



# Dell AI Factory

## Como a Mágica Acontece

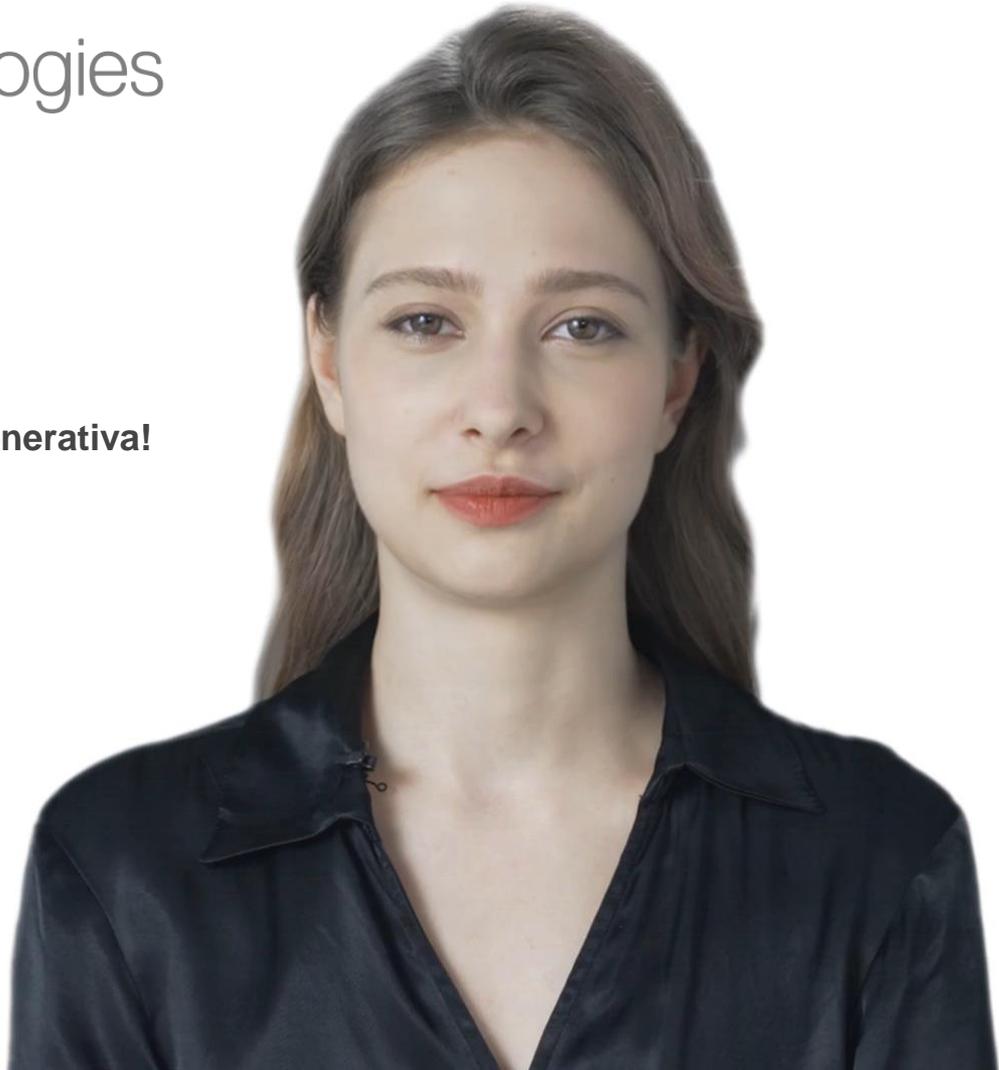
With NVIDIA & INTEL



Rogério Chola  
HPC&AI Solution Architect

**DELL**Technologies

**DELL** Technologies



**Bem-vindos a Era da IA Generativa!**



# Pilares da Tecnologia: Prontidão para o que vem a seguir

Principais tecnologias emergentes no mundo atual de dados em qualquer lugar

**5G**

**Segurança  
Intrinseca**

**Digital Twin  
Industria 4.0**

**IA Generativa  
IA Cognitiva**

**Computação  
Quântica**

**Edge  
Multicloud**

**Futuro do Trabalho**

# Generative AI

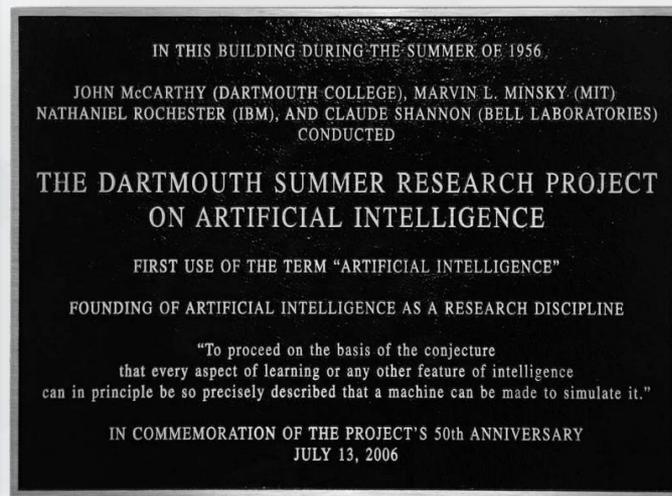


**DELL** Technologies

# Introdução à IA - Nascimento e Crescimento



John McCarthy @ Dartmouth College - Hanover, NH - 1956



# Generative AI

**Generative AI** refere-se a técnicas de inteligência artificial que aprendem uma representação de artefatos a partir de dados e os usam para gerar artefatos novos, completamente originais em escala, que preservam uma semelhança com os dados originais.

**Generative AI** permite que os computadores gerem variações de conteúdo totalmente originais (incluindo imagens, vídeo, música, fala e texto). Ele pode melhorar ou alterar o conteúdo existente, e pode criar novos elementos de dados e novos modelos de objetos do mundo real, como edifícios, peças, medicamentos e materiais.

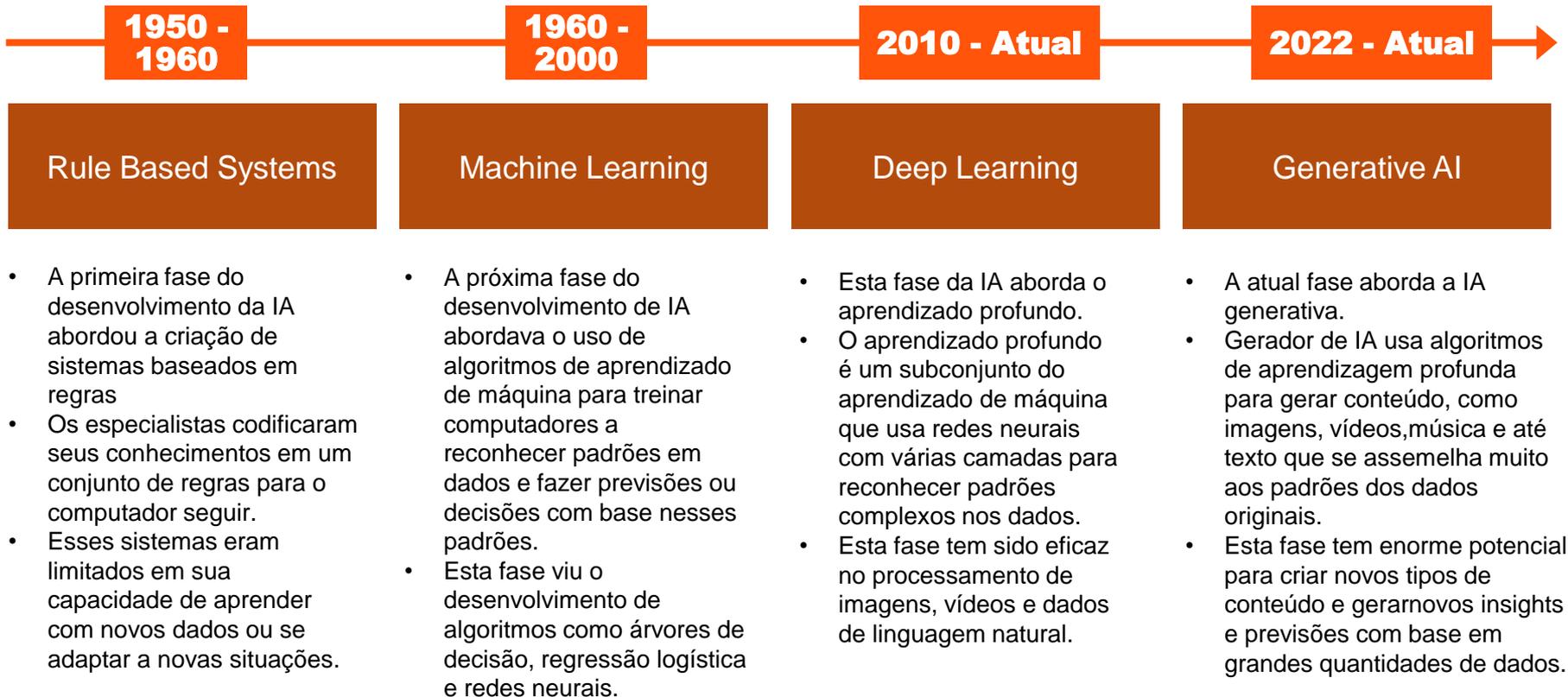


## IA Tradicional vs IA Generativa

A IA Tradicional visa resolver problemas específicos analisando dados existentes e fazendo previsões, enquanto a IA Generativa cria conteúdo novo e original com base em padrões aprendidos. A IA Tradicional é adequada para tarefas como classificação, recomendação e tomada de decisão, enquanto a IA Generativa libera o potencial para exploração criativa e geração de conteúdo.

