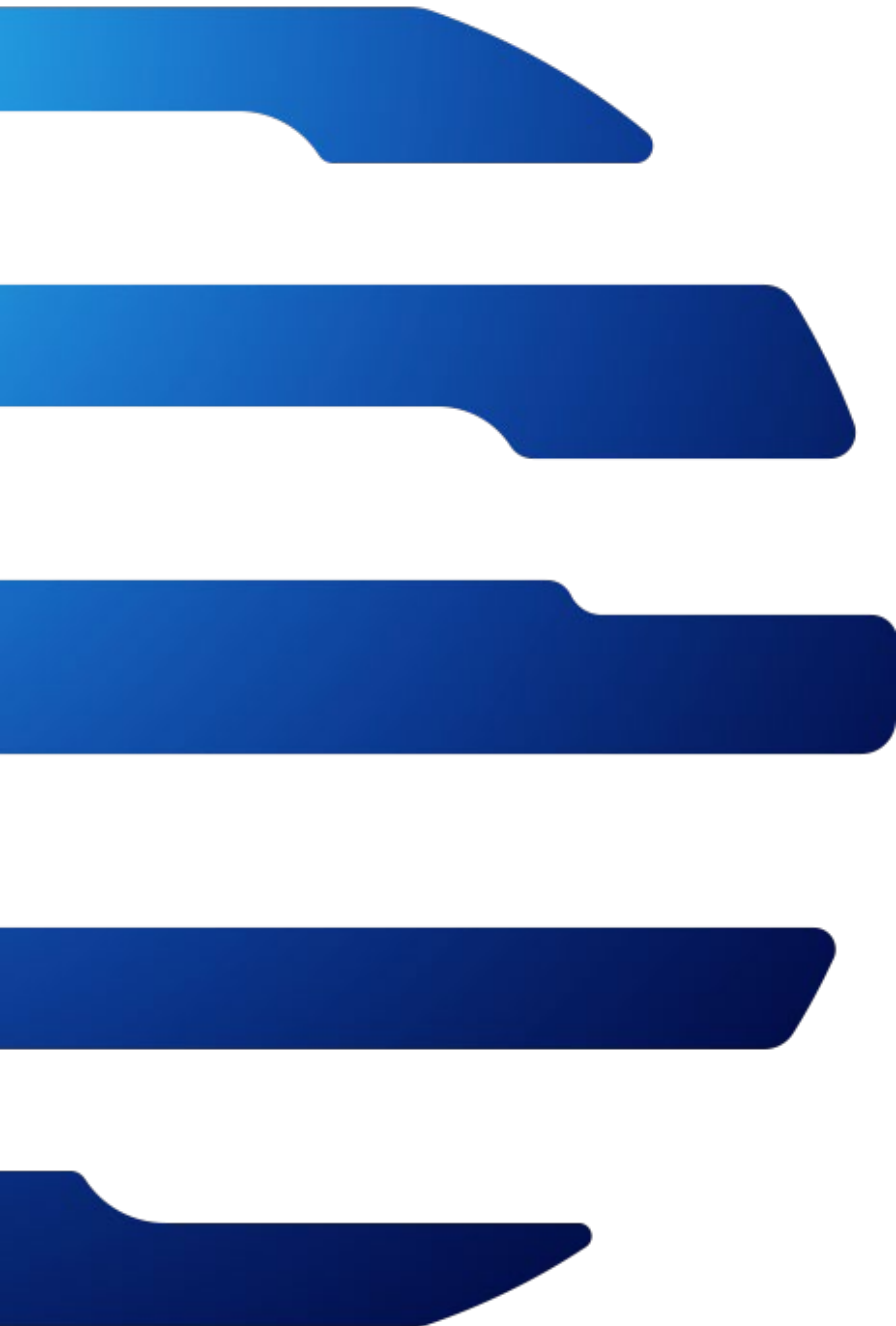# Performance Intensive Computing: Your AI-ready Infrastructure
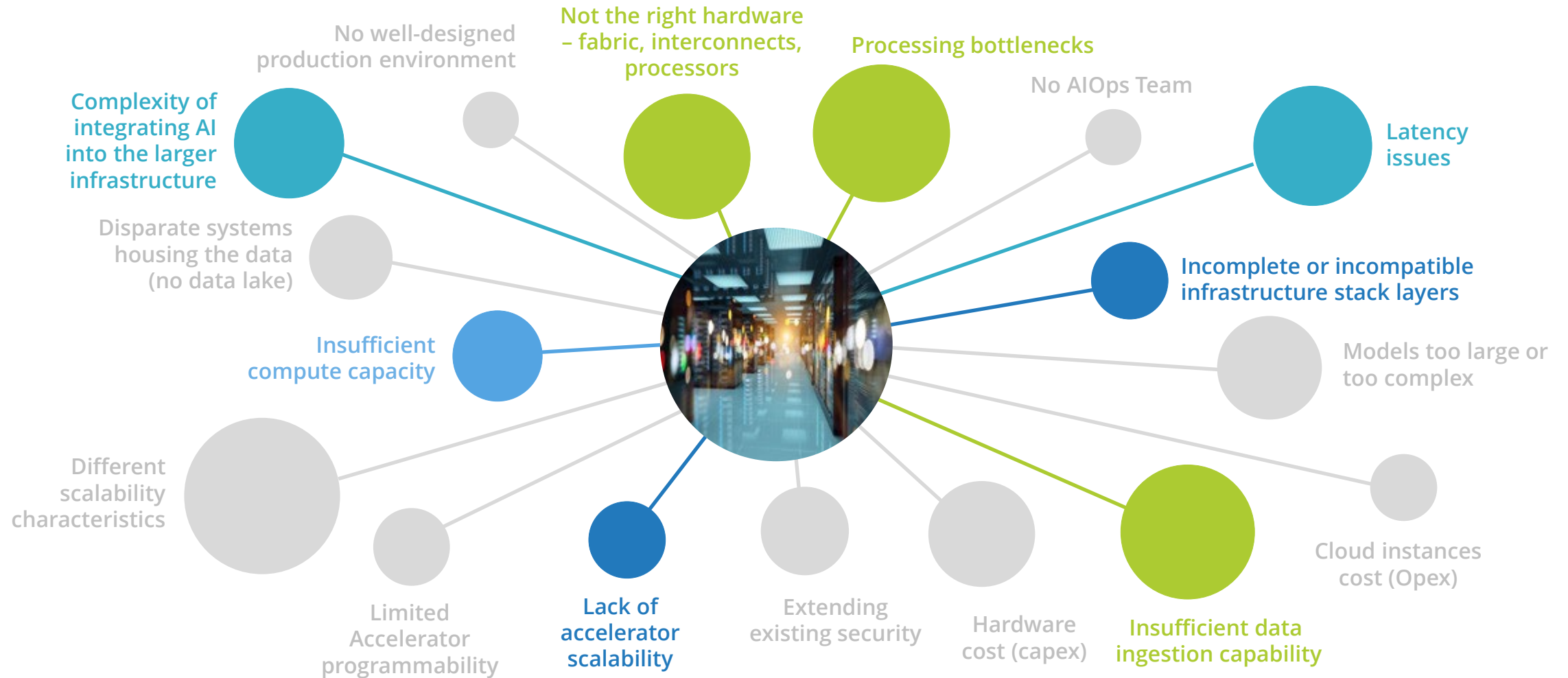
Extracting value from data in a timely manner

