

Essentials of Deploying AI in the Data Center



AI Use Cases

AI Concepts (ML, DL, Inferencing)

How GPUs revolutionized AI

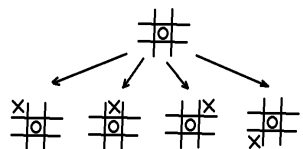
Nvidia AI Architecture

Infrastructure Planning

Data Center Facilities Planning

Infrastructure Provisioning and Management

AI to Revolutionize Future Data Centers



Artificial Intelligence

Using computers to do that requires human level Intelligence

1950s – 1980s



Machine Learning

Approach to AI that uses statistical learning algorithms to build Model from observed **data**

1980s – 2010s

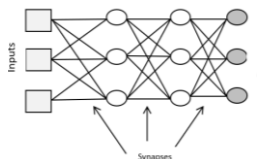
*Challenge:
Understanding and
extracting insights
from big data*



Deep Learning

Machine learning technique that is inspired by how human brain learns

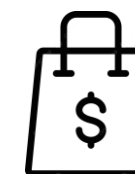
2010s – Today and the future



DNNs can achieve human level intelligence for many tasks but requires high computational power to train

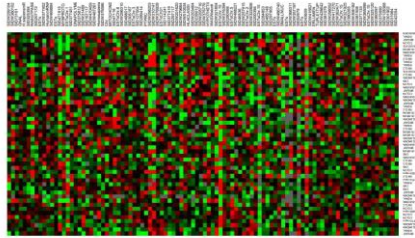


Data scientists uses DL to extract meaningful **data** from large datasets



Line of business owners use DL to optimize their business, reduce costs, Improve functionality and accuracy

AI in healthcare



Early Detection of diseases

DNA, Expression Microarrays, a regression problem facilitating the breakthrough study in biotechnology



Operational Efficiency

Predictive analytics, to identify patterns, and predict patient outcomes, enabling personalized and proactive treatment plans.



Precision Medicine

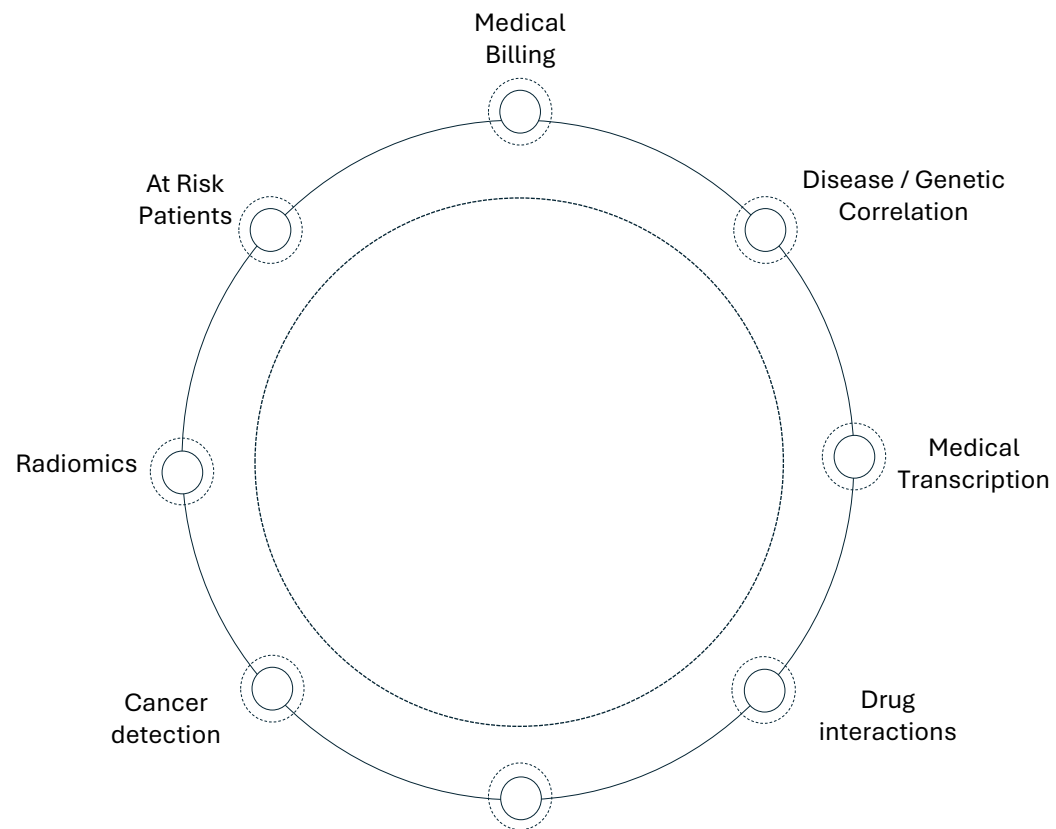
Descriptive analytics that facilitate the development of precise personalized medicines



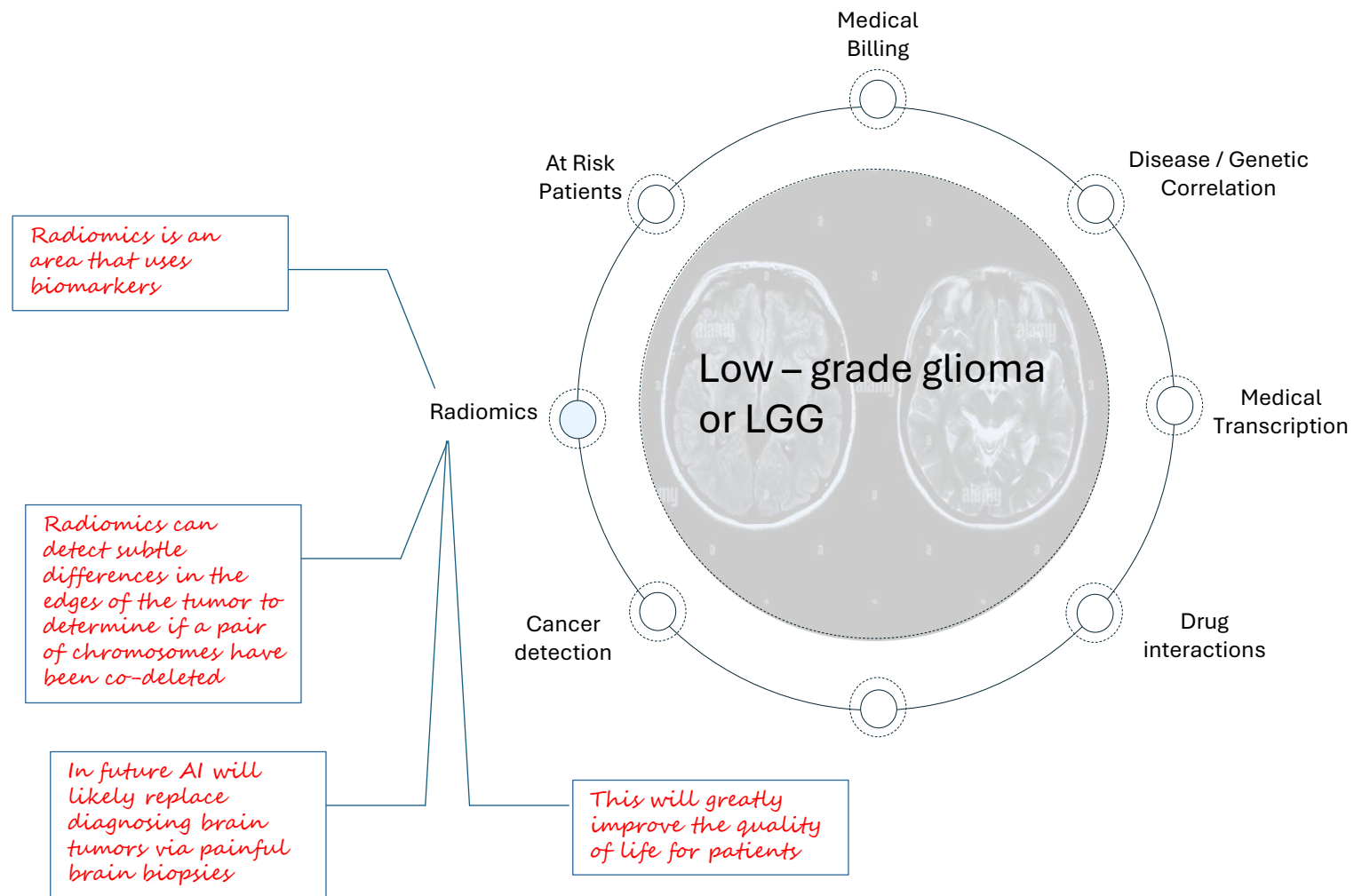
Drug Discovery

Researchers can bring new drugs to Market quicker and make sure that these drugs are monitored more efficiently