# AI TIMELINE

Desde a metade do século 20 estamos avançando no uso de IA.



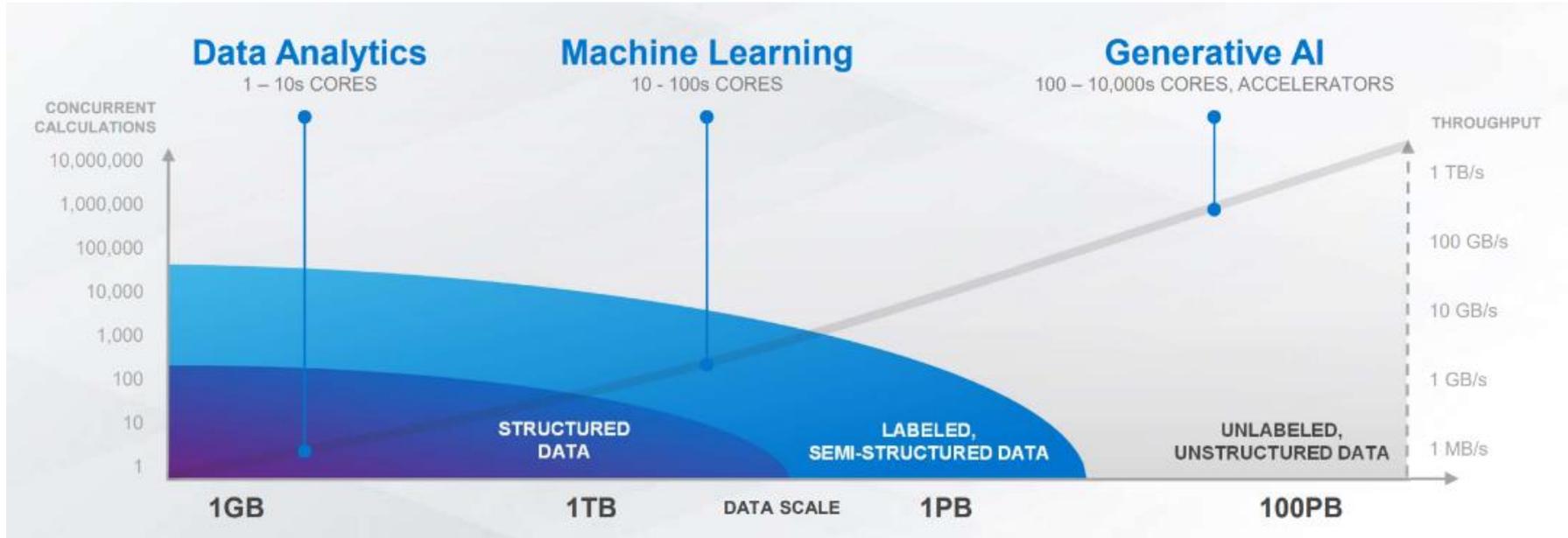
# AI TIMELINE

Uma série de avanços cada vez mais frequentes que dão sentido à linguagem natural



# Generative AI Evolution

Dell Technologies & Nvidia provides expertise and guidance to build, train, and model-tune Generative AI.



# Absolutamente todas as pessoas e negócios do mundo serão profundamente impactados pelo uso da Inteligência Artificial.

Agência de modelos v

A Deep Agency p  
utilizada para tira

Inteligên  
predição  
câncer

Estudo na Faculdade de Saúde Pública da USP revela benefícios na utilização de inteligência artificial para auxiliar tomadas de decisões médicas

Inovação

**Audi usa inteligência artificial para produzir rodas de carro; entenda como funciona**

Programa de IA é combinado por dois sistemas que competem para criar uma imagem do zero; montadora quer levar tecnologia para outras áreas

*"Artificial intelligence is the most important advancement of technology in decades. That will change the way how people work learn, travel, get medical assistance and communicate with each other others" Bill Gates*



tar à  
ificial

riar

Inspirações de maquiagem para o carnaval

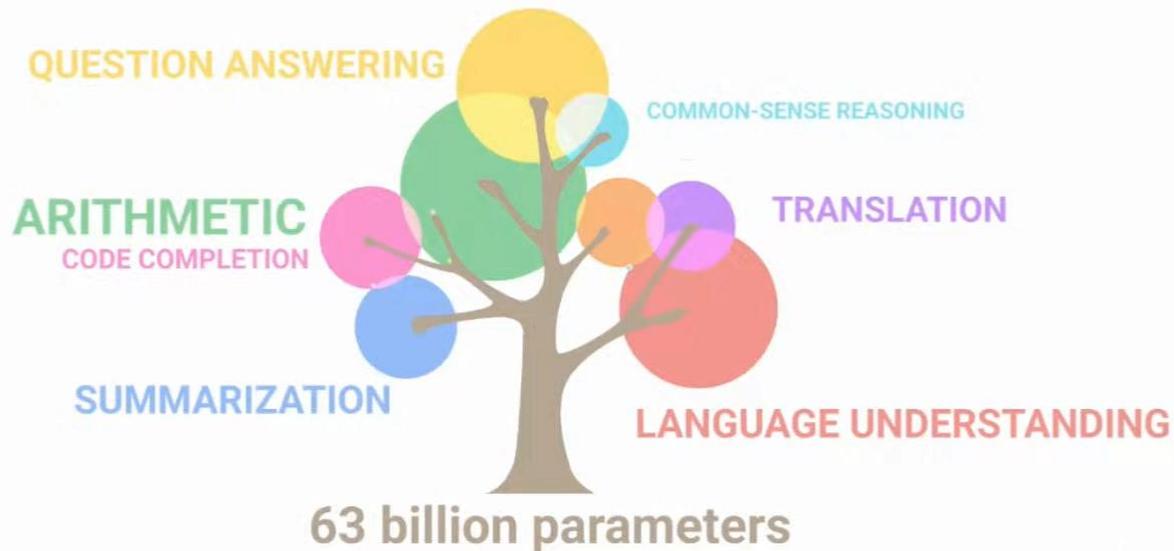
InCor incorpora tecnologia de imagem da Abbott que usa Inteligência Artificial para angioplastias

**Extensão para Chrome usa ChatGPT para responder e-mails do Gmail**

*"iPhone moment for AI" Jensen Huang, Nvidia's CEO about rise of Generative AI*



© marketoonist.com



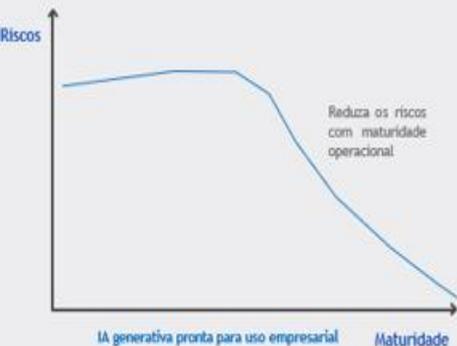
Made with

fliXier

# A Função Vital dos Dados na IA Generativa

## Gerenciando os riscos e aprimorando o valor

Ao iniciar sua jornada, você precisa fazer investimentos correspondentes em tecnologia e treinamento para reduzir os riscos e oferecer mais controle e valor crescente à sua organização.



Bring GenAI to Your Data

## Os Dados são o Grande Diferencial !

Dados e a qualidade desses dados são elementos cruciais para projetos de IA Generativa e quanto mais se utiliza os próprios dados, maior será o valor dos seus projetos para a organização

## Noções básicas sobre os riscos de dados na IA Generativa



Perda de propriedade intelectual



Vazamento de dados



Problemas de privacidade



Violações de conformidade



Perda de credibilidade e integridade



Viés



Violação de IP

# Problemas e Desafios comuns relacionados a adoção de IA

- ✓ Confusão sobre arquitetura
- ✓ Preparação de Dados, Qualidade e Organização
- ✓ Infraestrutura Compartilhada
- ✓ Ferramentas e Modelos
- ✓ Gestão de Recursos e Conhecimento
- ✓ Acumulação de Silos e Recursos

# Problemas e Desafios comuns relacionados a adoção de IA

- ✓ Comportamento e Viés Tóxico
- ✓ Alucinação e Respostas Incoerentes
- ✓ Conteúdo Congelado
- ✓ Conhecimento Específico
- ✓ Complexidade dos Algoritmos
- ✓ Infraestrutura de Larga Escala – Treinamento/Inferência

## Coesão – Coerência – Contexto

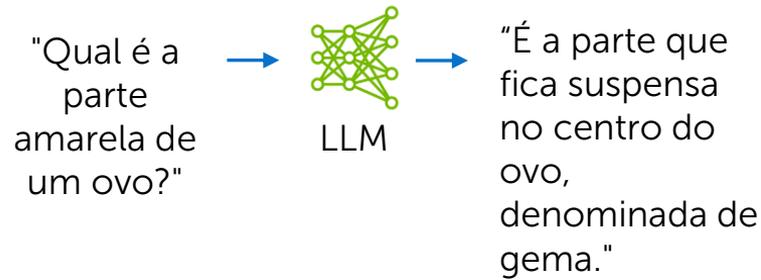
Os chatbots de IA Generativa, como o ChatGPT, utilizam um modelo de linguagem de larga escala (LLM) com capacidade de processamento de linguagem natural (PLN) para compreender e gerar textos coesos e contextualmente relevantes.