**Ashish Nadkarni**
**Worldwide Infrastructure Research**
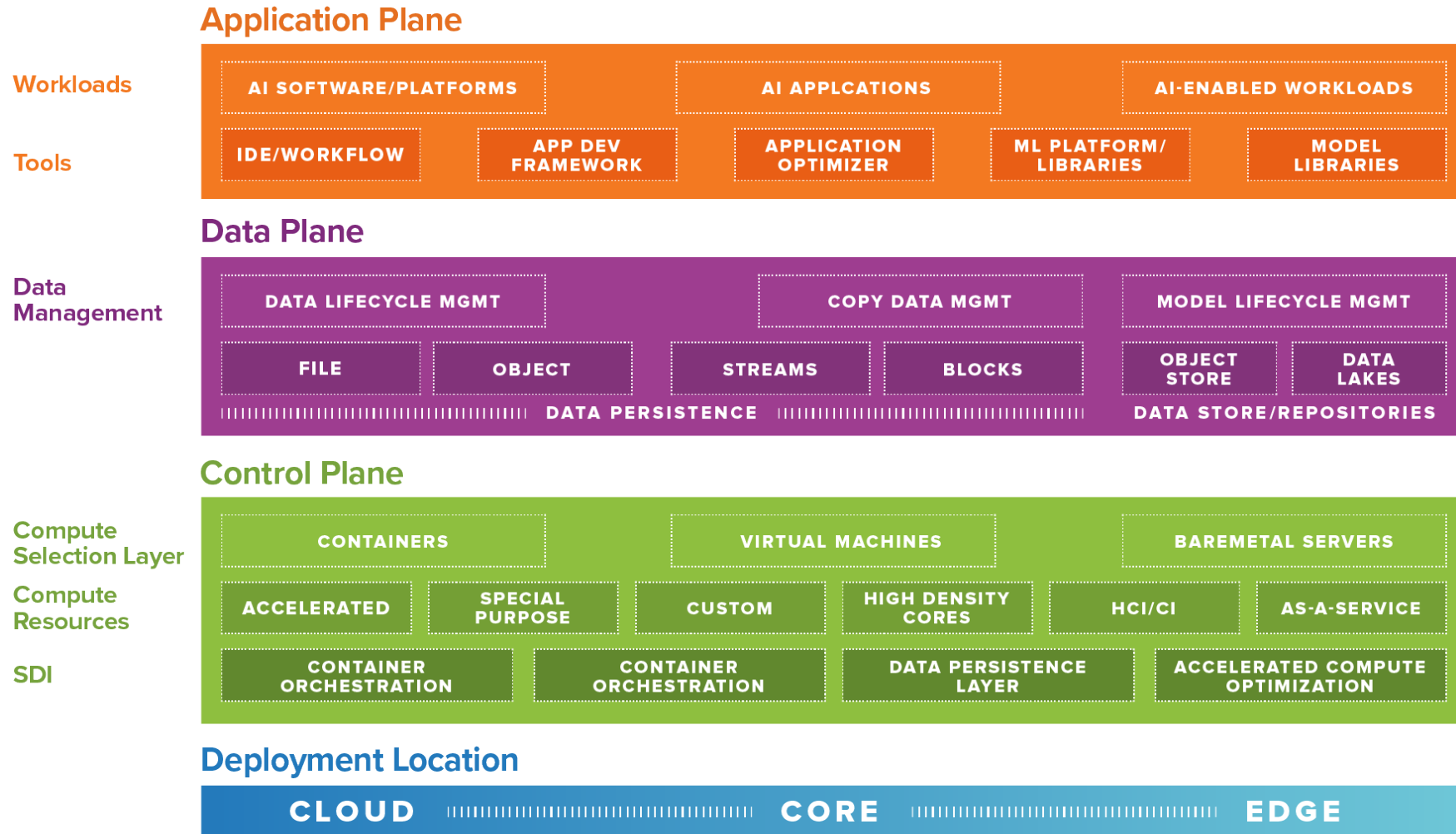
# The current state of Infrastructure for AI

# Many AI projects fail because of improper attention to infrastructure

No well-designed production environment

Not the right hardware – fabric, interconnects, processors

Processing bottlenecks

No AIOps Team

Complexity of integrating AI into the larger infrastructure

Latency issues

Disparate systems housing the data (no data lake)

Incomplete or incompatible infrastructure stack layers

Insufficient compute capacity

Models too large or too complex

Different scalability characteristics

Limited Accelerator programmability

Lack of accelerator scalability

Extending existing security

Hardware cost (capex)

Insufficient data ingestion capability

Cloud instances cost (Opex)