AI in healthcare



AI in healthcare



AI in Autonomous vehicles

Potential for Every Vehicle to be autonomous

- 1.5B vehicles in world today → 2B vehicles by 2035
- Mobility Services (Buses, Taxis, robotaxis, shuttles)
- 4M AV shipments by 2025 (New autonomous vehicles – Earthmovers, Forklifts, Delivery Bots, Tractors, Firetrucks etc)



**Reduced fatalities,
Injuries, and parking footprint**

Shared autonomous vehicles can substantially reduce parking space requirements



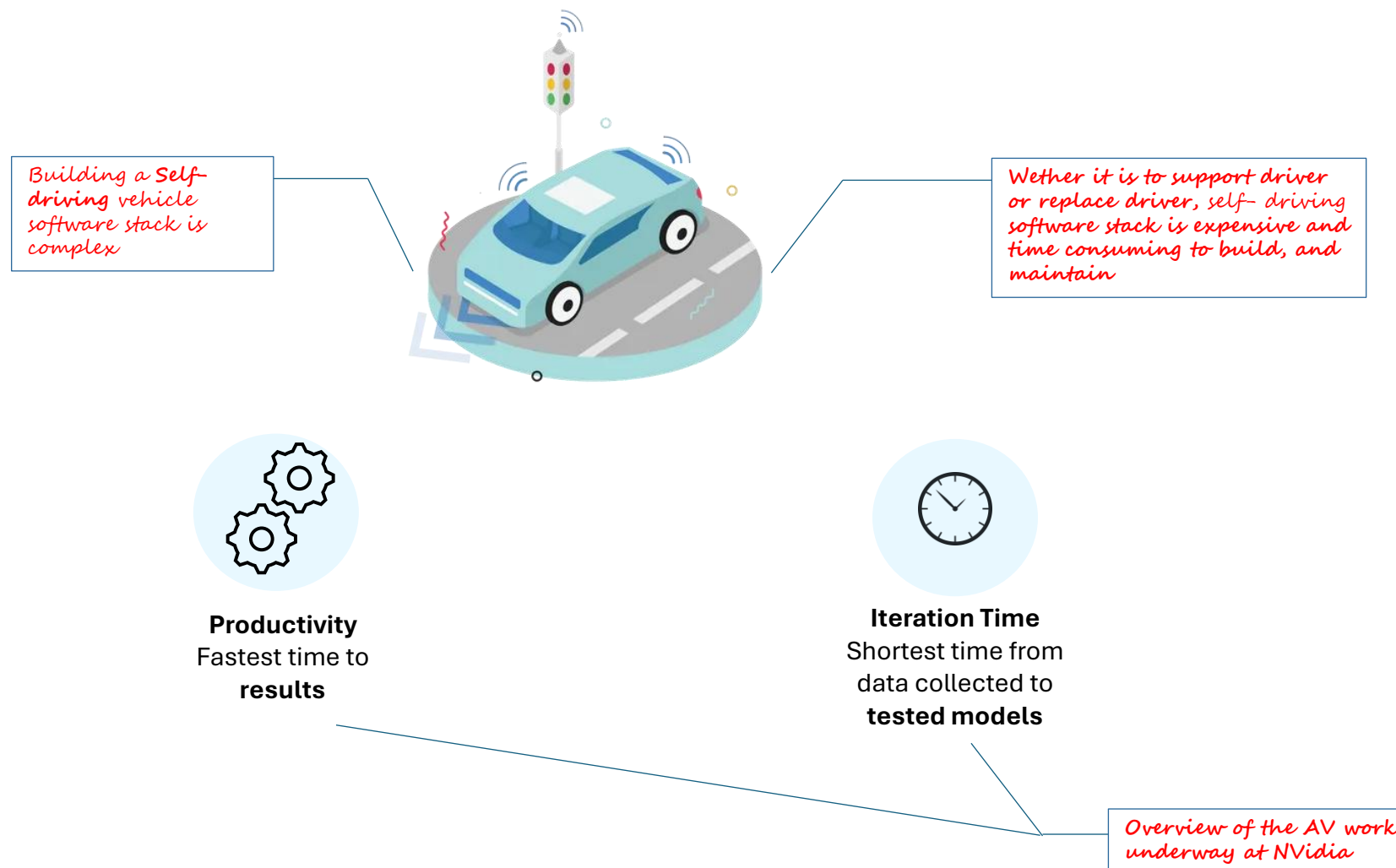
Mitigates shortage of delivery services and drivers

This shortage can be resolved with delivery bots and autonomous trucks

There is global shortage of 60000 truck drivers, this deficit is expected to triple by 2026