Para realizar essa tarefa, os modelos precisam ser treinados em grandes volumes de dados com o objetivo de aprender relações e padrões linguísticos.

A partir dessa base de informações, as IA's geram respostas em um processo de predição que busca determinar a continuação mais provável de um texto. Isso envolve analisar o contexto e o que foi dito anteriormente para elaborar uma resposta coerente e natural — sempre de acordo com os dados em que o modelo foi treinado.

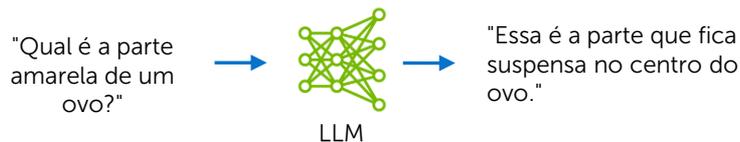
# LLM's – Exemplo de Resposta sem Contexto

## Resposta Sem Contexto (Coesão/Coerência)

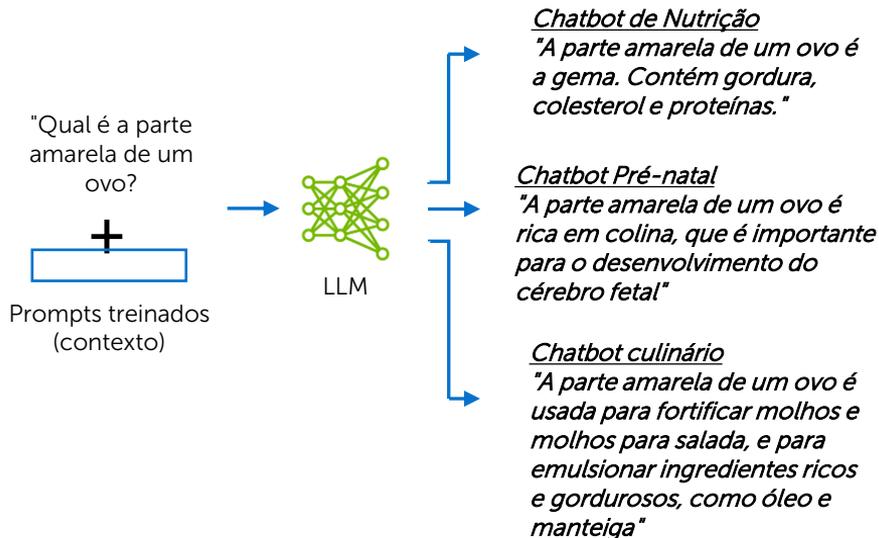


# A personalização é necessária para lidar com tarefas específicas do negócio - Contexto

Resposta Sem Contexto (com Coesão/coerência)



P-Tuned Response



# Casos de uso corporativos exigem conhecimento específico do domínio

Codifique e incorpore sua IA com as informações em tempo real da sua empresa para fornecer as respostas mais recentes

## Resposta de base/modelo personalizado

"Qual era a pressão do tanque às 23h de ontem à noite?"



"Fui treinado há 2 meses e não tenho os dados atuais"

## Modelo com técnicas RAG/C-RAG

"Qual era a pressão do tanque às 23h de ontem à noite?"



"A pressão às 23h da noite de ontem era de 345 psi"



Dataset (RAG/C-RAG)

**70%**

Dos dados corporativos é inexplorado  
Desbloqueie muitas novas oportunidades para maior inteligência



Retreinamento menos frequente  
Economia significativa de tempo e custo a longo prazo para manter LLMs

# Modelos de Implementação de IA: avaliando as vantagens e desvantagens de custo e valor

Os três primeiros tipos de modelos de implementação abaixo são o que 90% das organizações estão implementando agora. O modelo de IA que você escolher dependerá do nível de preparo da ciência de dados de sua organização, dos padrões de implementação e das implicações de cada um deles. O aumento de modelos é um ótimo ponto de partida, e a maioria das organizações migra dessa fase para modelos de ajuste fino.

## Inferência simples com um modelo pré-treinado

Conhecido como "prompt engineering", que faz uma pergunta a um modelo pré-treinado e recebe um resultado. O uso do ChatGPT é um exemplo de inferência simples.

## Aumento de modelo

Inferência mais acesso a seus dados, para o qual o RAG é um caso de uso comum. É um acesso fácil para tornar um modelo de IA generativa mais inteligente com base em seus dados.

## Modelos de ajuste fino

Envolve alterar a ponderação de um modelo e informá-lo com seus dados. Aproveitado à medida que as organizações começam a dimensionar seus modelos, ele oferece melhores resultados, mas também exige mais esforço para ser estabelecido.

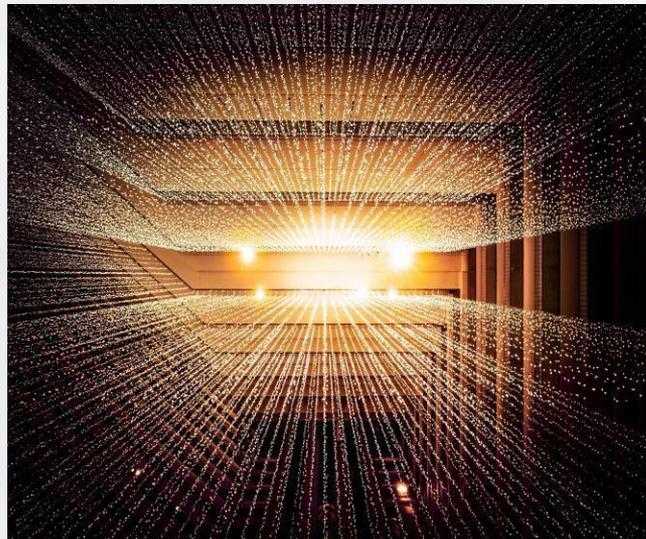
## Treinamento de modelos

Inclui a criação de um modelo muito específico e o treinamento com um conjunto de dados. Geralmente, isso requer a maioria dos recursos e trabalho e normalmente é reservado para resolver os problemas mais complexos.