# A homegrown AI Infrastructure Stack is not often designed for scale

## Application Plane

**Workloads**

| AI SOFTWARE/PLATFORMS | AI APPLCATIONS | AI-ENABLED WORKLOADS |

**Tools**

| IDE/WORKFLOW | APP DEV FRAMEWORK | APPLICATION OPTIMIZER | ML PLATFORM/ LIBRARIES | MODEL LIBRARIES |

## Data Plane

**Data Management**

| DATA LIFECYCLE MGMT | COPY DATA MGMT | MODEL LIFECYCLE MGMT |

| FILE | OBJECT | STREAMS | BLOCKS | OBJECT STORE | DATA LAKES |

DATA PERSISTENCE   DATA STORE/REPOSITORIES

## Control Plane

**Compute Selection Layer**

| CONTAINERS | VIRTUAL MACHINES | BAREMETAL SERVERS |

**Compute Resources**

| ACCELERATED | SPECIAL PURPOSE | CUSTOM | HIGH DENSITY CORES | HCI/CI | AS-A-SERVICE |

**SDI**

| CONTAINER ORCHESTRATION | CONTAINER ORCHESTRATION | DATA PERSISTENCE LAYER | ACCELERATED COMPUTE OPTIMIZATION |

## Deployment Location
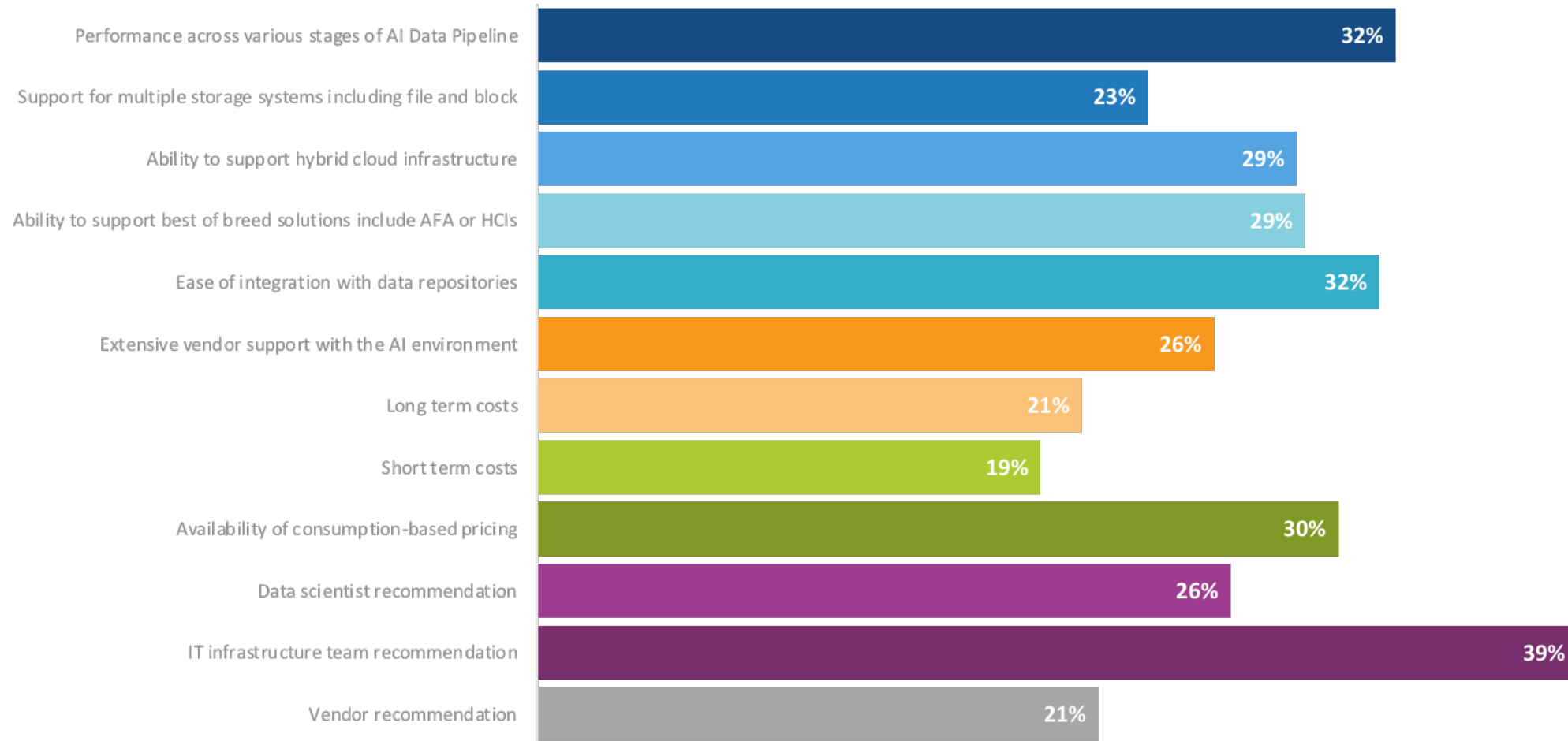
**CLOUD**   **CORE**   **EDGE**

IDC

# Lack of proper computing considerations can cause problems with outcomes

Less than half (45%) of the compute infrastructure currently used for running AI workloads is virtualized or a combination of virtualized and containerized. This is expected to increase to 74% in the next 18 months

Bare metal
40%

Virtualized
33%

Both
12%

Containerized
37%

Today

Bare metal
30%

Virtualized
30%

Both
44%

Containerized
27%

In 18 months

# Storage Infrastructure requirements are often not well understood



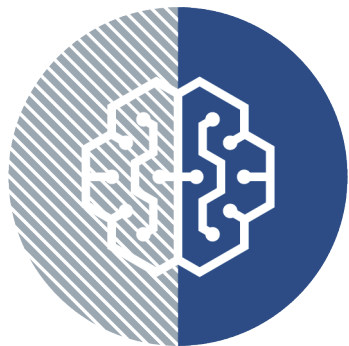| Requirement | Percentage |
|---|---|
| Performance across various stages of AI Data Pipeline | 32% |
| Support for multiple storage systems including file and block | 23% |
| Ability to support hybrid cloud infrastructure | 29% |
| Ability to support best of breed solutions include AFA or HCIs | 29% |
| Ease of integration with data repositories | 32% |
| Extensive vendor support with the AI environment | 26% |
| Long term costs | 21% |
| Short term costs | 19% |
| Availability of consumption-based pricing | 30% |
| Data scientist recommendation | 26% |
| IT infrastructure team recommendation | 39% |
| Vendor recommendation | 21% |

Q. What were the key requirements when selecting storage infrastructure for AI in you datacenter/colocation provider/edge location?