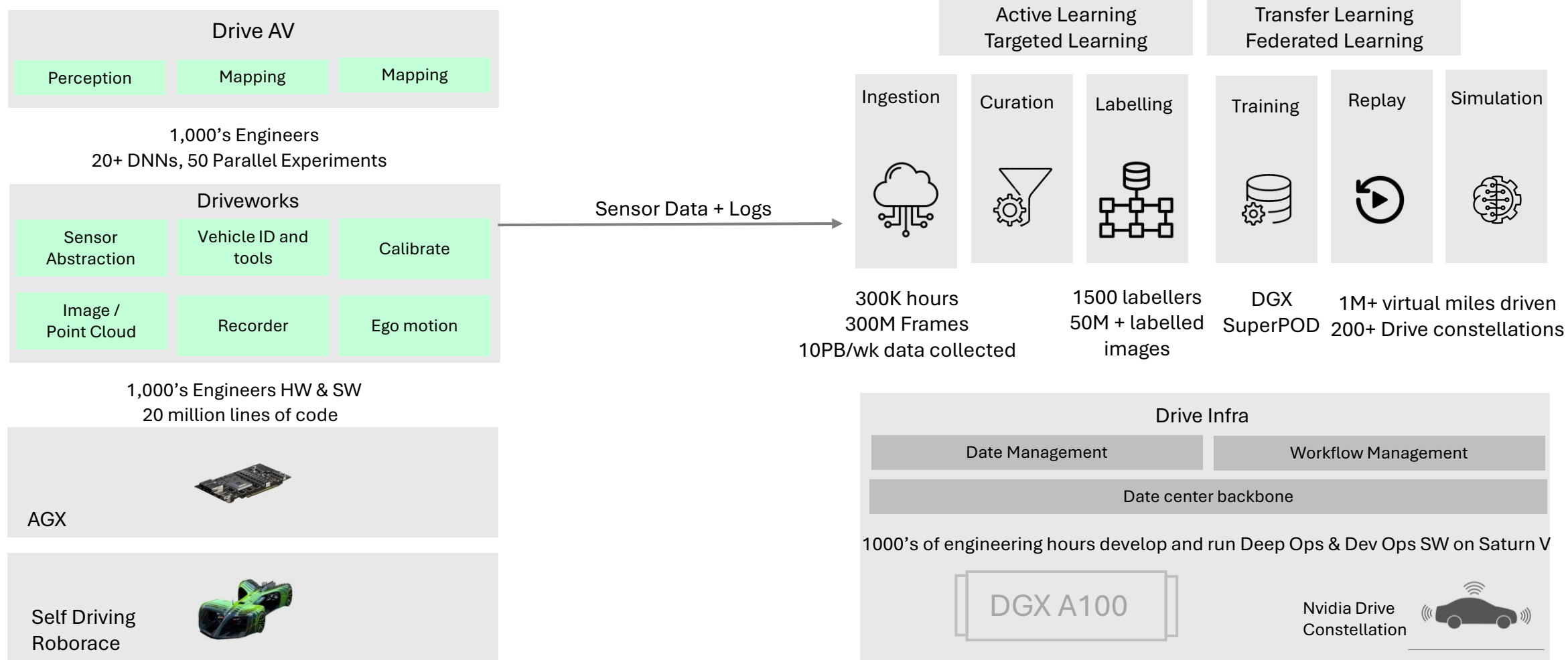
Nvidia Drive

E2E AV Solution to Enable Rapid, Large Scale AI Development & Testing



Nvidia Drive

E2E AV Solution to Enable Rapid, Large Scale AI Development & Testing

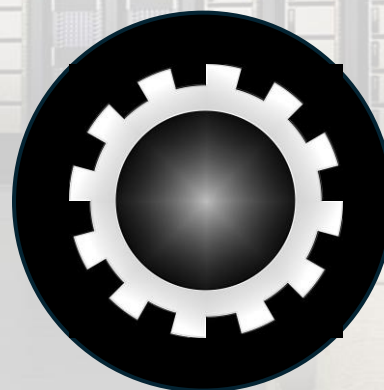


Nvidia Drive

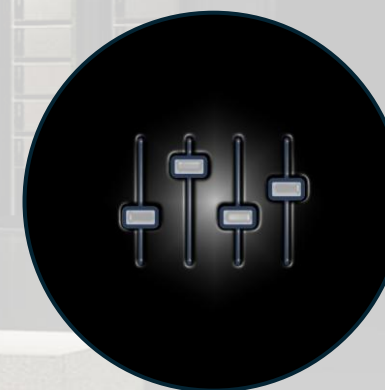
E2E AV Solution to Enable Rapid, Large Scale AI Development & Testing



Deploying



Managing



Optimizing

Large scale DGX SuperPOD systems

AV Workflow – Challenges and Pain Points



Accelerating Digital Transformation in Finance

AI / ML Optimizes Performance and Outcomes



Banking



Capital Markets



Insurance



Payments



Fintech

Fraud Detection

Using a range of AI techniques to better verify identity for Anti-Money Laundering (AML) and know your Customer (KYC) requirements

Risk Simulations

Using HPC to run Monte Carlo risk simulations and derivatives pricing

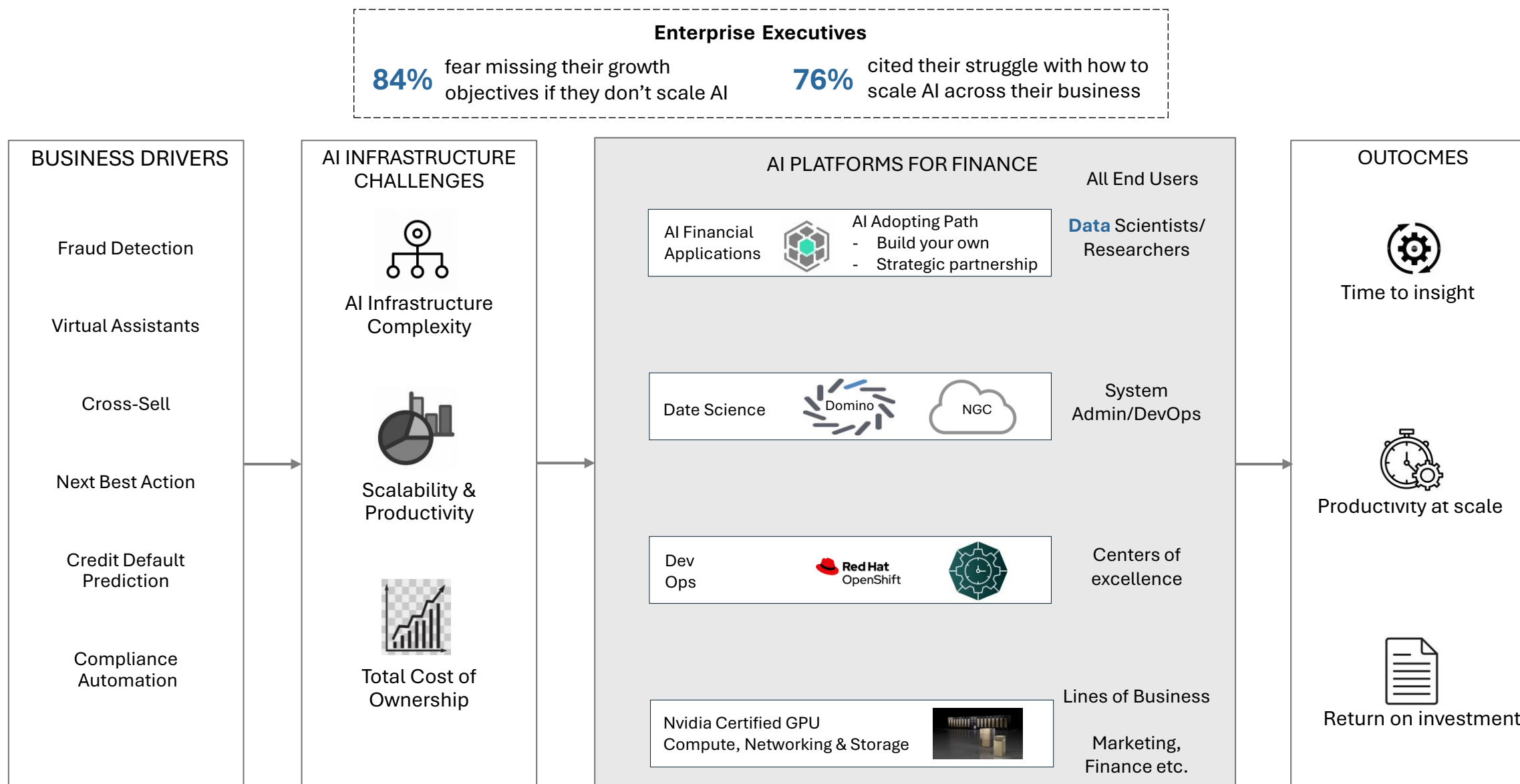
Underwriting

Enables lenders to integrate more sophisticated modelling techniques and alternative data into lending decisions or insurance pricing

Algorithmic Trading

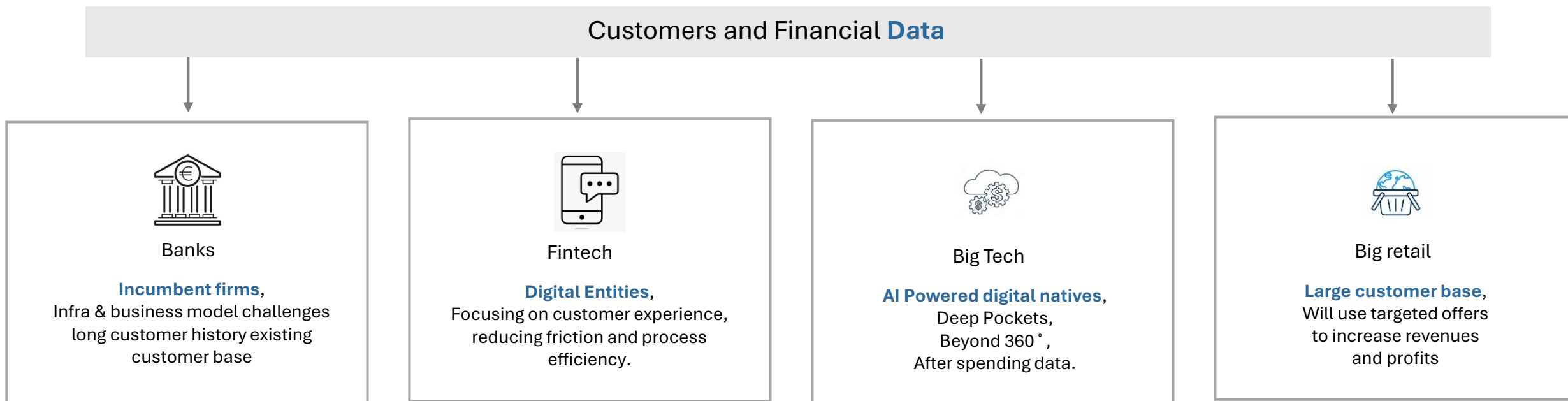
Create competitive advantages across a range of investment types

AI Platform for AI Driven Finance



Where is the Industry Going?