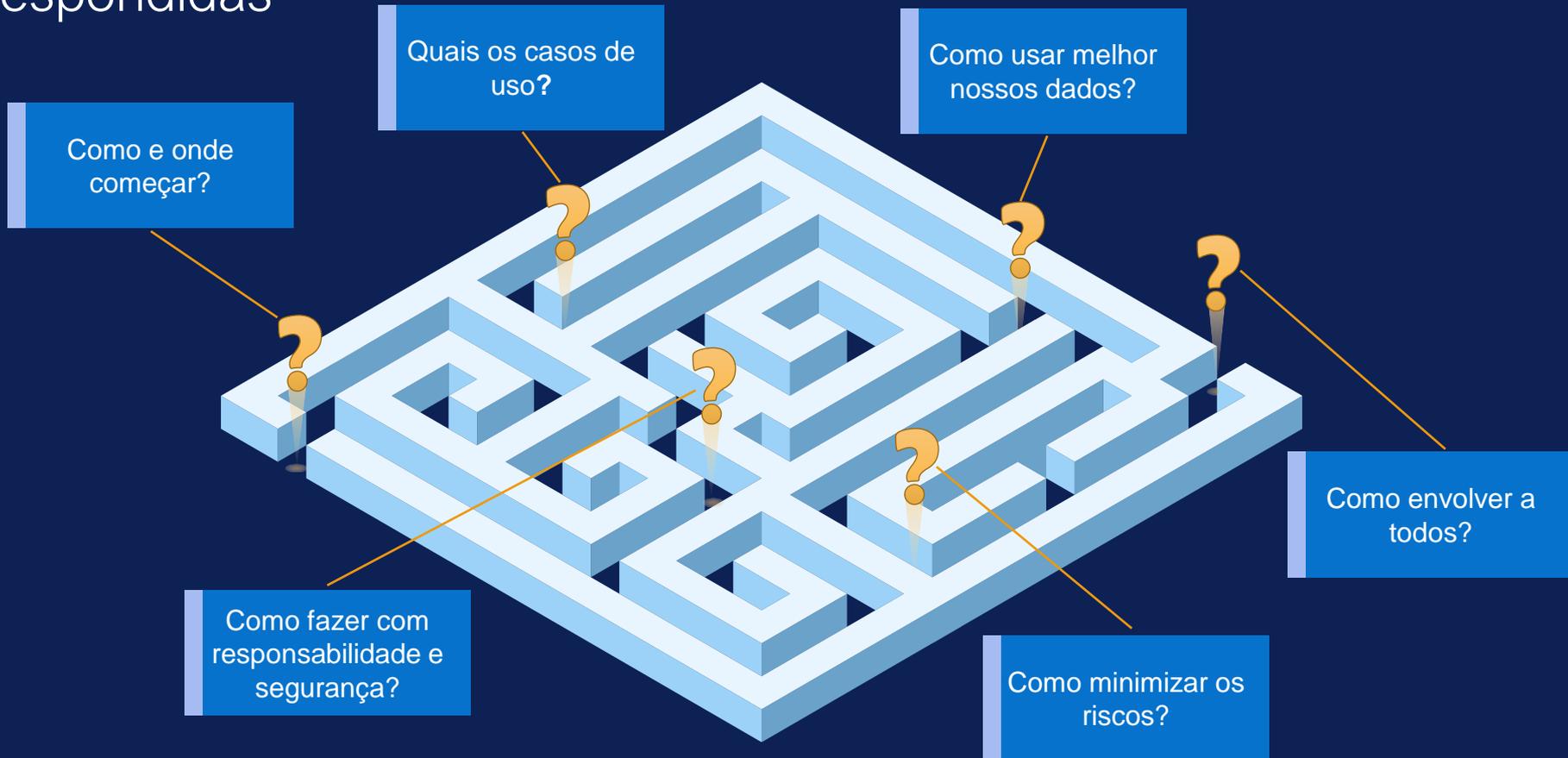
	Modelo pré-treinado	Aumento do modelo (RAG, PAL)	Modelos de ajuste fino	Treinamento de modelos	
<b>Definição do modelo</b>	Modelo inicialmente treinado ou pré-treinado  Modelo do transformador				<ul style="list-style-type: none"> <li> Quantidade mínima</li> <li> Quantidade baixa</li> <li> Quantidade média</li> <li> Quantidade alta</li> <li> Quantidade muito alta</li> </ul>
<b>Esforço</b>					
<b>Custo</b>					
<b>Valor e diferenciação</b>					
<b>Integração de dados</b>					
<b>Infraestrutura</b>	Client – servidor	Client – servidor	Otimizado para a GPU	Implementação de GPU grande	
<b>Habilidades</b>	Operações de TI	Desenvolvedor	Cientista(s) de dados	Cientista(s) de dados	

# Fundamentos para Construção - IA Generativa

- ✓ Bons Modelos de Dados (LLM's) - Qualidade dos Dados
- ✓ Acesso a Grandes Quantidades de Dados Brutos (ETL)
- ✓ Infraestrutura de Alto Desempenho
- ✓ Arquitetura - Robusta - Escalável - Orientado ao Caso de Uso
- ✓ Custo Acessível para Processamento Rápido e Eficiente
- ✓ Software de Ciência de Dados - Framework
- ✓ Segurança - Guardrails
- ✓ Expertise