# Artificial intelligence is not always standalone; It is being infused into Enterprise Applications

## 35%

**The average number of applications using some form of AI/ML or DL TODAY**
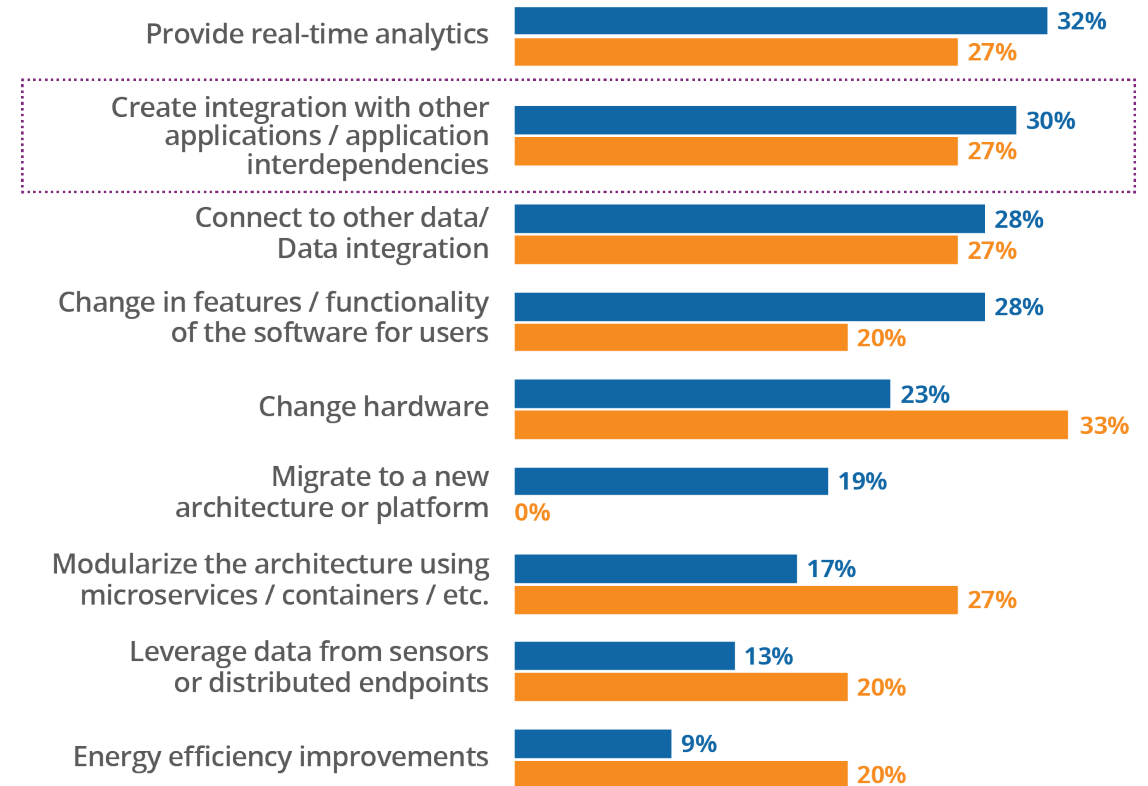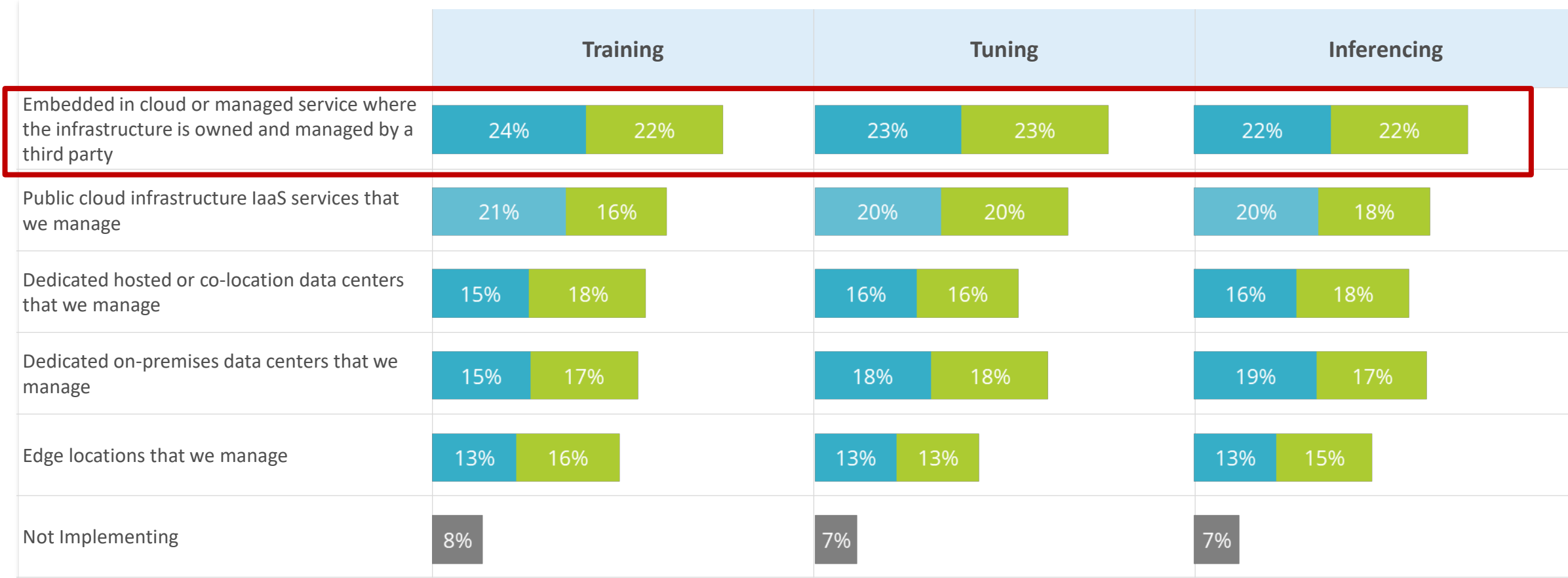
## 50%

**IN TWO YEARS**

AI applications will be integrated with other applications across the cloud portfolio

**What type of transformation do you expect for your**
**■ AI software services and     ■ AI lifecycle applications?**

| Transformation | AI software services | AI lifecycle applications |
|---|---|---|
| Provide real-time analytics | 32% | 27% |
| Create integration with other applications / application interdependencies | 30% | 27% |
| Connect to other data/ Data integration | 28% | 27% |
| Change in features / functionality of the software for users | 28% | 20% |
| Change hardware | 23% | 33% |
| Migrate to a new architecture or platform | 19% | 0% |
| Modularize the architecture using microservices / containers / etc. | 17% | 27% |
| Leverage data from sensors or distributed endpoints | 13% | 20% |
| Energy efficiency improvements | 9% | 20% |