Your AI Transformative Strategy is incomplete without NVidia



AI-Based underwriting – Risk Prediction

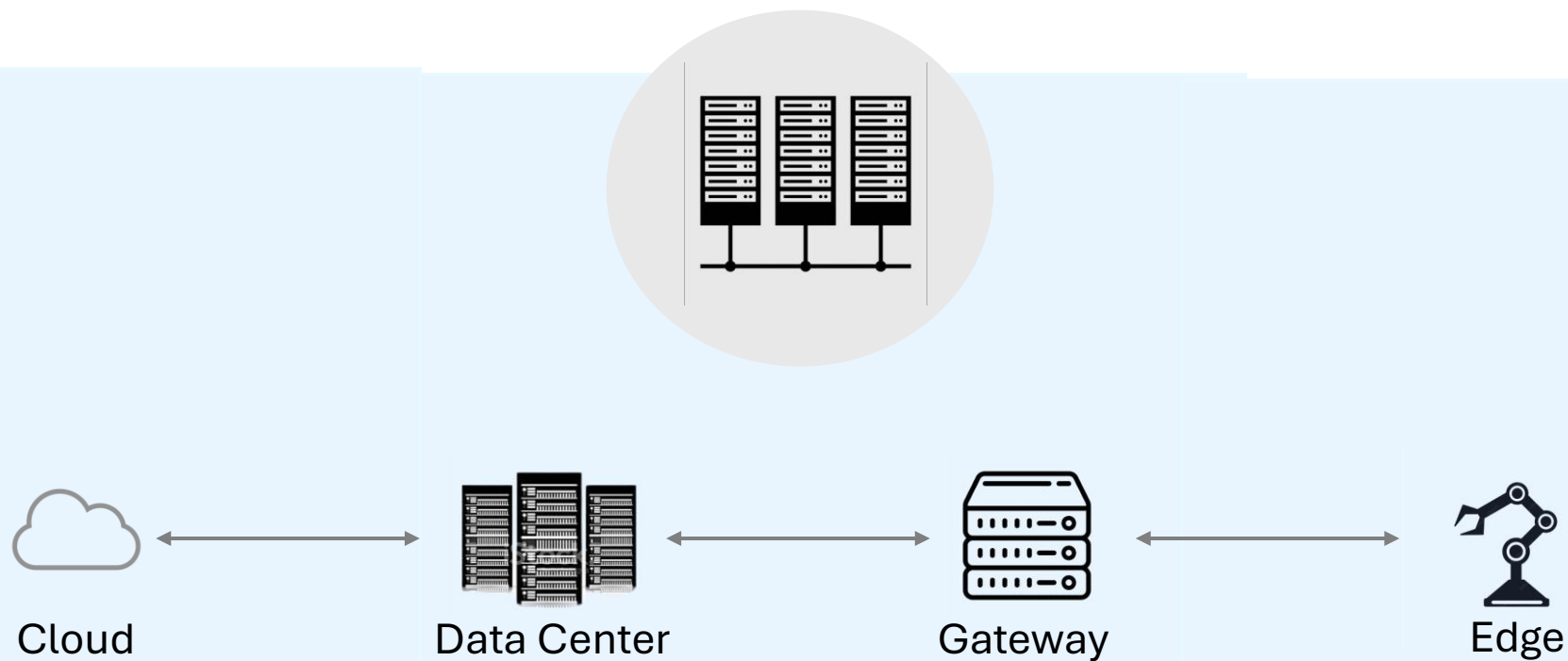
AI-enabled chatbots – Superior Customer Service

AI for fraud detection – Greater Security

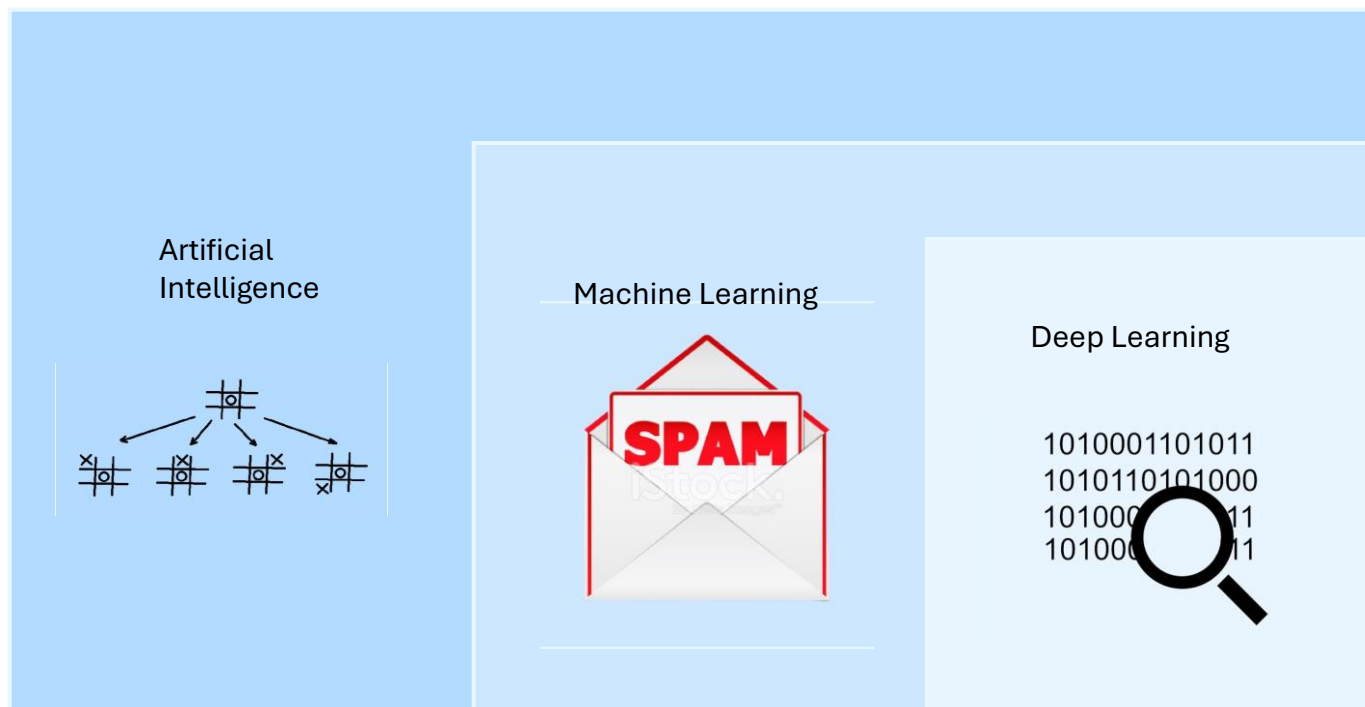
Industry Overview

From the cloud to the Edge

AI required huge amount of
compute power



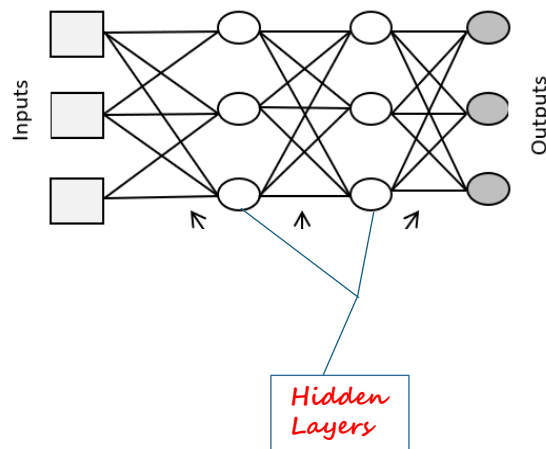
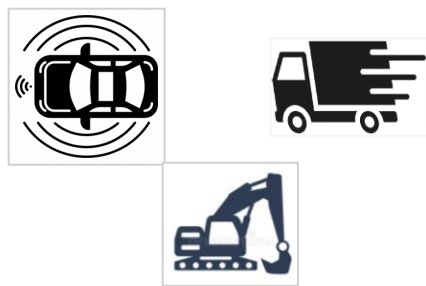
Deep Learning Approach



