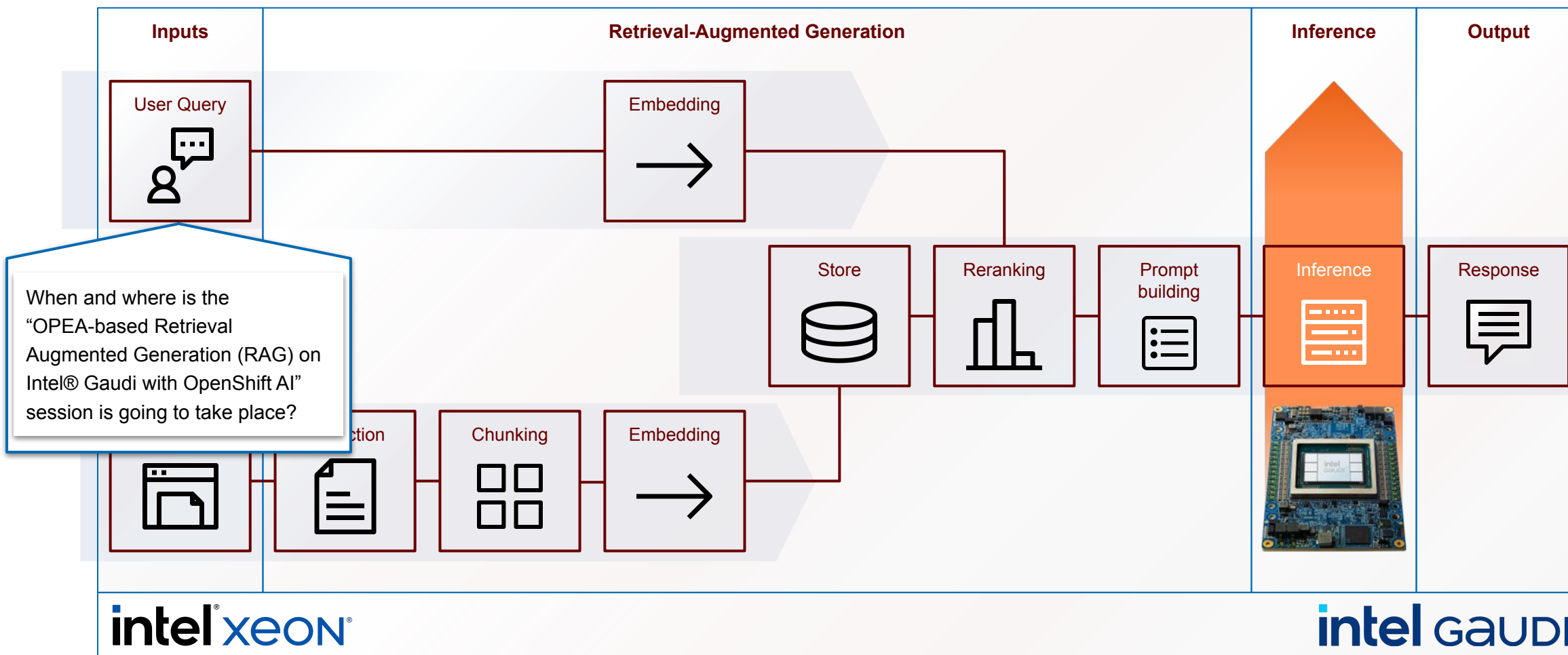
# Retrieval Augmented Generation (RAG)