# Public Cloud is the de facto approach for a range of Generative AI activities (but it may not always be the best one)

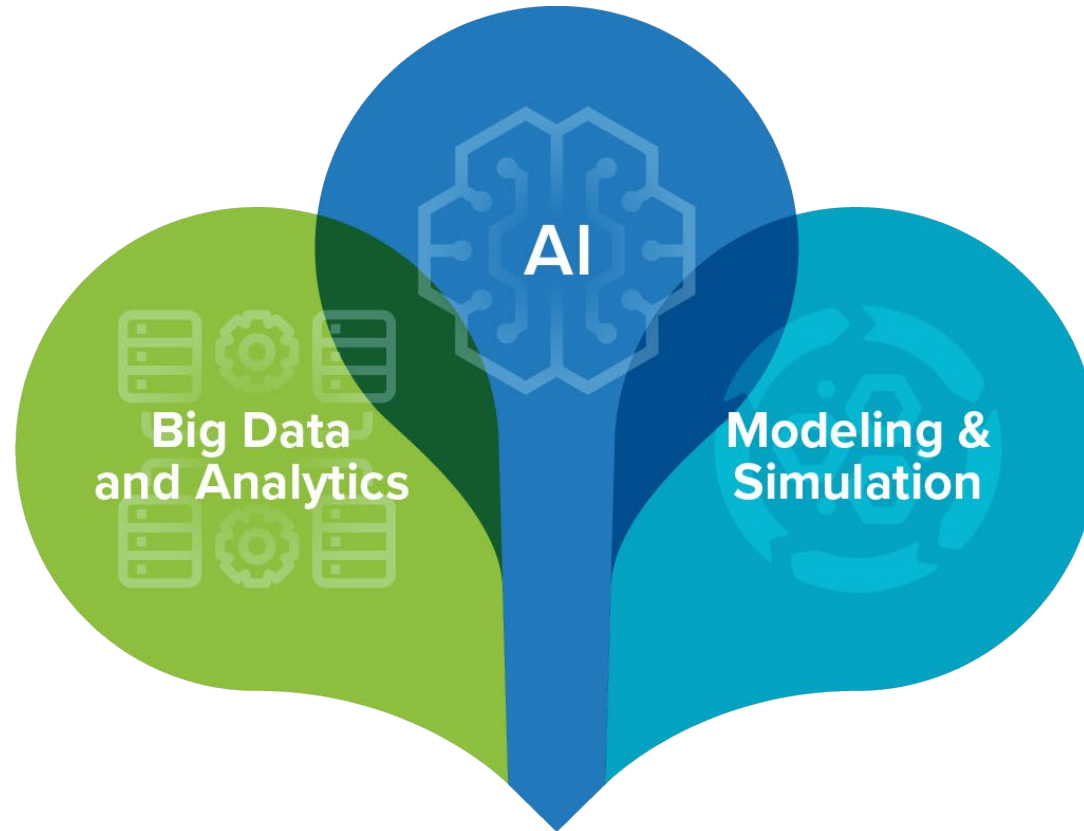| Over the next 18 months, what will be the primary approach which your organization deploys and manages infrastructure resources for GenAI? What will be your secondary approach? | Training | | Tuning | | Inferencing | |
|---|---|---|---|---|---|---|
| Embedded in cloud or managed service where the infrastructure is owned and managed by a third party | 24% | 22% | 23% | 23% | 22% | 22% |
| Public cloud infrastructure IaaS services that we manage | 21% | 16% | 20% | 20% | 20% | 18% |
| Dedicated hosted or co-location data centers that we manage | 15% | 18% | 16% | 16% | 16% | 18% |
| Dedicated on-premises data centers that we manage | 15% | 17% | 18% | 18% | 19% | 17% |
| Edge locations that we manage | 13% | 16% | 13% | 13% | 13% | 15% |
| Not Implementing | 8% | | 7% | | 7% | |

■ Primary approach   ■ Secondary approach

# Considering an AI Ready Infrastructure

# IDC is seeing a Convergence of Three "Workload Groups" onto One Infrastructure Approach

**AI**

**Big Data and Analytics**

**Modeling & Simulation**

Borrowing scaling approaches from HPC (modeling and simulation) infrastructure for AI workloads