Deep Learning Approach

Deep Neural Network «Model»

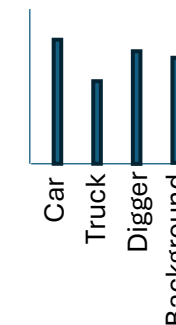
Labelled Trained Data



Prediction



Object class Predictions



error ~ [label - prediction]

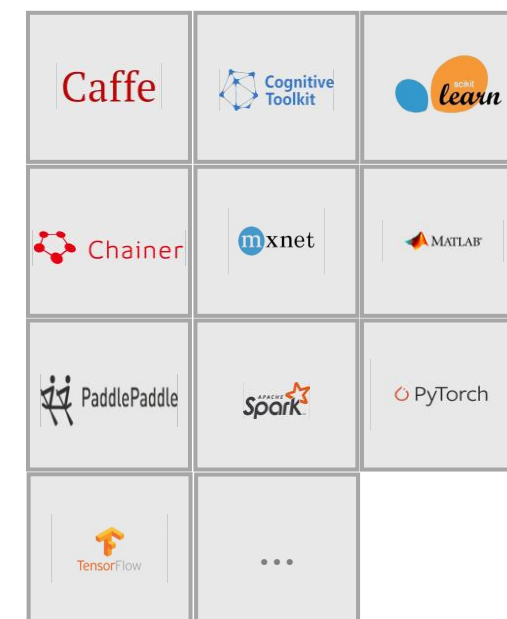
Back propagate errors for
parameter update

Machine/Deep Learning Frameworks

Essential Tools for Data Scientists, Researchers, and Engineers



Computer Vision
Natural Language Processing
Speech and audio Processing
Robot Learning
more...



Machine/Deep Learning Frameworks

Essential Tools for Data Scientists, Researchers, and Engineers

The MXNet logo consists of a blue circle containing a white lowercase 'm', followed by the word 'xnet' in a black sans-serif font.

MXNet is a modern open-source deep learning framework used to train and deploy deep neural networks. It is scalable, allowing for fast model training, and supports a flexible programming model and multiple languages. The MXNet Library is portable and can scale to multiple GPUs and multiple machines



TensorFlow is a popular open-source software library for **dataflow** programming across a range of tasks. It is a symbolic math library and is commonly used for deep learning applications.



Scikit-learn is a free software machine learning library for the python programming language. It features various classification, regression and clustering algorithms, and is designed to interoperate with the Python numerical and scientific libraries: NumPy and SciPy

Nvidia Deep Learning Software Stack

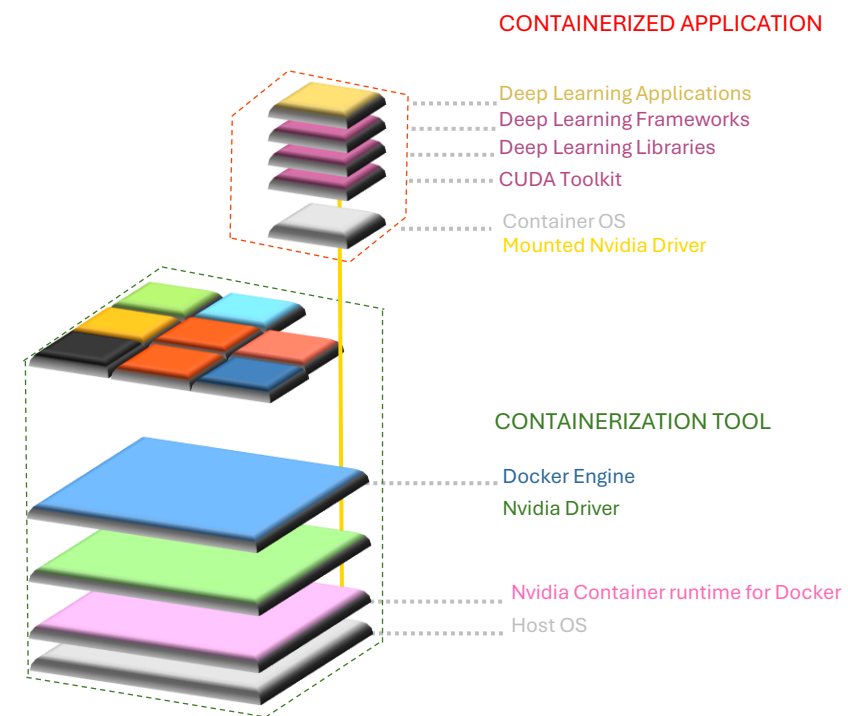
Host OS and Nvidia Driver – Enables the deep learning framework to use the GPU Functions

NGC Containers – Publicly available containers optimized to run on Nvidia GPUs

DL Frameworks – Popular deep learning frameworks available inside the containers

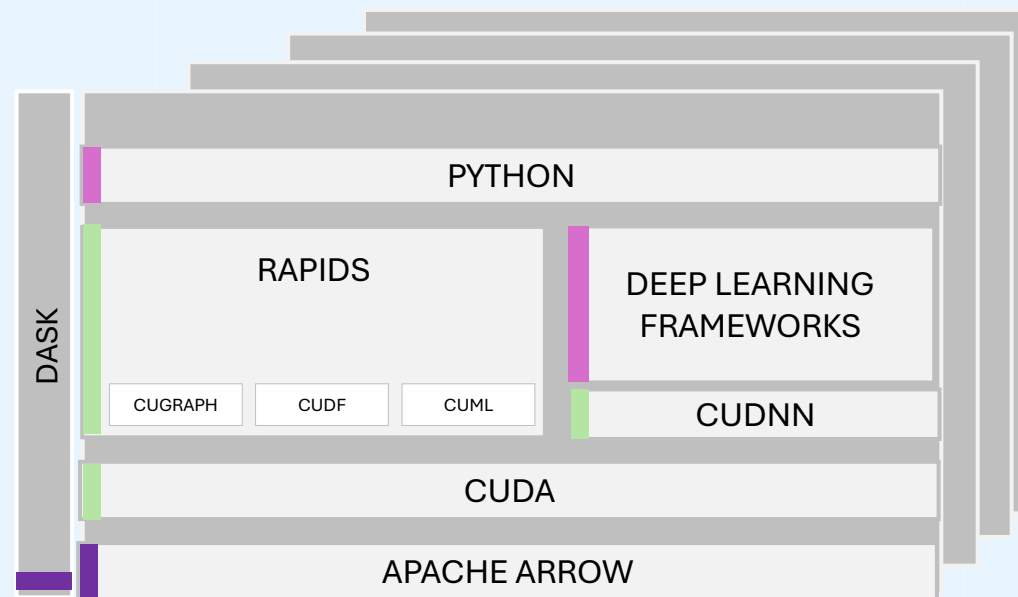
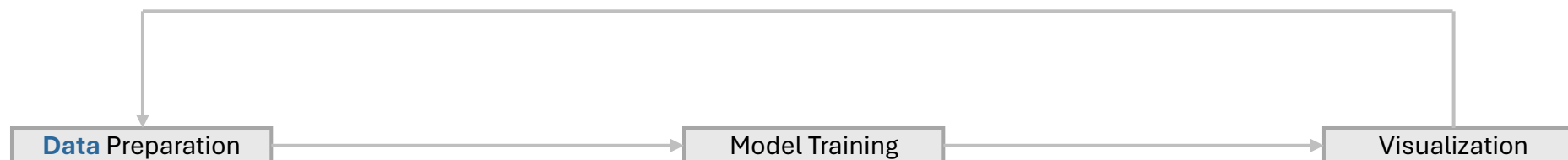
CUDA – Nvidia's ground breaking parallel programming model.