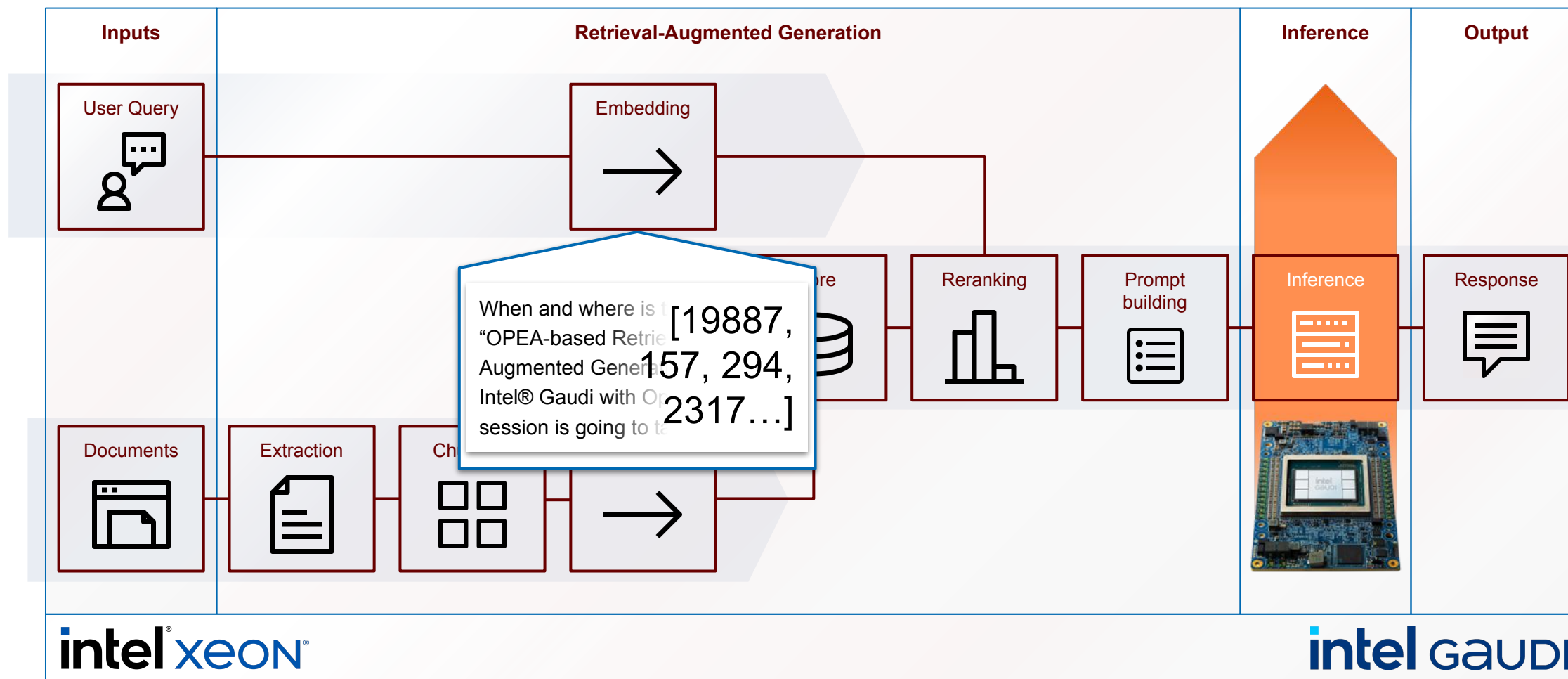
# Retrieval Augmented Generation (RAG)



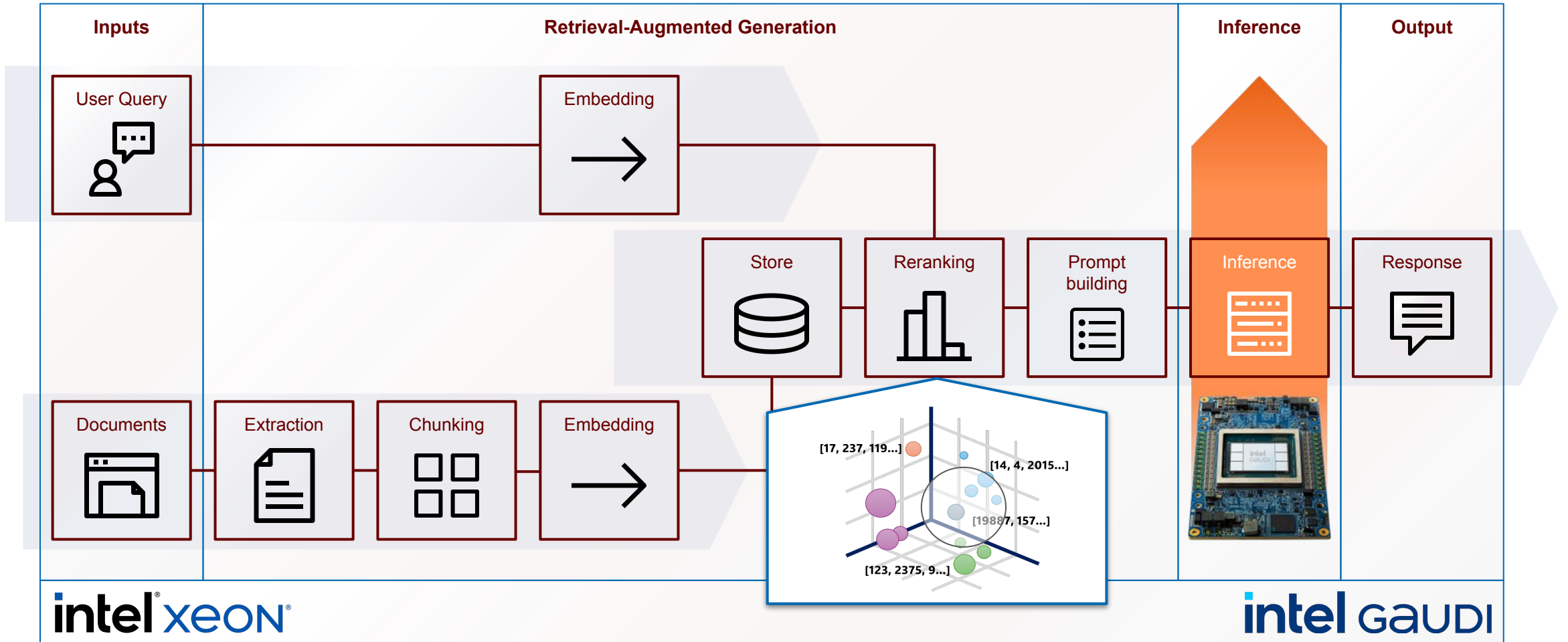
# Retrieval Augmented Generation (RAG)