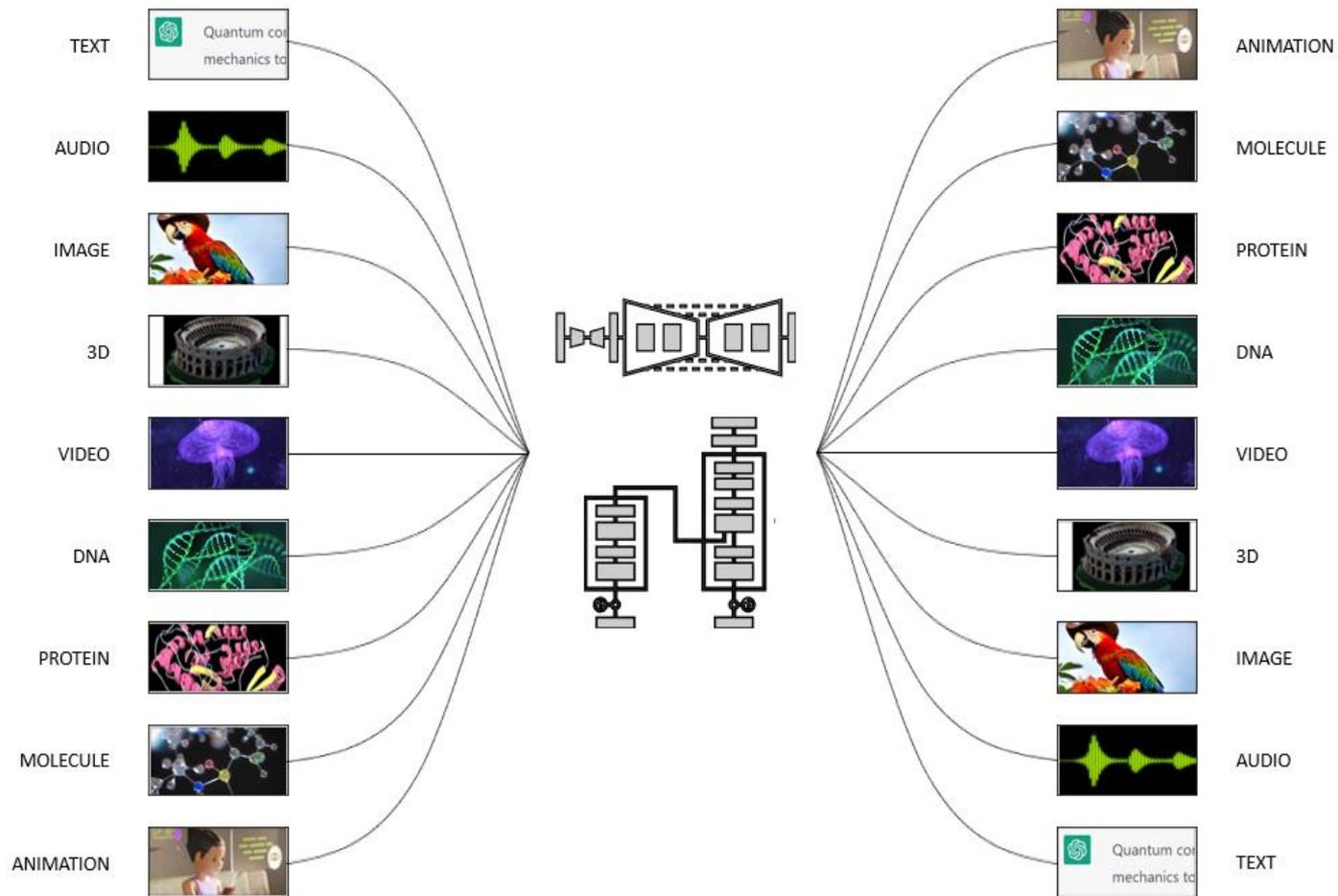


# A jornada para o GenAI tem muitas perguntas que precisam ser respondidas

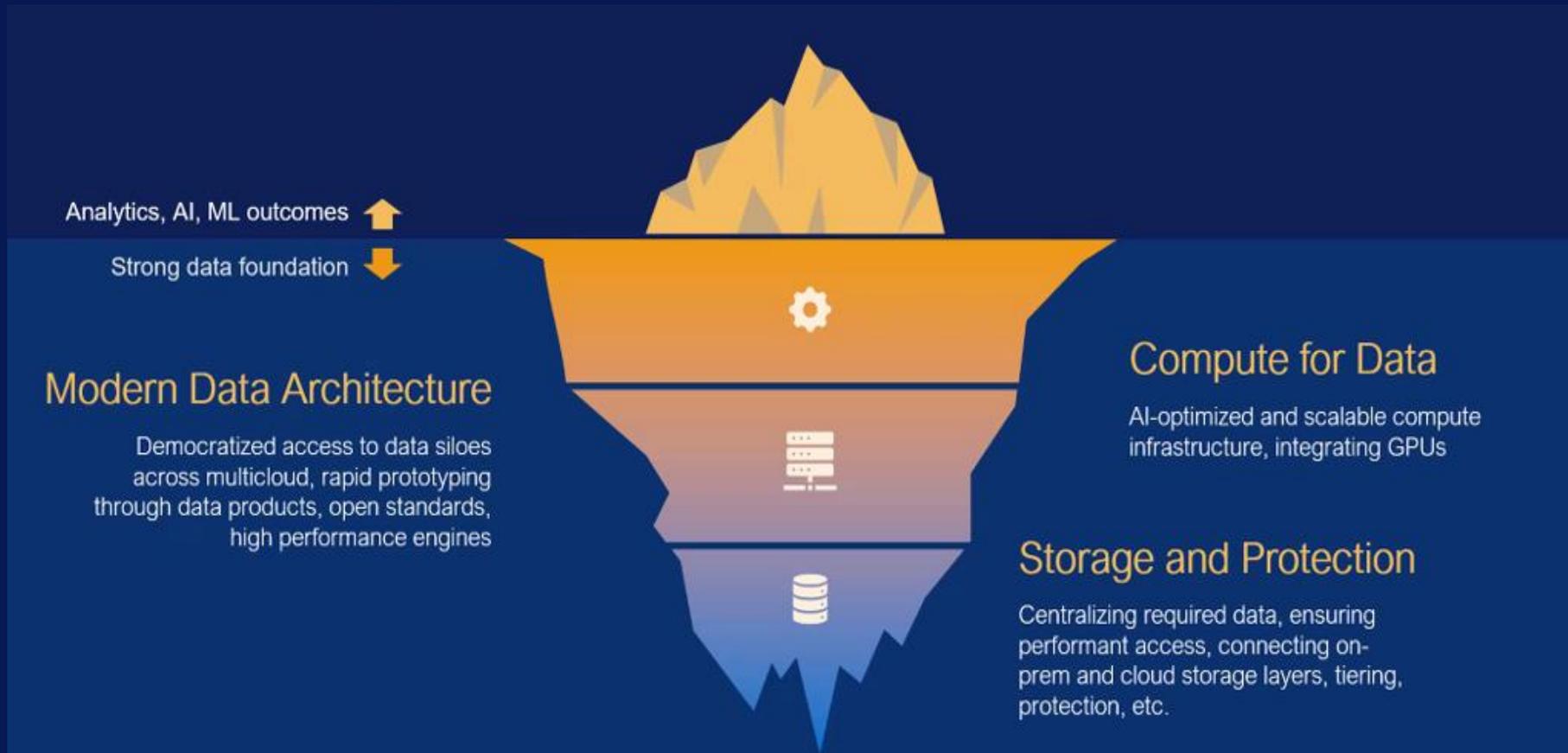


# What's Generative AI

Bring Your Own ChatGPT  
Bring GenAI to Your Data

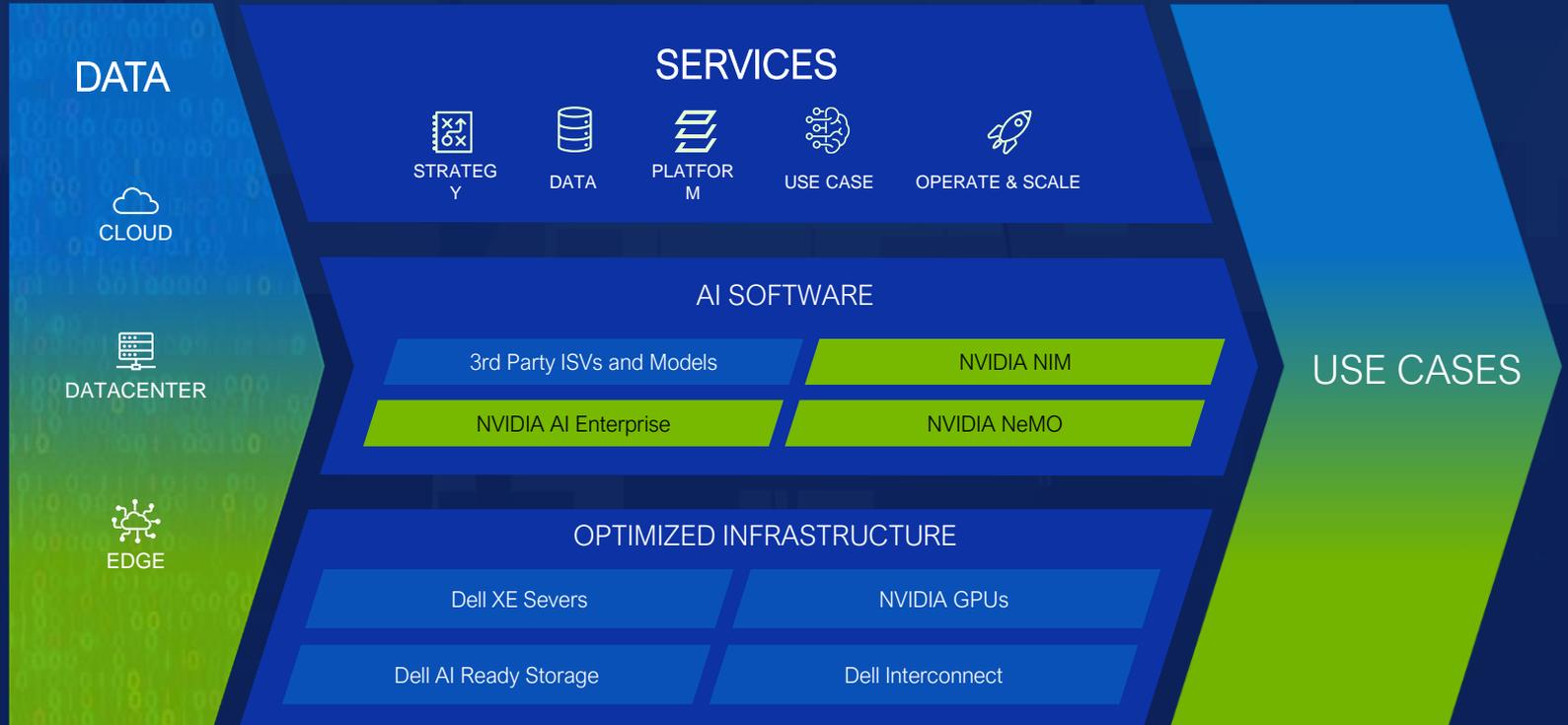


# Uma base sólida e estruturada é Fundamental para IA e Analytics



# The Dell AI Factory with NVIDIA

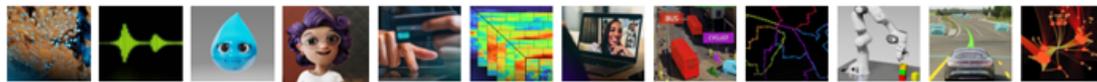
INDUSTRY'S FIRST END-TO-END OPERATIONAL AI PLATFORM



# NVIDIA AI ENTERPRISE SOFTWARE

The Proven Standard for Enterprise AI

AI Workflows, Frameworks and Pretrained Models\*



Medical imaging    Speech AI    Conversational AI    Customer Service    Recommenders    Physics ML    Communications    Video Analytics    Logistics    Robotics    Autonomous Vehicles    Cybersecurity

AI and Data Science Development and Deployment Tools

Cloud Native Management and Orchestration

Infrastructure Optimization

Accelerated Infrastructure



Cloud



Data Center



Edge

\*NVIDIA NGC public catalog provides a complete listing of over 50 supported frameworks and pretrained models.

## NVIDIA Base Command

Provides AI workload and cluster management, libraries that accelerate compute, storage, and networking infrastructure, and an operating system optimized for AI workloads.

## Optimize Performance

NVIDIA NIM and CUDA-X microservices provided an optimized runtime and easy to use building blocks to streamline generative AI development

## NVIDIA API Catalog

Pre-trained models help deliver results faster, more cost-efficiently

# NVIDIA NIM - Inference Microservices for Generative AI

The best way to deploy AI models on Dell AI Factory

## Inference Microservice

- > Industry-standard APIs
- > Pre-configured container for simplified deployment
- > Optimized inference engines built on Triton™ Inference Server, TensorRT™, and TensorRT-LLM
- > Enterprise management data, including identity, metrics, health checks, and monitoring
- > Part of NVIDIA AI Enterprise

## Optimized AI Models

Large language models (LLM), image, video, 3D models, automatic speech recognition (ASR), text-to-speech (TTS), video language models (VLM), biology, and retrieval models

## Accelerated Infrastructure

The ability to deploy with a single command or orchestrate and auto-scale with Kubernetes on NVIDIA-accelerated infrastructure anywhere



# Portfólio INTEL para IA

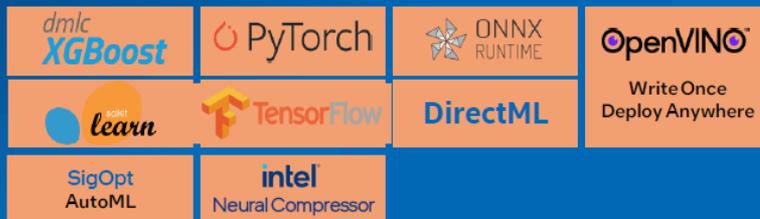
Engineer Data

Create Models

Optimize & Deploy



Data Analytics at Scale\*



Machine & Deep Learning Frameworks, Optimization and Deployment Tools\*

1  
oneAPI

Intel® oneAPI Deep Neural Network Library

Intel® oneAPI Collective Communications Library

Intel® oneAPI Math Kernel Library

Intel® oneAPI Data Analytics Library

Open, cross-architecture programming model for CPUs, GPUs, and other accelerators

CLOUD & ENTERPRISE



CLIENT & WORKSTATION



EDGE



Accelerate end-to-end data science and AI



Intel® Developer Cloud and Intel® Developer Catalog

Try the latest Intel tools and hardware, and access optimized AI Models

cnvrg.io

Full stack ML operating system

Intel® Geti

Annotation/training/optimization platform



Intel optimizations and fine-tuning recipes, optimized inference models, and model serving