Provides essential optimizations for deep learning, machine learning, and high performance computing [HPC] leveraging Nvidia GPUs

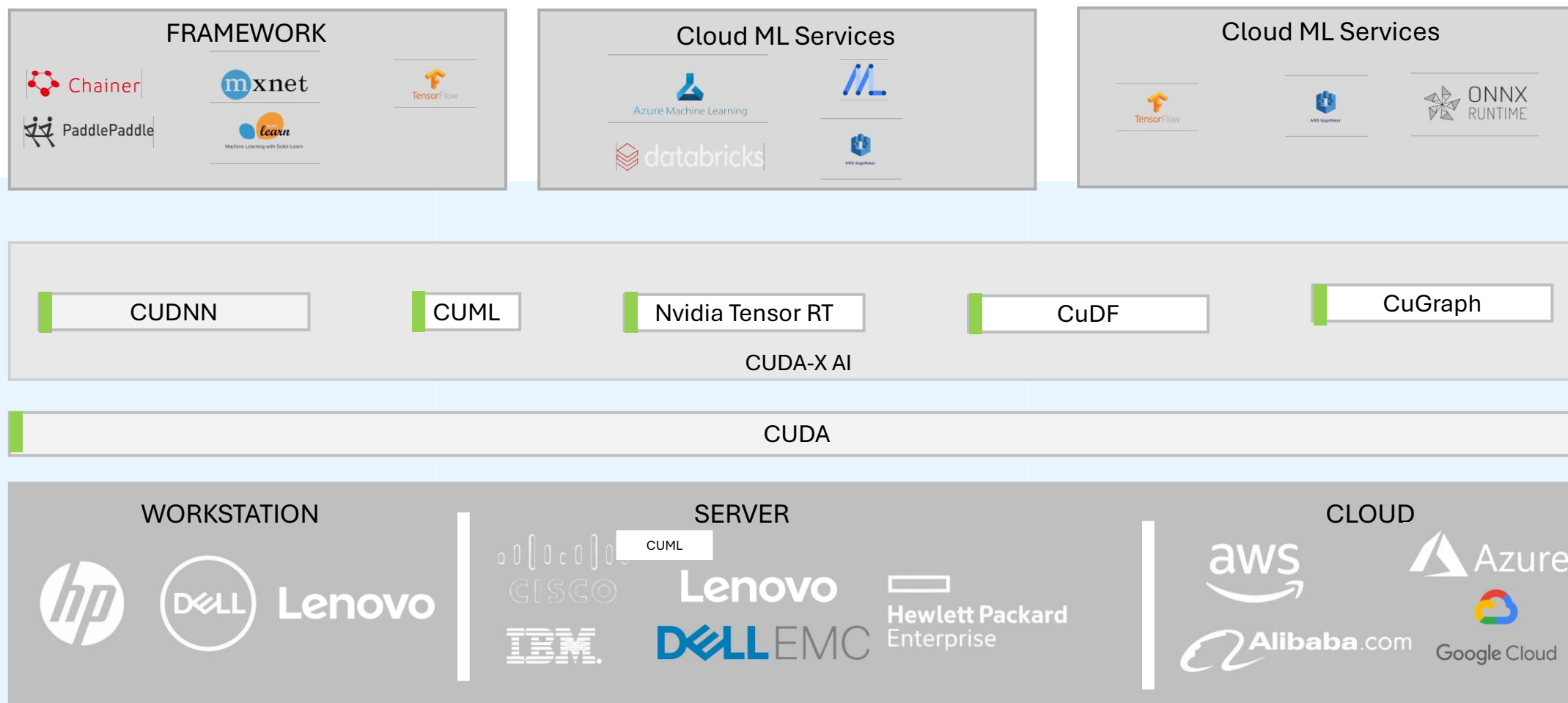


Nvidia GPU software stack

Machine Learning Software Stack



Nvidia CUDA – X AI Ecosystem




Nvidia CUDA – X AI Libraries

Data Processing	nvJPEG	https://developer.nvidia.com/nvjpeg
	DALI	https://github.com/NVIDIA/DALI
	cuDF	https://github.com/rapidsai/cudf
	OpticalFlow	https://developer.nvidia.com/opticalflow-sdf
	NPP	https://developer.nvidia.com/npp
Deep Learning Training	AMP	https://developer.nvidia.com/automatic-mixed-precision
	Apex	https://github.com/NVIDIA/apex
	cuBLAS	https://developer.nvidia.com/cublas
AI Specific Acceleration Libraries	cuDNN	https://developer.nvidia.com/cuDNN
	cuxFilter	https://github.com/rapidsai/cuxfilter
	cuML	https://github.com/rapidsai/cuml
	cuGRAPH	https://github.com/rapidsai/cudgraph
	cuTensor	https://developer.nvidia.com/cuTensor
	NCCL	https://developer.nvidia.com/nccl
Deployment	TensorRT	https://developer.nvidia.com/tensorrt
	Inference Server	https://github.com/NVIDIA/tensorrt-inference-server
High Level Constructs	Transfer Learning Tool	https://developer.nvidia.com/transfer-learning-toolkit
	Rapids	https://github.com/rapidsai
	DeepStream	https://developer.nvidia.com/deepstream-sdk

Key Technologies for Deployment

Containers



DOCKER

Package app


- Libraries
- Compilers
- Network Drivers
- Other Components

Simplifies Deployments
Eliminates complex, time-consuming builds and installs

Quick Start
Simply download and run the app

Portable
Deploy across various environments, from test to production, with minimal changes

Kubernetes




Orchestration tool to easily deploy containers on various nodes

Automates Deployments
Launch apps on appropriate nodes

Scaling
Automatically spins up nodes to meet demand

Monitoring
Automatically restarts if application crashes

Slurm



slurm
workload manager

Job scheduler to manage allocation of resources and launching jobs on a cluster

Restore Allocation
Automatically allocates resources for job

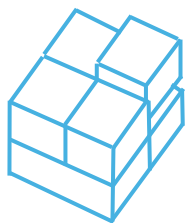
Large Clusters
Supports running on small to vary large clusters

Distributed Jobs
Allows multi-nodes jo to be launched for faster distributed training

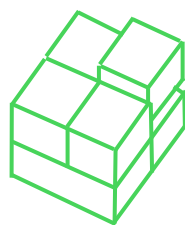
NGC Catalog – GPU Optimized Hub for AI & HPC software

Simplify and Accelerate End-to-End Workflows

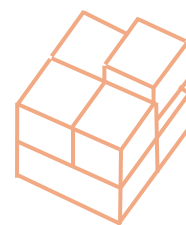
Software stacks vary End-to-End Workflow