## Performance-Intensive Computing Infrastructure

IDC

# What are Performance Intensive Computing Workloads?

## IDC defines workloads are applications and their associated datasets

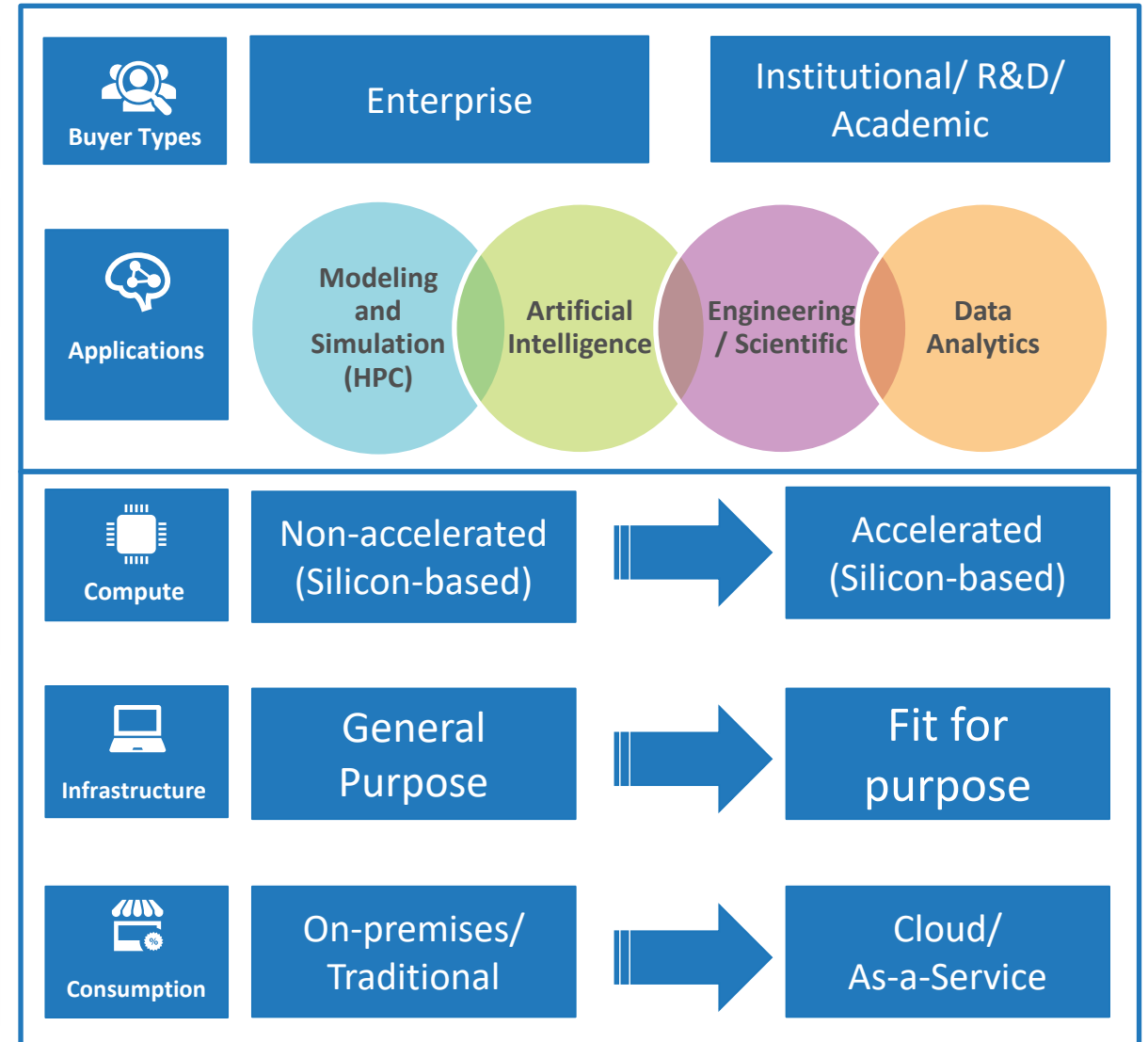Workloads that **perform large-scale mathematically intensive computations**

Workloads that **process large volumes of data**

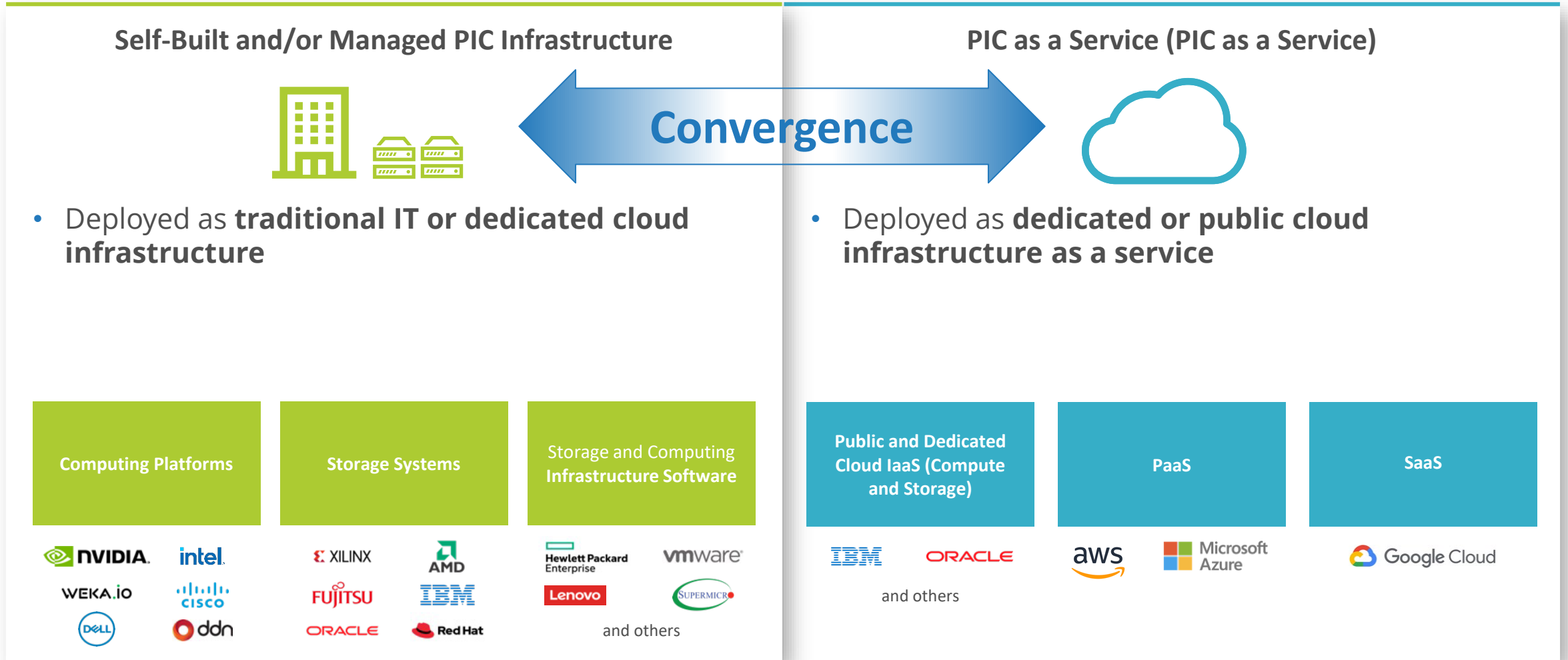Workloads that **complex instruction sets to be executed in the shortest amount of time**

Workloads that **are deployed with compressed time-to-insights objectives**

Use cases

- **Artificial Intelligence and Machine Learning (AI/ML)**
- **Modeling and simulation (M&S)**
- **Big Data and Analytics (BDA)**
- **Engineering, technical and industry specific**

| Buyer Types | Enterprise | Institutional/ R&D/ Academic |
|---|---|---|
| Applications | Modeling and Simulation (HPC) · Artificial Intelligence · Engineering / Scientific · Data Analytics | |

| | | |
|---|---|---|
| Compute | Non-accelerated (Silicon-based) | ➡ Accelerated (Silicon-based) |
| Infrastructure | General Purpose | ➡ Fit for purpose |
| Consumption | On-premises/ Traditional | ➡ Cloud/ As-a-Service |