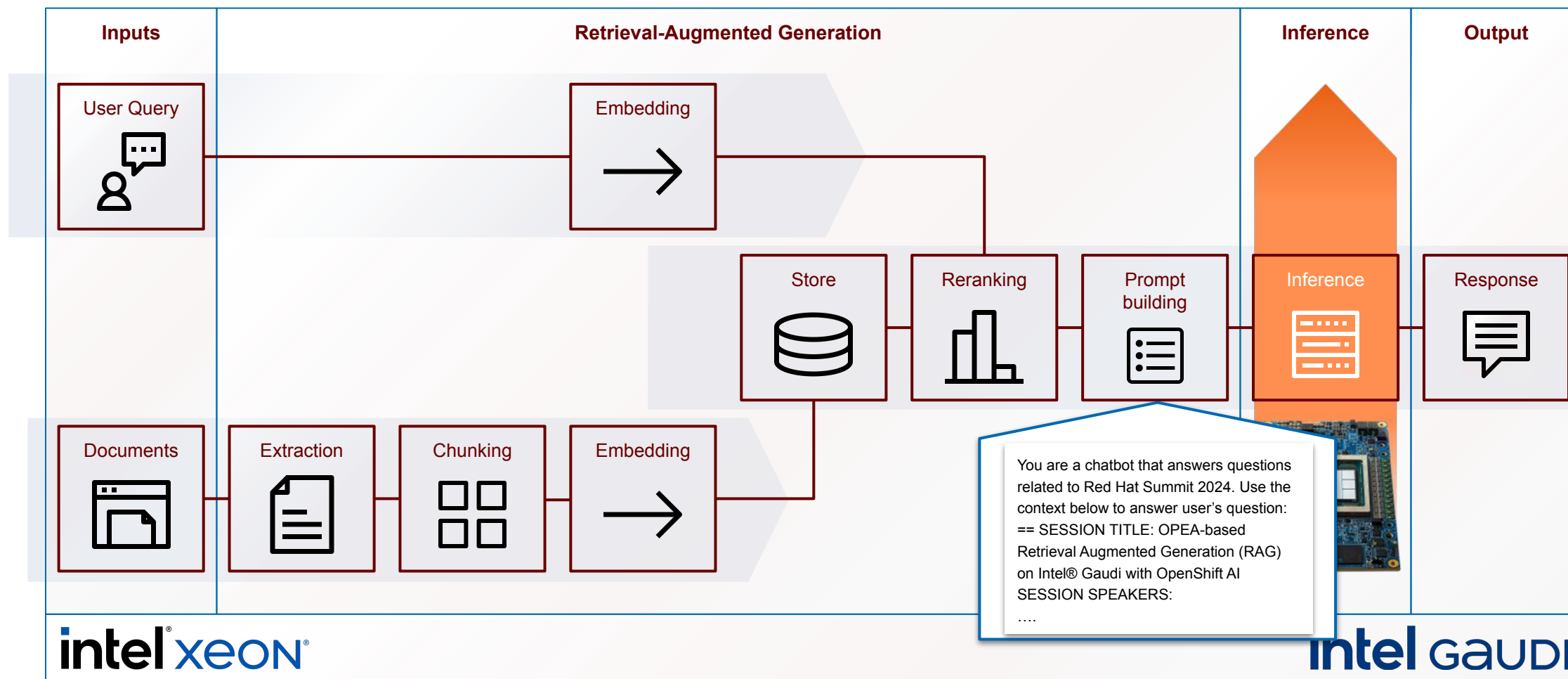
# Retrieval Augmented Generation (RAG)