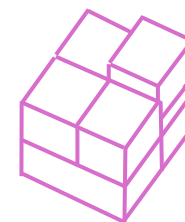
AI Stack



HPC Simulation apps



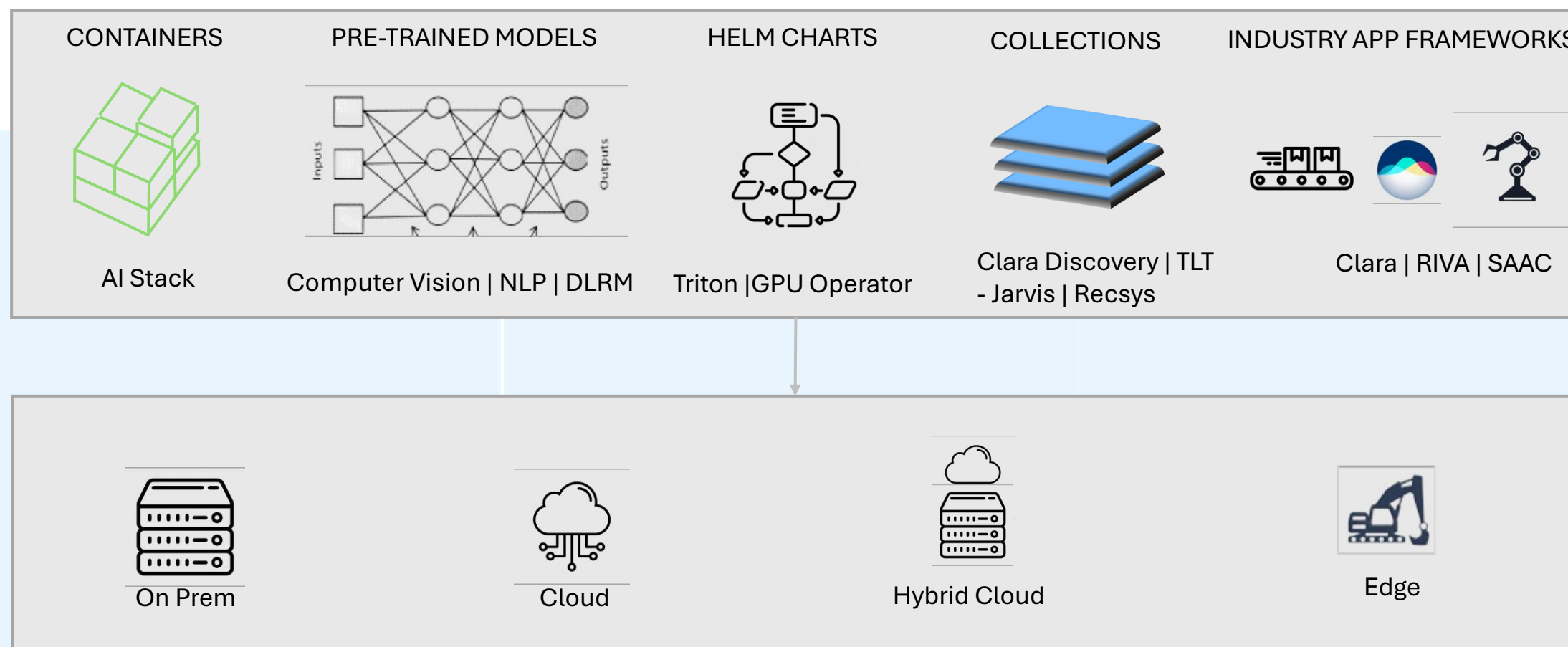
Genomics Stack



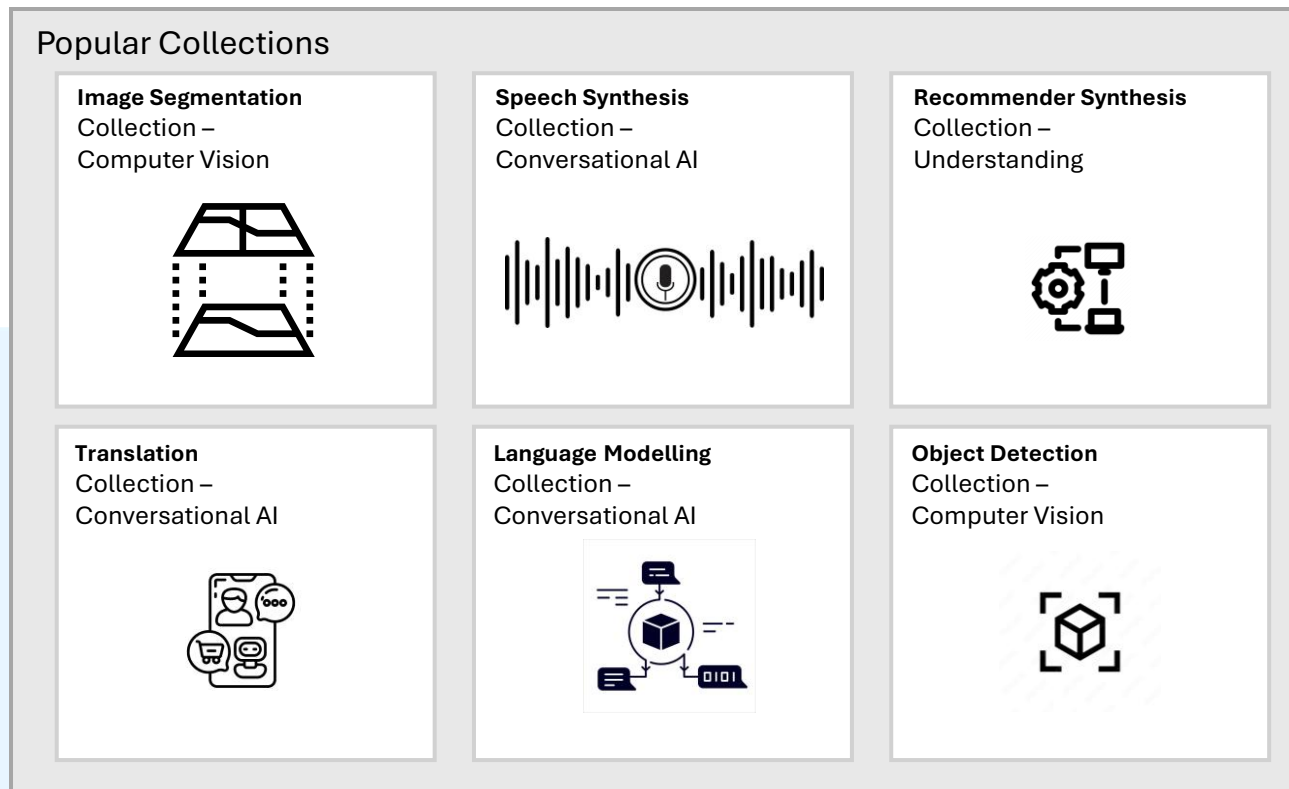
Visualization app

NGC Catalog – GPU Optimized Hub for Ai & HPC Software

Simplify and Accelerate End-to-End Workflows



Fast Track AI with Pre-Trained Models from NGC



Production Quality

- Trained and Continuously updated by experts
- Model resumes to find the right fit

Wide Range of Use Cases

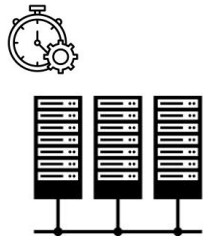


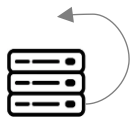
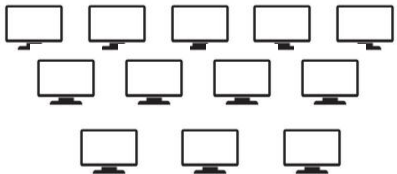

- People Detection, Vehicle Detection, Gaze Estimation
- Instant Classification, Question-Answering, Speech recognition, and Text-to-speech

Adapt & Integrate

- Adapt to your domain with your custom data
- Integrate easily into industry SDKs

Benefits of GPU Virtualization

Nvidia in Virtualization Environments – Industry Leading Innovations

<p>Bare Metal Performance</p> 	<p>Operational Management</p> 	<p>Insight & Tools</p> 
<p>Business Continuity and Workload Balancing</p> 	<p>Resource sharing and improved utilization</p> 	<p>Infrastructure & Data Security</p> 