# Two Principal Deployments, or as a Hybrid Approach

**Self-Built and/or Managed PIC Infrastructure**

**PIC as a Service (PIC as a Service)**

**Convergence**

- Deployed as **traditional IT or dedicated cloud infrastructure**

- Deployed as **dedicated or public cloud infrastructure as a service**

| Computing Platforms | Storage Systems | Storage and Computing Infrastructure Software |
|---|---|---|

| Public and Dedicated Cloud IaaS (Compute and Storage) | PaaS | SaaS |
|---|---|---|

**NVIDIA**   **intel**    **XILINX**   **AMD**    Hewlett Packard Enterprise   **vmware**

**WEKA.io**   **cisco**    **FUJITSU**   **IBM**    **Lenovo**   **SUPERMICRO**

**DELL**   **ddn**    **ORACLE**   **Red Hat**    and others

**IBM**    **ORACLE**     **aws**   **Microsoft Azure**     **Google Cloud**

and others

# A Framework for selecting the right infrastructure stack
# Part 1 – Cultural/ Deployment (What and How)

| For Lower Costs >>> | Prefer on-premises or Collocation (Private Cloud) if | Prefer off-premises (Public Cloud) if |
|---|---|---|
| **AI Initiatives** | On-going AI initiatives from a significantly busy team | On- and off AI initiatives |
| **System utilization** | The ability to keep utilization rates very high (keep expensive processors busy) | No ability to keep utilization rates high |
| **IT Skills** | Possess in-house skills for for complex AI deployments | Limited IT skills for AI deployments, leave alone complex |
| **Facilities** | No limitations on datacenter floorspace, power, and cooling capabilities | Limited floor space, power, and cooling |
| **Opex friendly options** | System vendor can provide consumption-based pricing | System vendor can provide capital only pricing |

# A Framework for selecting the right infrastructure stack
## Part 2 – Use case specific (Model considerations)

| For Lower Costs >>> | Prefer on-premises or Collocation (Private Cloud) if | Prefer off-premises (Public Cloud) if |
|---|---|---|
| **Model Iteration** | Many model training iterations | Fewer model iterations |
| **Model Scaling** | High scaling needs | Lower scaling needs |
| **Model Accuracy** | Highly customized | Little to no customization |
| **Model customization** | Heavily customized or Inference only | No API changes or customization |
| **Model Performance** | High performance requirements | Lower performance requirements |

# A Framework for selecting the right infrastructure stack
# Part 3 – Use case specific (Data considerations)

| For Lower Costs >>> | Prefer on-premises or Collocation (Private Cloud) if | Prefer off-premises (Public Cloud) if |
|---|---|---|
| **Data sensitivity** | Highly sensitive data, strict data compliance requirements, proprietary data | Data is not proprietary, no compliance requirements or has been completely sanitized |
| **Data isolation** | Model data <u>cannot</u> mix with public data, requires isolation | Model data Data can safely mix with public data, does not require isolation |
| **Time to Value** | Not time critical | Highly time critical |

# When choosing a Cloud Provider for AI
## Building differentiation in AI Infrastructure will be the next battle ground for cloud providers

**New architectures**
Internally developed and 3rd party CPUs and GPUs

**Edge computing**
Extending cloud services to smaller, remote locations

**Sovereign clouds**
Addressing regulatory compliance for data and operations

**Multicloud integration**
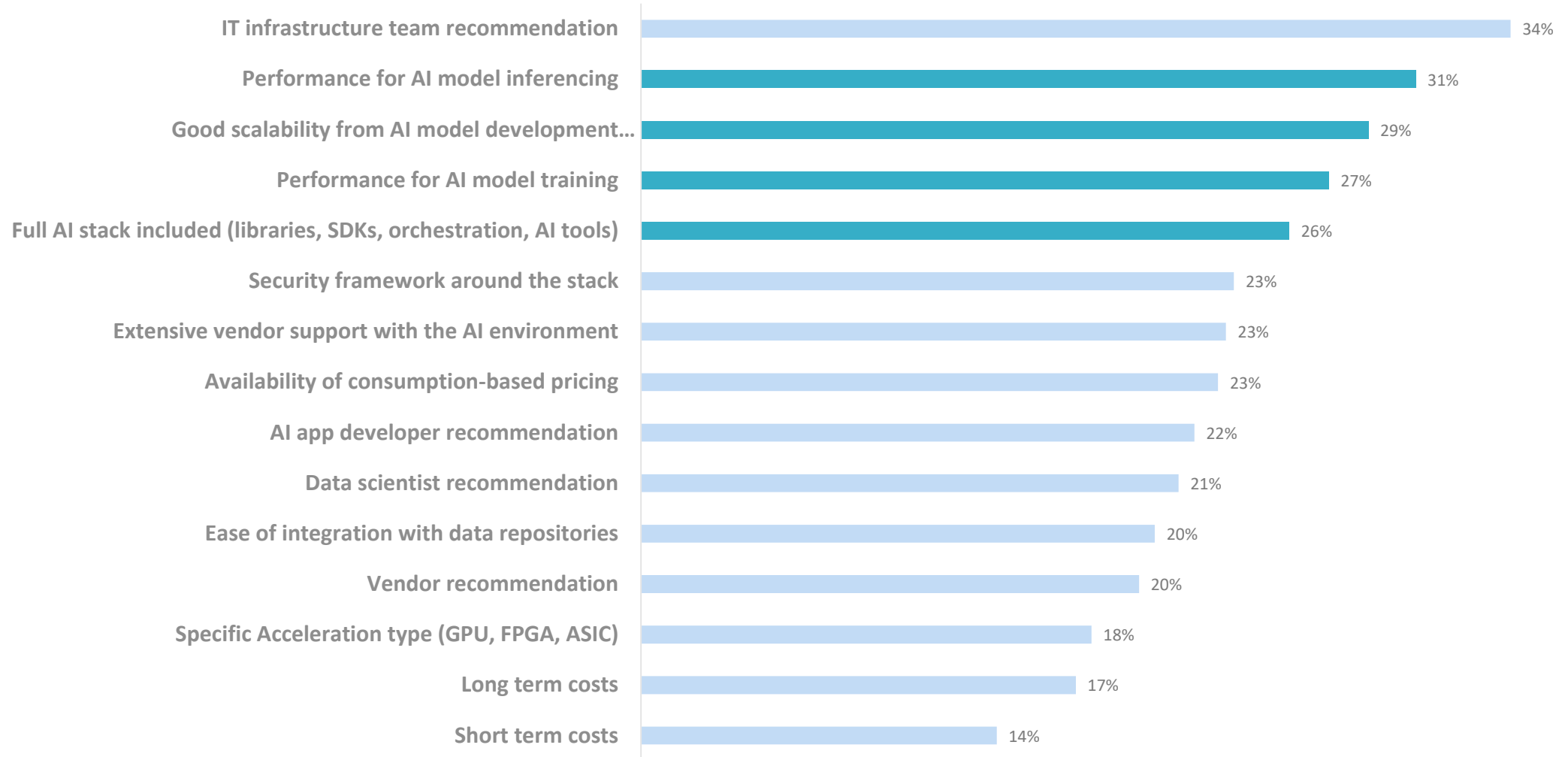Tools and commercial terms to facilitate management and security