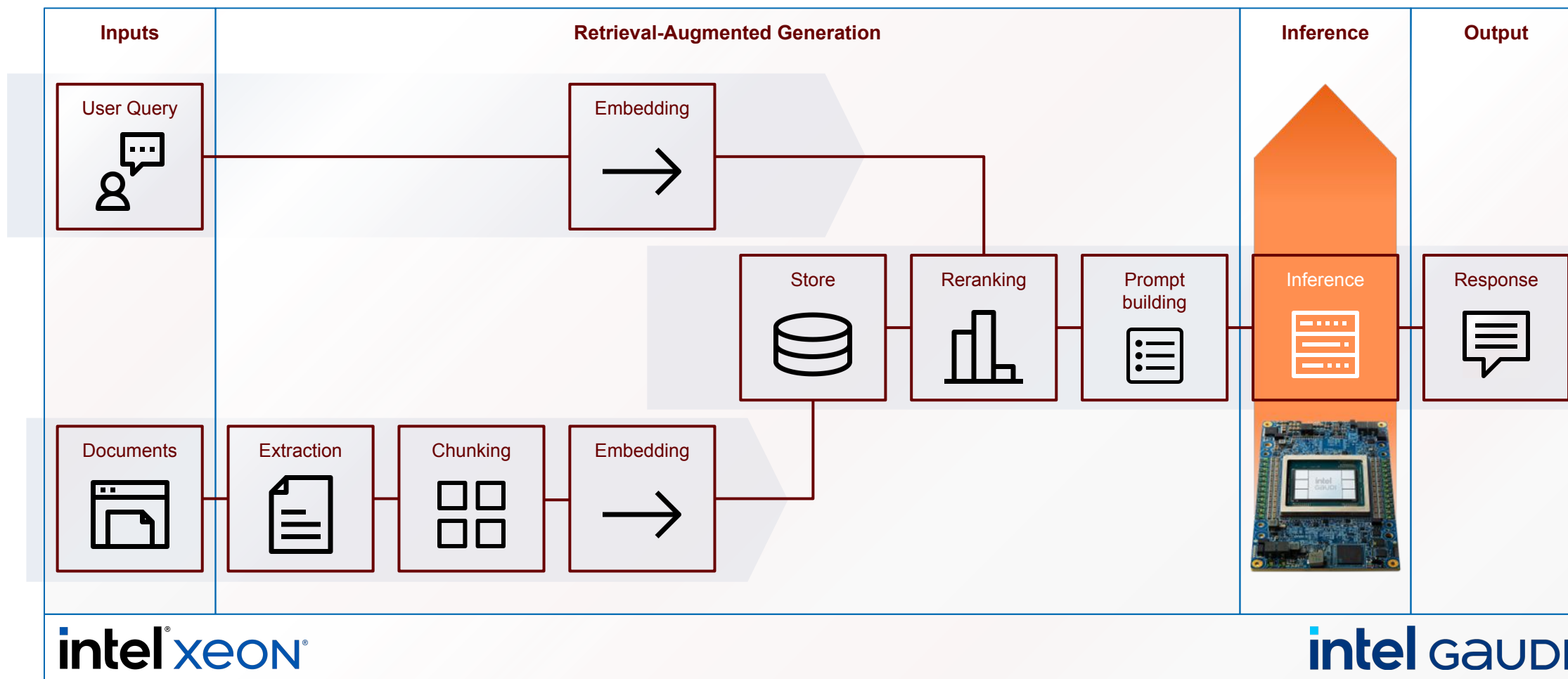
# Retrieval Augmented Generation (RAG)



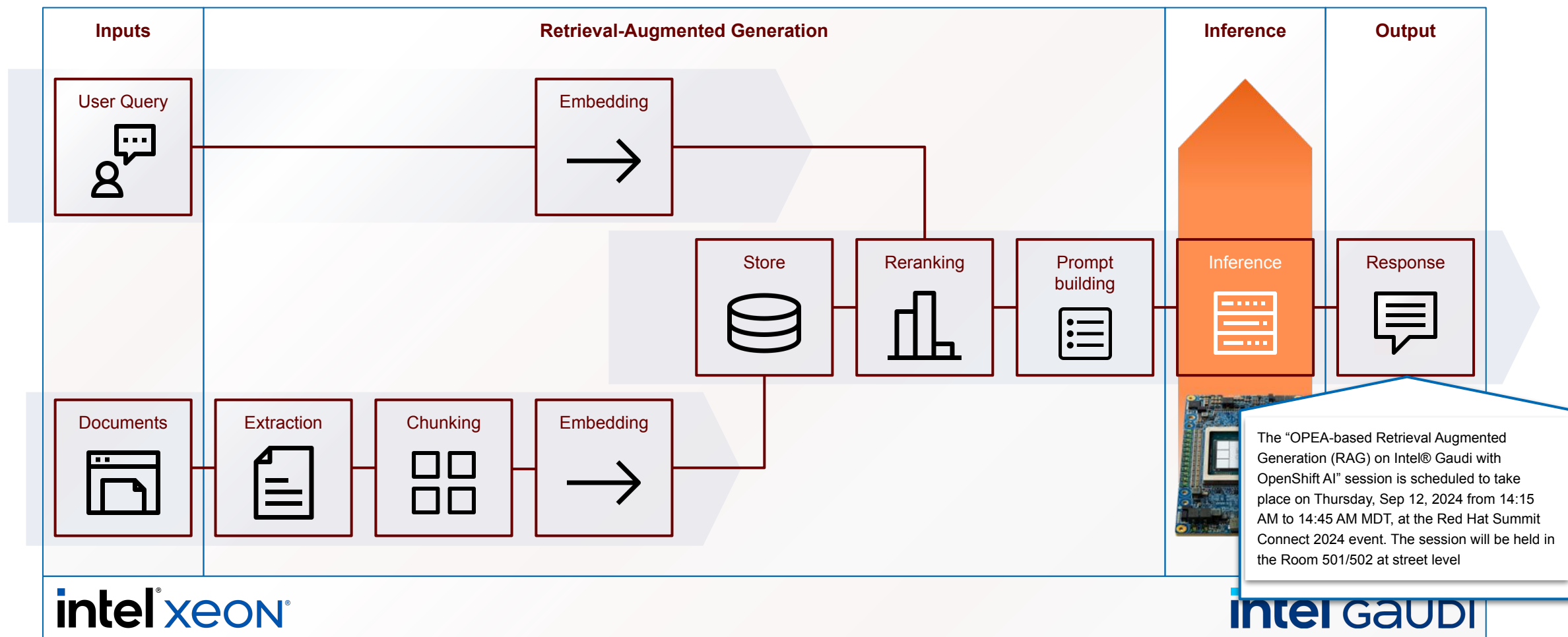
# Retrieval Augmented Generation (RAG)