Nvidia End-to-End AI software

Optimized for Nvidia-Certified Servers

PRE-TRAINED MODELS



Automatic Speech
Recognition



Object
Detection



Speech
Synthesis



Machine
Translation



Image
Classification



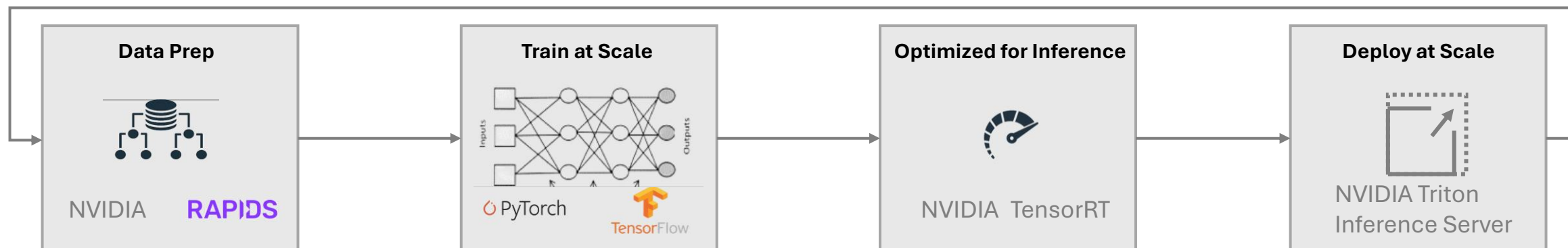
Language
Modelling



Recommender
Systems

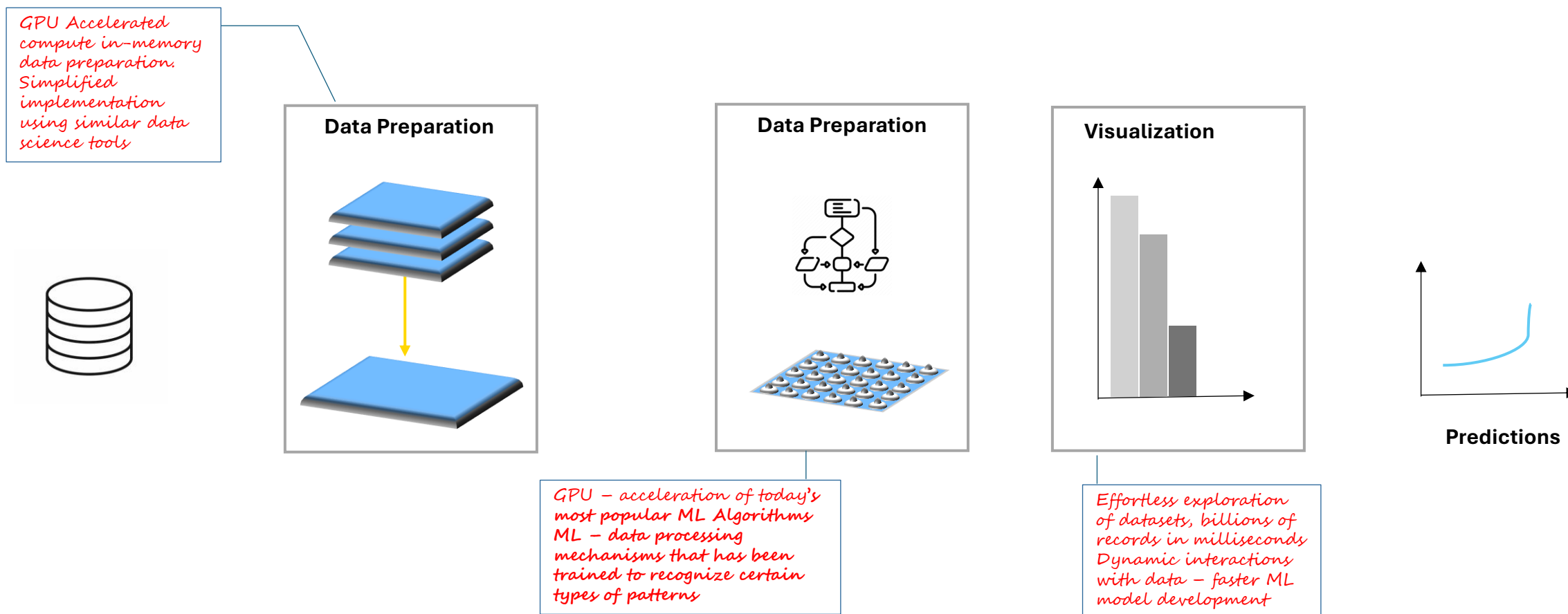


Segmentation

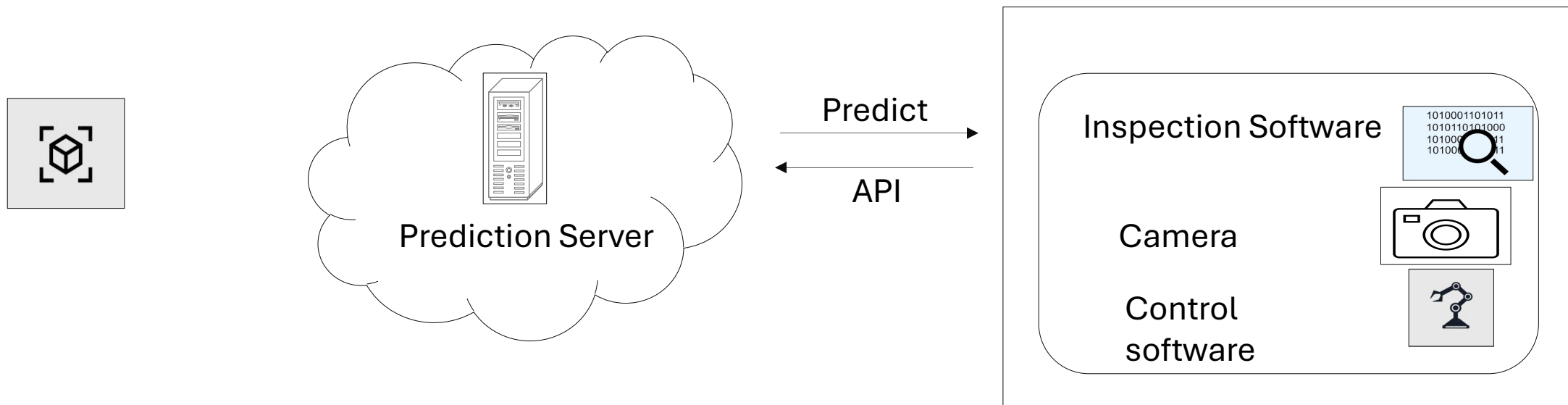


AI Workflow

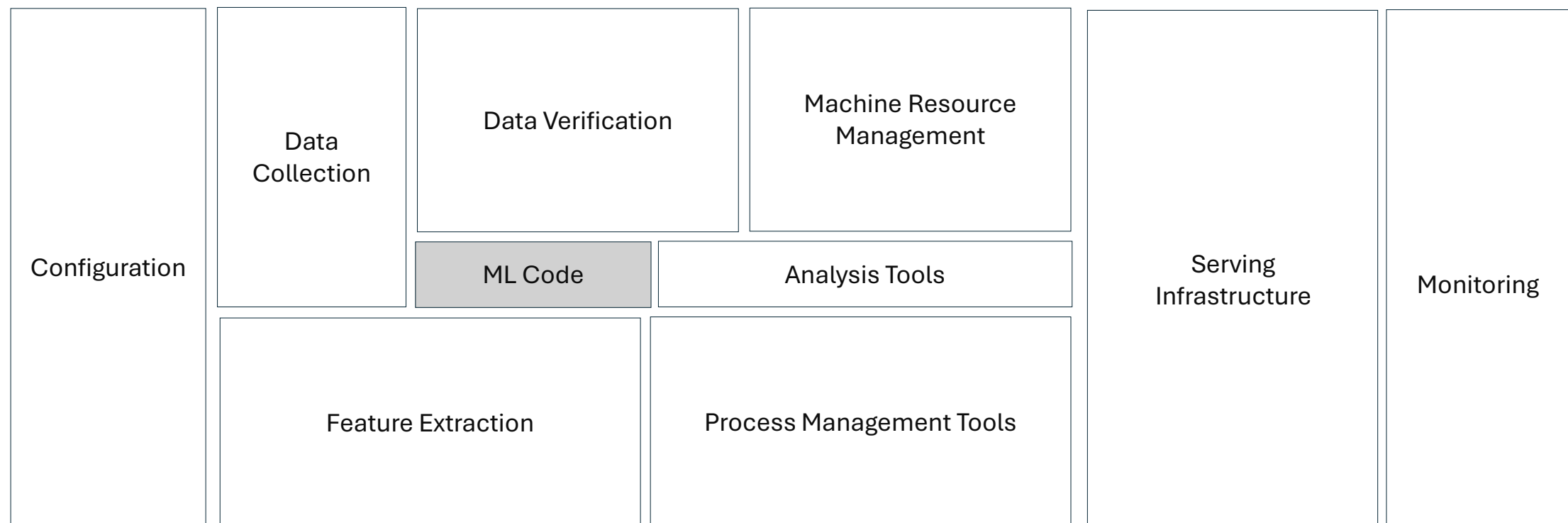
Open Source, End-to-end GPU-acclerated Workflow



The Machine Learning Project Lifecycle



The requirements surrounding ML Infrastructure



[D. Sculley et. al NIPS 2015: Hidden Technical Debt in Machine Learning Systems]

The ML Project Lifecycle