# Construindo Pipelines de Dados de IA

## Establish Strategy



Define Vision /  
Solution Arch



Identify  
Use Cases / Models

## Prepare Data



Data  
Management



Data  
Engineering

## Gen AI Platform



Gen AI  
Software



Gen AI  
Hardware

## Deploy & Test Model



Language  
Modeling



Implement AI  
Use Cases

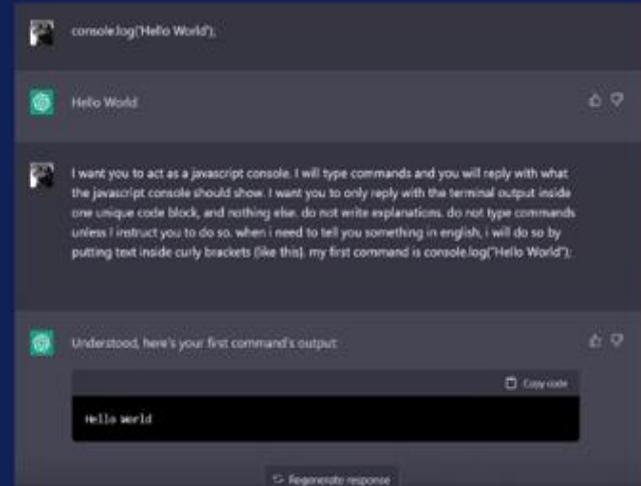
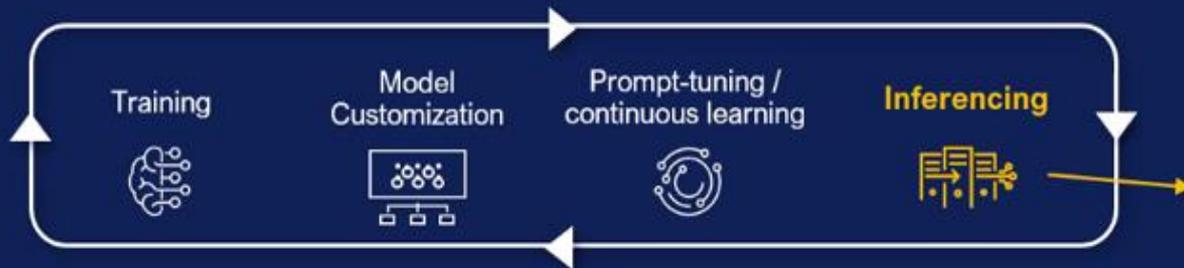
## Operate & Scale



AI/Data  
Governance



Gen AI  
Operate / Optimize



# Dell AI Factory Accelerates AI Adoption

**SIMPLE**



Comprehensive AI optimized solutions make AI easier, augmenting skills gaps and addressing data readiness

**SECURE**



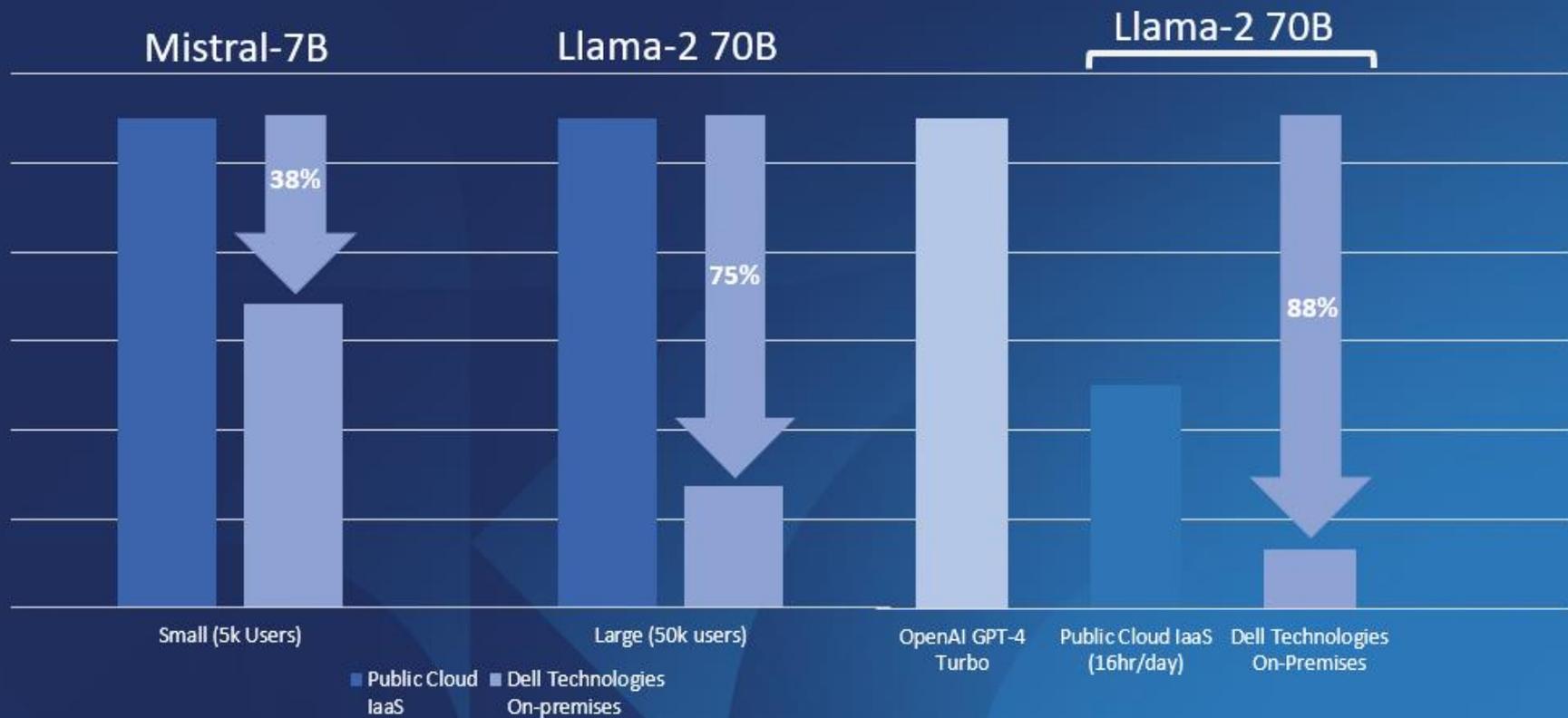
Trusted infrastructure, full-stack solutions and services ensure valuable data is protected, managed and accessible

**ECONOMICAL**



Right sizing AI investment to fit their use cases requirements with Dell extensive AI portfolio to lower costs.

# Cost-effectiveness of Dell Technologies Solutions for Gen AI



# One-stop shop for AI

1

## Dell Validated Design

- Validated architecture
- Accelerated deployment
- On-prem & APEX

2

## AI / GenAI toolkit bundle

- Complete set of tools
- Opensource ecosystem
- Intuitive IDEs
- Customizable sample use cases