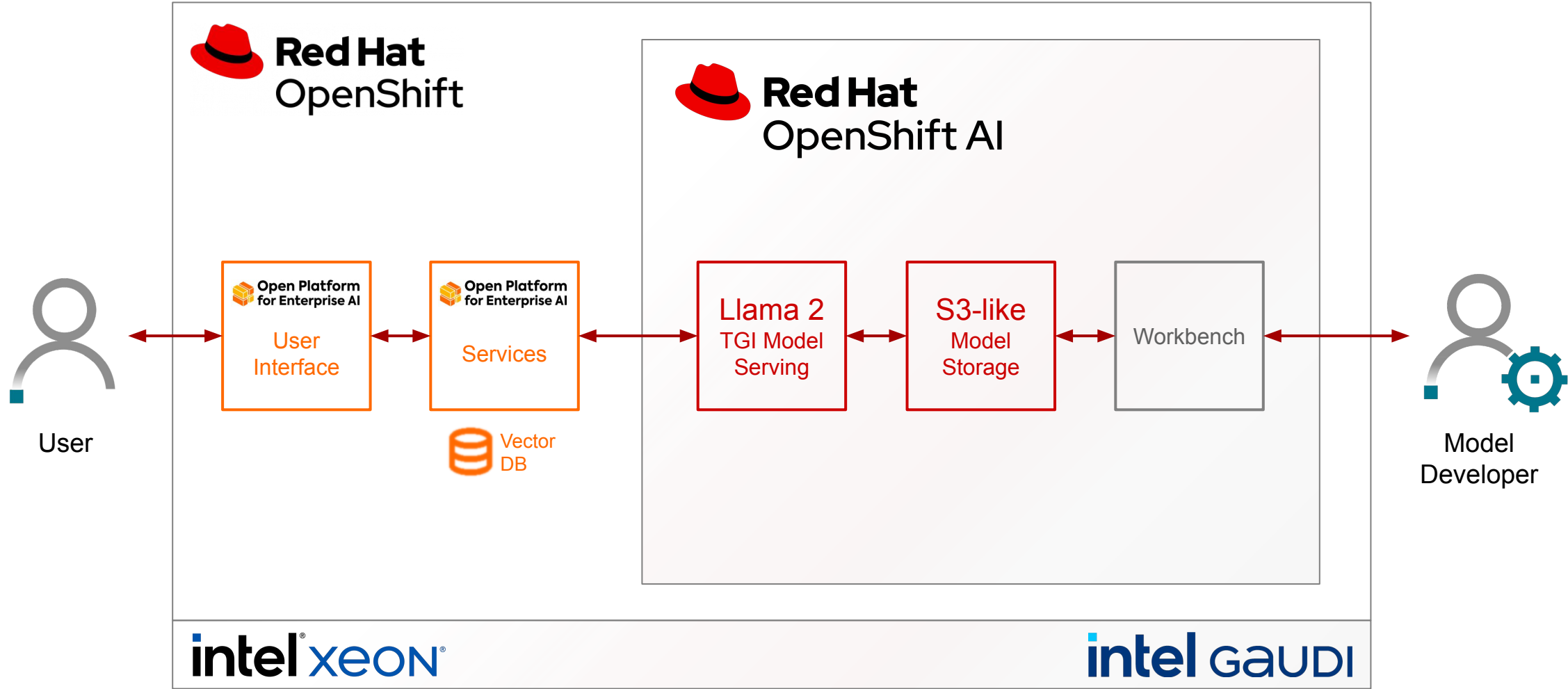
# Retrieval Augmented Generation (RAG)



# Retrieval Augmented Generation (RAG)



# Retrieval Augmented Generation (RAG) Chatbot Demo



Red Hat

OpenShift

Project: All Projects

Administrator

Home

Operators

OperatorHub

Installed Operators

Workloads

Serverless

Networking

Storage

Builds

Observe

Compute









User Management

Administration

Installed Operators

Name

Search by name...

Name	Namespace	Managed Namespaces	Status	Last updated	Provided APIs
<div></div> <div><div>Habana AI</div><div>115.0-479 provided by Habana Labs Ltd.</div></div>	<div>NS</div> habana-ai-operator	<div>NS</div> habana-ai-operator	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 30, 2024, 6:10 PM	<a href="#">Device Config</a> <div></div>
<div></div> <div><div>Kernel Module Management</div><div>21.0 provided by Red Hat</div></div>	<div>NS</div> openshift-kmm	All Namespaces	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 30, 2024, 11:54 AM	<a href="#">PreflightValidation</a> <a href="#">PreflightValidationOCP</a> <a href="#">Module</a> <a href="#">NodeModulesConfig</a> <div></div>
<div></div> <div><div>LVM Storage</div><div>4.14.4 provided by Red Hat</div></div>	<div>NS</div> openshift-storage	<div>NS</div> openshift-storage	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 29, 2024, 3:17 PM	<a href="#">LVMCluster</a> <div></div>
<div></div> <div><div>Node Feature Discovery Operator</div><div>4.14.0-202404161544 provided by Red Hat</div></div>	<div>NS</div> openshift-nfd	<div>NS</div> openshift-nfd	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 30, 2024, 6:10 PM	<a href="#">NodeFeatureDiscovery</a> <a href="#">NodeFeatureRule</a> <div></div>
<div></div> <div><div>Package Server</div><div>0.0.1-snapshot provided by Red Hat</div></div>	<div>NS</div> openshift-operator-lifecycle-manager	<div>NS</div> openshift-operator-lifecycle-manager	<div>✓</div> Succeeded	<div>🕒</div> Apr 29, 2024, 3:17 PM	<a href="#">PackageManifest</a> <div></div>
<div></div> <div><div>Red Hat OpenShift AI</div><div>2.8.1 provided by Red Hat</div></div>	<div>NS</div> redhat-ods-operator	All Namespaces	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 29, 2024, 3:17 PM	<a href="#">Data Science Cluster</a> <a href="#">DSC Initialization</a> <a href="#">FeatureTracker</a> <div></div>
<div></div> <div><div>Red Hat OpenShift Serverless</div><div>1.32.1 provided by Red Hat</div></div>	<div>NS</div> openshift-serverless	All Namespaces	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 29, 2024, 3:18 PM	<a href="#">Knative Serving</a> <a href="#">Knative Eventing</a> <a href="#">Knative Kafka</a> <div></div>
<div></div> <div><div>Red Hat OpenShift Service Mesh</div><div>2.5.1-0 provided by Red Hat, Inc.</div></div>	<div>NS</div> openshift-operators	All Namespaces	<div>✓</div> Succeeded Up to date	<div>🕒</div> Apr 30, 2024, 11:54 AM	<a href="#">Istio Service Mesh Control Plane</a> <a href="#">Istio Service Mesh Member</a> <a href="#">Istio Service Mesh Member Roll</a> <div></div>

5

kube:admin

- Administrator
- Home
- Operators
  - OperatorHub
  - Installed Operators
- Workloads
- Serverless
- Networking
- Storage
- Builds
- Observe
- Compute
- User Management
- Administration

Project: All Projects

## Installed Operators

Name
Search by name...

Name	Namespace	Managed Namespaces	Status	Last updated	Provided APIs
<b>Habana AI</b> 115.0-479 provided by Habana Labs Ltd.	NS habana-ai-operator	NS habana-ai-operator	Succeeded Up to date	Apr 30, 2024, 6:10 PM	<a href="#">Device Config</a>
<b>Kernel Module Management</b> 21.0 provided by Red Hat	NS openshift-kmm	All Namespaces	Succeeded Up to date	Apr 30, 2024, 11:54 AM	<a href="#">PreflightValidation</a> <a href="#">PreflightValidationOCPModule</a> <a href="#">NodeModulesConfig</a>
<b>LVM Storage</b> 4.14.4 provided by Red Hat	NS openshift-storage	NS openshift-storage	Succeeded Up to date	Apr 29, 2024, 3:17 PM	<a href="#">LVMCluster</a>
<b>Node Feature Discovery Operator</b> 4.14.0-202404161544 provided by Red Hat	NS openshift-nfd	NS openshift-nfd	Succeeded Up to date	Apr 30, 2024, 6:10 PM	<a href="#">NodeFeatureDiscovery</a> <a href="#">NodeFeatureRule</a>
<b>Package Server</b> 0.01-snapshot provided by Red Hat	NS openshift-operator-lifecycle-manager	NS openshift-operator-lifecycle-manager	Succeeded Up to date	Apr 29, 2024, 3:17 PM	<a href="#">PackageManifest</a>
<b>Red Hat OpenShift AI</b> 2.8.1 provided by Red Hat	NS redhat-ods-operator	All Namespaces	Succeeded Up to date	Apr 29, 2024, 3:17 PM	<a href="#">Data Science Cluster</a> <a href="#">DSC Initialization</a> <a href="#">FeatureTracker</a>
<b>Red Hat OpenShift Serverless</b> 1.32.1 provided by Red Hat	NS openshift-serverless	All Namespaces	Succeeded Up to date	Apr 29, 2024, 3:18 PM	<a href="#">Knative Serving</a> <a href="#">Knative Eventing</a> <a href="#">Knative Kafka</a>
<b>Red Hat OpenShift Service Mesh</b> 2.5.1-0 provided by Red Hat, Inc.	NS openshift-operators	All Namespaces	Succeeded Up to date	Apr 30, 2024, 11:54 AM	<a href="#">Istio Service Mesh Control Plane</a> <a href="#">Istio Service Mesh Member</a> <a href="#">Istio Service Mesh Member Roll</a>