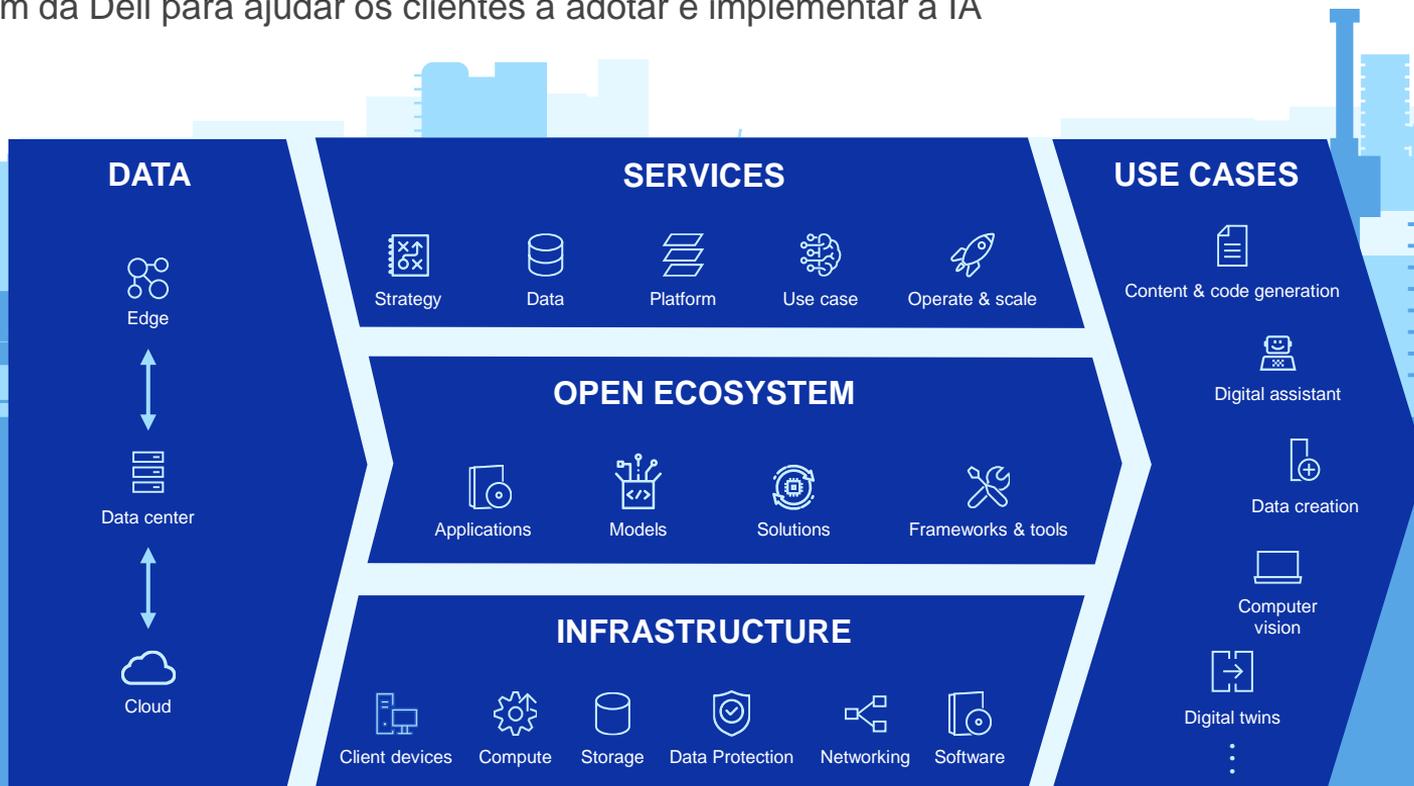
3

## Dell Professional Services

- Forbes & Forrester: Awarded consulting organization
- Proven methodology
- Accelerate AI / Gen adoption at scale
- Flexible services bundles