Red Hat

OpenShift

Project: All Projects

Administrator

Home

Operators

OperatorHub

Installed Operators

Workloads

Serverless

Networking

Storage

Builds

Observe

Compute









User Management

Administration

Installed Operators

Name

Search by name...

Name	Namespace	Managed Namespaces	Status	Last updated	Provided APIs
 <div><b>Habana AI</b> 115.0-479 provided by Habana Labs Ltd.</div>	NS habana-ai-operator	NS habana-ai-operator	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 6:10 PM	Device Config
 <div><b>Kernel Module Management</b> 2.1.0 provided by Red Hat</div>	NS openshift-kmm	All Namespaces	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 11:54 AM	PreflightValidation PreflightValidationOCP Module NodeModulesConfig
 <div><b>LVM Storage</b> 4.14.4 provided by Red Hat</div>	NS openshift-storage	NS openshift-storage	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 29, 2024, 3:17 PM	LVMCluster
 <div><b>Node Feature Discovery Operator</b> 4.14.0-202404161544 provided by Red Hat</div>	NS openshift-nfd	NS openshift-nfd	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 6:10 PM	NodeFeatureDiscovery NodeFeatureRule
 <div><b>Package Server</b> 0.0.1-snapshot provided by Red Hat</div>	NS openshift-operator-lifecycle-manager	NS openshift-operator-lifecycle-manager	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 6:10 PM	
 <div><b>Red Hat OpenShift AI</b> 2.8.1 provided by Red Hat</div>	NS redhat-ods-operator	All Namespaces	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 6:10 PM	
 <div><b>Red Hat OpenShift Serverless</b> 1.32.1 provided by Red Hat</div>	NS openshift-serverless	All Namespaces	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 6:10 PM	
 <div><b>Red Hat OpenShift Service Mesh</b> 2.5.1-0 provided by Red Hat, Inc.</div>	NS openshift-operators	All Namespaces	<div>✓ Succeeded</div> <div>Up to date</div>	🕒 Apr 30, 2024, 11:54 AM	Istio Service Mesh Control Plane Istio Service Mesh Member Istio Service Mesh Member Roll

Operators are necessary for Gaudi® to run properly on the Red Hat® OpenShift platform.

## Serving runtimes

Manage your model serving runtimes.

Single-model serving enabled

Multi-model serving enabled

Add serving runtime

Name	Enabled	Serving platforms supported	API protocol
<div><div></div>Text Generation Inference on Habana Gaudi</div>	<div><div></div></div>	<div>Single-model</div>	<div>REST</div>
<div><div></div>Caikit TGIS ServingRuntime for KServe</div> <div>Pre-installed</div>	<div><div></div></div>	<div>Single-model</div>	<div>REST</div>
<div><div></div>OpenVINO Model Server</div> <div>Pre-installed</div>	<div><div></div></div>	<div>Single-model</div>	<div>REST</div>
<div><div></div>OpenVINO Model Server</div> <div>Pre-installed</div>			
<div><div></div>TGIS Standalone ServingRuntime for KServe</div> <div>Pre-installed</div>			

To accelerate your OpenShift AI model with Intel® Gaudi® 2, you need a suitable Serving runtime

- Applications >
- Data Science Projects
- Data Science Pipelines
- Model Serving
- Resources
- Settings ▾
  - Notebook images
  - Cluster settings
  - Accelerator profiles
  - Serving runtimes
  - User management



# Accelerator profiles

Manage accelerator profile settings for users in your organization

Name ▾

 Find by name

Create accelerator profile

Name ↑	Identifier ↑ ⓘ ⓘ	Enable ⓘ ⓘ	Last modified ↑	
Gaudi 2	habana.ai/gaudi		9 days ago	

and an adequate  
Accelerator profile.