# DELL AI FACTORY E USE CASES

A abordagem da Dell para ajudar os clientes a adotar e implementar a IA



# Open Ecosystem HUB for AI Innovation

## GLOBAL SYSTEMS INTEGRATORS



## SILICON PROVIDERS



## COLOCATION PROVIDERS



## HOSTING PROVIDERS



HARDWARE



## MODELS



## FRAMEWORKS



## ORCHESTRATION & TOOLING



## DATABASES & INFO RETRIEVAL



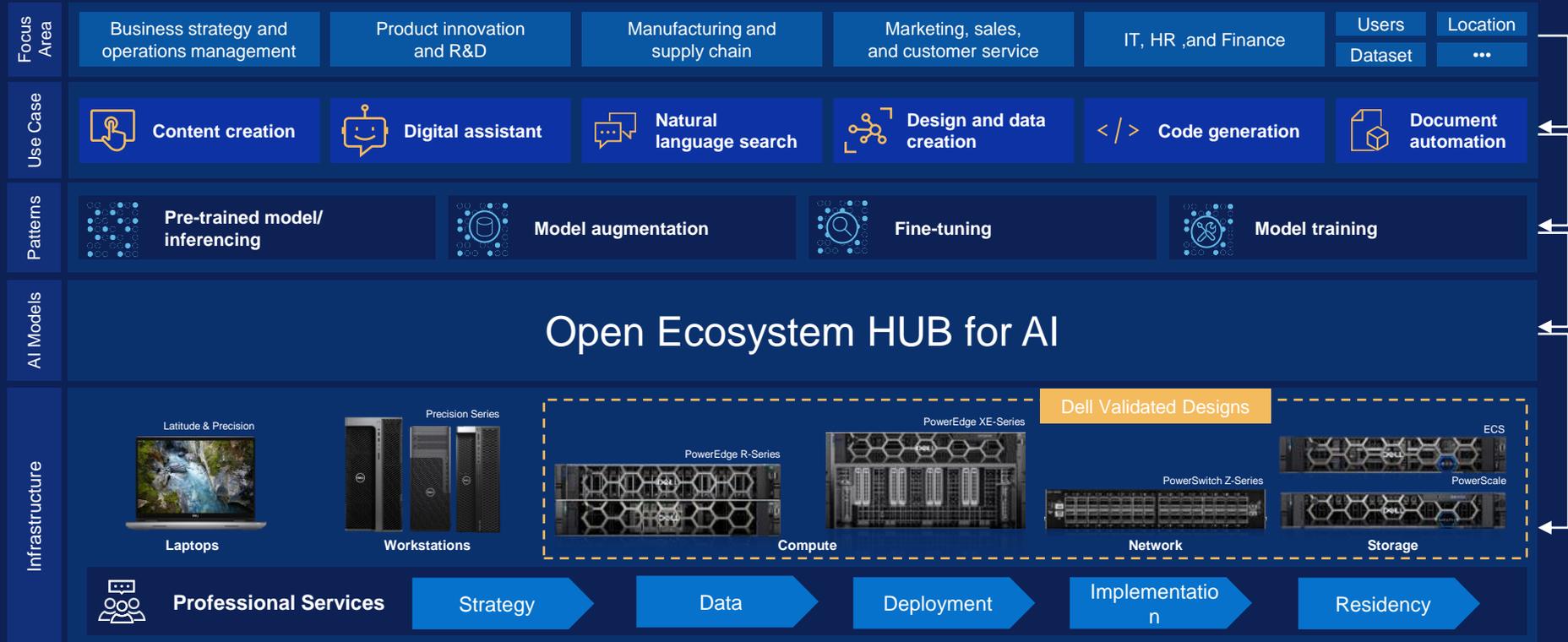
## DATA MANAGEMENT



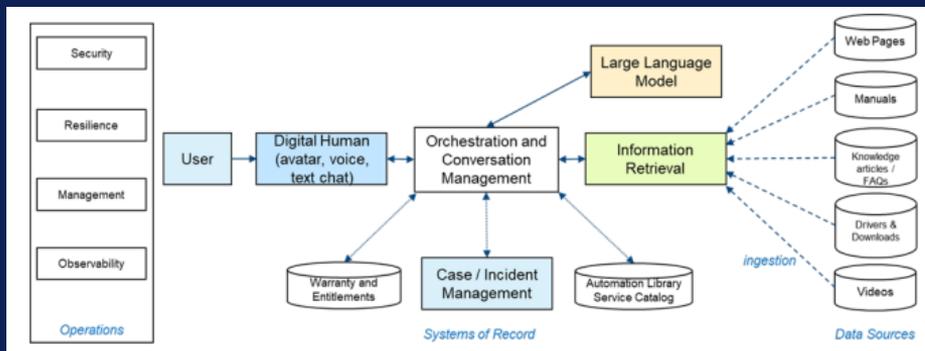
SOFTWARE

Ecosystem

# Right-sizing your AI investment

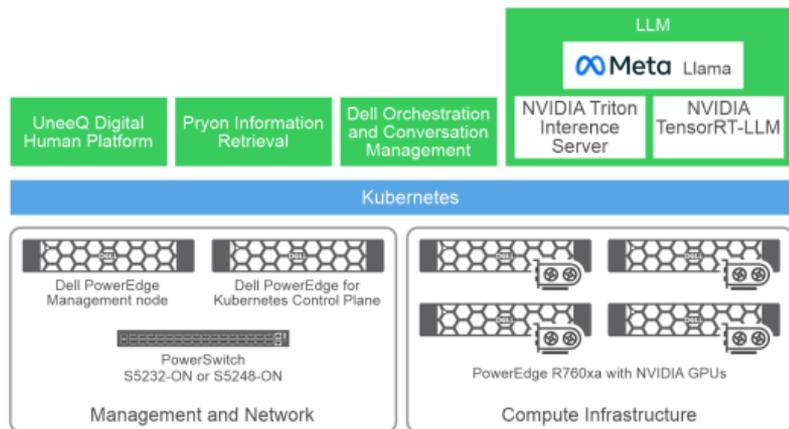


# Digital Assistant Conceptual Architecture

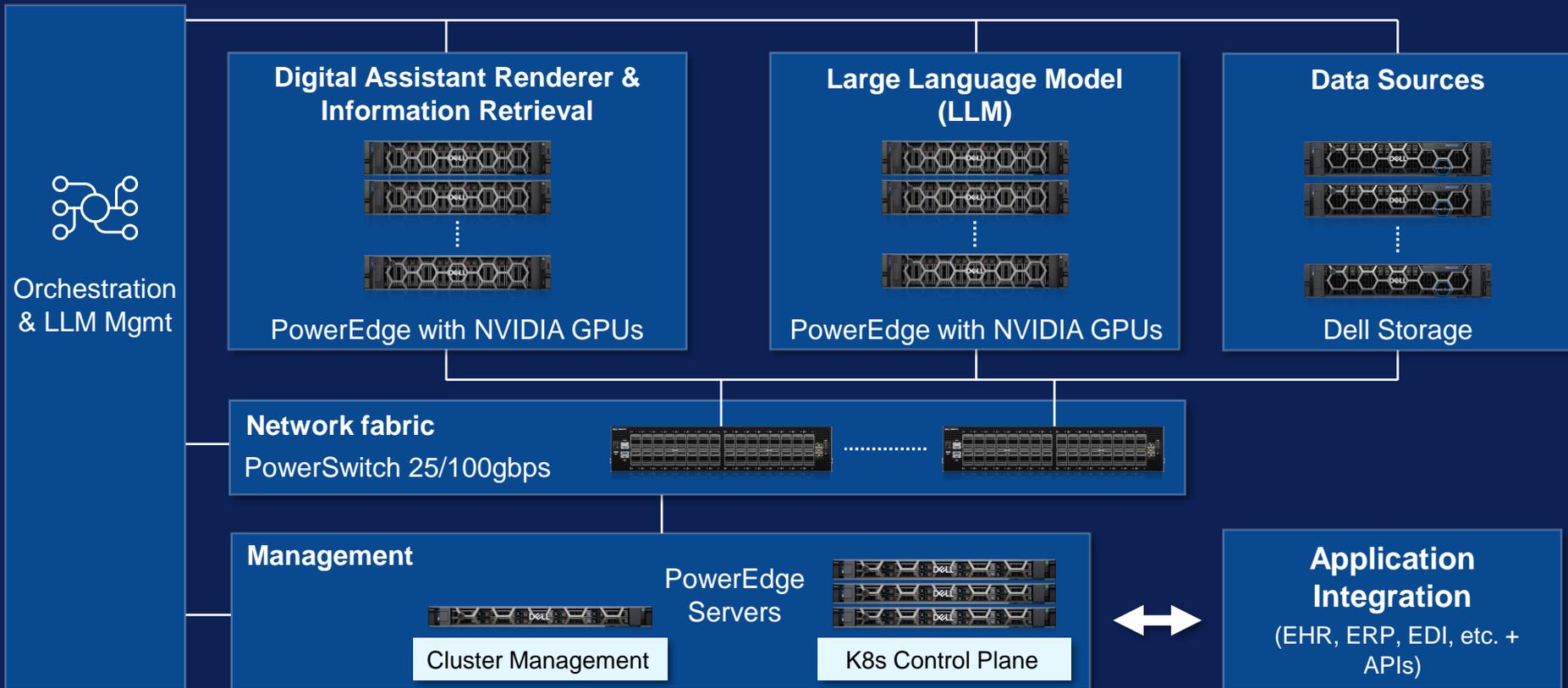


## Orchestration and Conversation Management

- Speech to text (STT)**—STT is a technology that converts spoken language into written text. It is commonly used in voice recognition systems, transcription services, and various applications where it is more convenient or efficient for the user to speak rather than type.
- Text to Speech (TTS)**—TTS is the opposite of speech to text. It converts written text into spoken words. It is often used in assistive technologies for visually impaired users, in navigation systems, or in applications for which audio content is more suitable. The text-to-speech component is only able to accept complete sentences (compared to streamed tokens) to achieve proper pronunciation. Our LLM observability must add a 'time to first sentence' metric to the standard set of LLM metrics.
- Natural Language Processing (NLP)**—NLP is the interaction between computers and humans through natural language. The goal is to enable computers to understand, interpret, and generate human language in a valuable way. NLP performs several tasks such as language understanding, language generation, translation, and sentiment analysis. Generative LLMs perform most NLP tasks well, however our orchestration layer maintains the ability to perform classical NLP or keyword recognition tasks if needed.
- Prompt engineering**—Prompt engineering is the process of designing and optimizing prompts to guide an AI model's responses effectively. It is a crucial aspect of interacting with AI models, especially in language models where the quality and structure of the prompt can significantly influence the model's output.
- Retrieval Augmented Generation (RAG)**—RAG is a method used in AI language models in which the model retrieves relevant documents or information from a database before generating a response. This method allows the model to retrieve real-world facts and details, enhancing the accuracy and richness of its responses.
- Translation**—The digital assistant solution can be equipped with the ability to translate text from one language to another. The task involves understanding the semantics and context of the source language and accurately reproducing the meaning in the target language. This process is crucially important in multinational environments in which original content is created in a single language and translated on demand by a digital assistant.
- Integration hub**—An integration hub facilitates the seamless integration of various systems of record (see Systems of record) into the overall solution. It acts as a bridge, connecting disparate systems and enabling them to communicate and share data

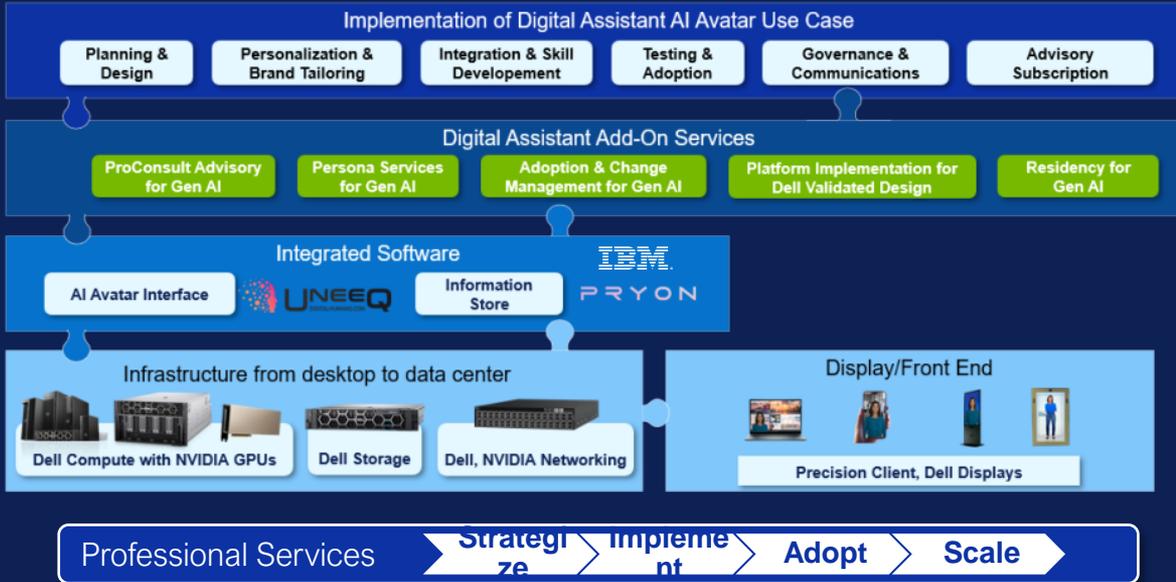


# Optimized Enterprise-Scale Architecture



# Dell Professional Services

Bringing the complete solution together, Dell Professional Services for Digital Assistants provide strategic guidance for customer use cases, personalized AI Avatar development, platform implementation with data source integrations and ongoing support to drive continued success, upskilling and optimization.



## Implementation Services for Digital Assistants - AI Avatar Use Case:

Create a custom Digital Assistant featuring a humanistic AI Avatar interface that is tailored to customer use cases, integrated with their data, and personalized for their brand and audience.

## Implementation Services for Digital Assistants - AI Avatar Platform:

Deploy the Digital Assistant on Dell hardware and integrated with their monitoring ecosystem, leveraging the Dell Validated Design for Digital Assistants.

# Serving Generative AI future use cases, today





## Manufacturing

- Structural analysis
- Fluid dynamics
- Packaging
- Impact modeling



## Geosciences Oil & Gas

- Seismic
- Reservoir modeling



## Government Classified

- Homeland security
- Defense
- Nuclear safety



## Government & Academic Research

- Particle physics
- Life sciences
- Humanities
- Climate modeling



## Financial Services

- Financial analytics
- High frequency trading
- Fraud detection



## Electronics Design

- Electrical design
- Circuit verification
- Board layout



## Life Sciences Pharma/Biotech

- Drug design
- Genomics
- Bio informatics



## Media & Entertainment

- Rendering
- Gaming



## Enterprise HPC

- Analytics – HPDA
- Marketing
- Simulations

# AI já está Transformando cada Segmento e Industria

**CREDIT CARD FRAUD**  
1.1B Credit Transactions / Day



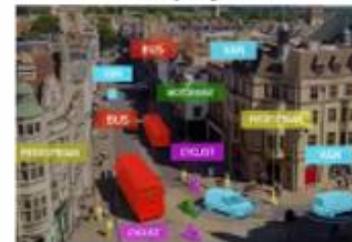
**CONTACT CENTER AI**  
500M Calls / Day



**MEETING  
TRANSCRIPTION**  
15B Meeting Minutes / Day



**PUBLIC SAFETY**  
> 1B Smart City Cameras Deployed



**PRODUCT  
RECOMMENDATIONS**  
300M E-commerce Visitors / Day



**RETAIL ASSET  
PROTECTION**  
\$275M Inventory Loss / Day



**MEDICAL IMAGING**  
10M Diagnostic Scans / Day



**INDUSTRIAL INSPECTION**  
94M Vision Sensors Installed by 2025



# Why Dell for Generative AI?

Dell's unique approach leverages simplified, tailored and trusted solutions that combine your data with the power of Generative AI to drive innovation and tangible business value.



Dell is the #1 worldwide provider in AI server plus storage infrastructure<sup>1</sup>



We offer full-stack scalable Generative AI solutions, in collaboration with NVIDIA



Dell Services provide deep AI expertise at every stage to accelerate tangible time-to-value

## #1 Worldwide

provider in AI server plus storage infrastructure

<sup>1</sup> Source: IDC Semiannual AI Tracker: Worldwide Server and Storage Revenue, 2021 and 2022 H1



Esta é a

**DELL** Technologies

**OBRIGADO**