- Applications >
- Data Science Projects
- Data Science Pipelines
- Model Serving
- Resources
- Settings ▾
- Notebook images
- Cluster settings
- Accelerator profiles
- Serving runtimes
- User management

Data Science Projects > Red Hat Summit LLM w/ RAG Demo

# Red Hat Summit LLM w/ RAG Demo

Components

Permissions

Jump to section

Workbenches

Cluster storage

Data connections

Models and model servers

## Workbenches

Create workbench

Name	Notebook image	Container size	Status	
> wb1	HabanaAI	Small	<div><div></div>Running</div>	<div>Open</div>

## Cluster storage

Add cluster storage

Name	Type	Connected workbenches	
> wb1	<div>Persistent storage</div>	wb1	

## Data connections

Add data connection

Name	Type	Connected workbenches	
bucket-70b-chat	<div>Object storage</div>	No connections	
minio	<div>Object storage</div>	No connections	

## Models and model servers

Deploy model

Single-model serving enabled

Model name	Serving runtime	Inference endpoint	API protocol	Status	
> tgi-70b-1	Text Generation Inference on Habana Gaudi	<div>https://tgi-70b-1-rag-on-gaudi.apps.cluster1.gaudi.internal</div>	REST	<div></div>	
> tgi-70b-2	Text Generation Inference on Habana Gaudi	<div>https://tgi-70b-2-rag-on-gaudi.apps.cluster1.gaudi.internal</div>	REST	<div></div>	<div>Edit</div> <div>Delete</div>

What's on your mind?



Stock Llama2 70b

Stock Llama2 70b with RAG



redis



What is Red Hat Summit?



## Stock Llama2 70b

Hi!

Hello! What can I help you with today?



## Stock Llama2 70b with RAG

Hi!

Hi! Welcome to Red Hat Summit 2024! We're excited to have you join us for this year's event. Can we help you with something specific? Are you looking for information on a particular topic or session? Or perhaps you'd like to know more about our keynote speakers or networking opportunities? Let us know, and we'll do our best to assist you. See you soon!

### Sources



Red Hat Summit 2024 s...



# Summary

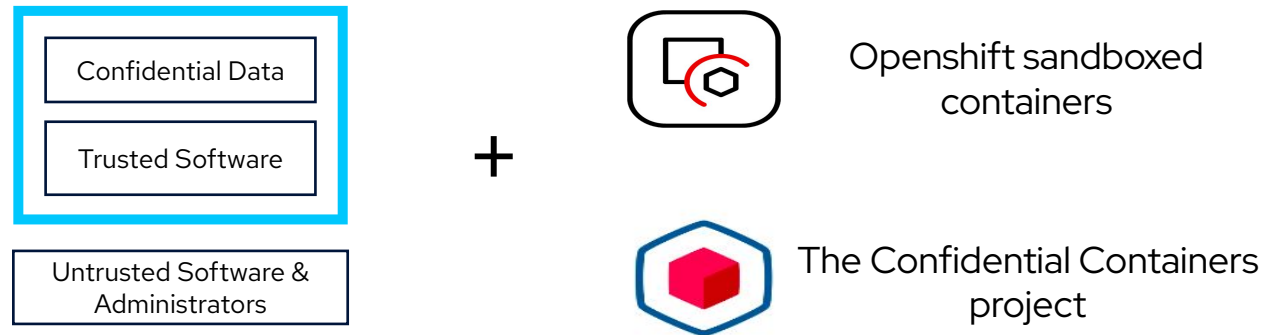
# Key Takeaways

- ▶ RAG enhances AI development
- ▶ OPEA simplifies AI deployment
- ▶ OpenShift AI integrates into DevOps workflow
- ▶ Intel Gaudi 3 accelerates AI training and inference

# Confidential AI Helps Protect Data & Models In-Use

## Utilizing Confidential Computing for Containers with Intel TDX

Hardware-Based Protection of Data In-Use  
With Intel Trusted Domain Extensions (TDX)



Confidential Computing is about **protecting data in-use**.  
You do not **have to trust** the system admins of the providers any longer.

# Q&A

Red Hat  
**Summit**

**Connect**

# Thank you



[linkedin.com/company/red-hat](https://linkedin.com/company/red-hat)



[facebook.com/redhatinc](https://facebook.com/redhatinc)



[youtube.com/user/RedHatVideos](https://youtube.com/user/RedHatVideos)



[twitter.com/RedHat](https://twitter.com/RedHat)

## CODRIN BUCUR

Principal AI Specialist Solution Architect  
Red Hat EMEA



**Bio:** As an Principal AI Specialist Solution Architect, Codrin is supporting Red Hat customers and partners in EMEA with their data science, AI/ML and MLOps needs and best practices. Previously, as Architect and TSM in Red Hat Consulting Alps for 7+ years, Codrin has supported customers with their adoption of Red Hat container platform, integration and middleware technologies.

**Contact:** [cbucur@redhat.com](mailto:cbucur@redhat.com)

<https://www.linkedin.com/in/codrin>