**High performance**
Compute and storage services optimized for HPC and AI

**Deeper partnerships**
Creating bridges with traditional enterprise OEMs and ISVs.

# Key requirements when selecting compute infrastructure for AI in datacenters, at colocation providers and edge locations

| Requirement | Percentage |
|---|---|
| IT infrastructure team recommendation | 34% |
| Performance for AI model inferencing | 31% |
| Good scalability from AI model development... | 29% |
| Performance for AI model training | 27% |
| Full AI stack included (libraries, SDKs, orchestration, AI tools) | 26% |
| Security framework around the stack | 23% |
| Extensive vendor support with the AI environment | 23% |
| Availability of consumption-based pricing | 23% |
| AI app developer recommendation | 22% |
| Data scientist recommendation | 21% |
| Ease of integration with data repositories | 20% |
| Vendor recommendation | 20% |
| Specific Acceleration type (GPU, FPGA, ASIC) | 18% |
| Long term costs | 17% |
| Short term costs | 14% |

# Key requirements when selecting <u>storage</u> infrastructure for AI in datacenters, at colocation providers and edge locations

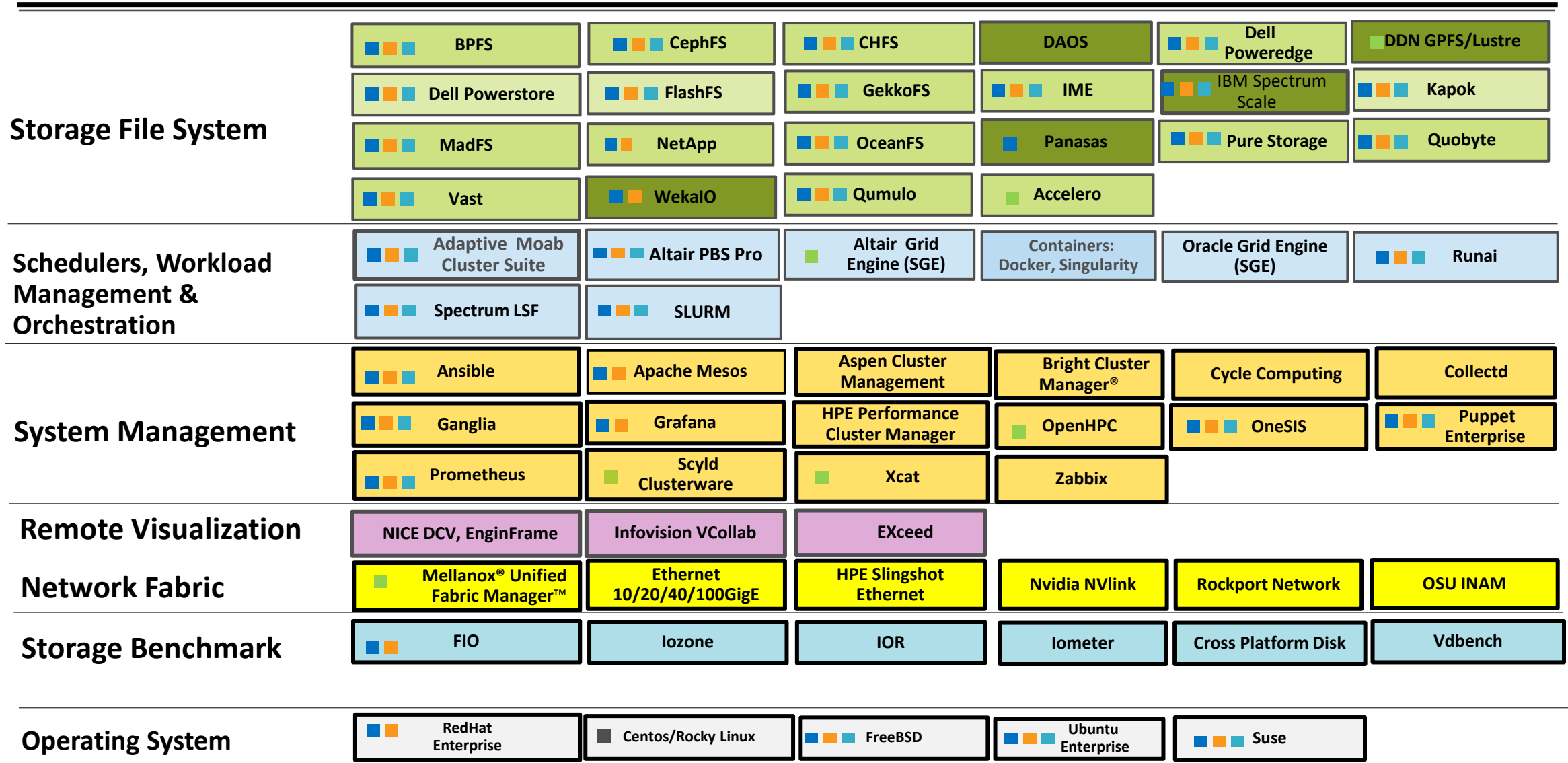| Requirement | Percentage |
|---|---|
| IT infrastructure team recommendation | 36% |
| Good scalability from AI model development to production | 29% |
| Availability of consumption-based pricing | 27% |
| Performance for AI model inferencing | 26% |
| Data scientist recommendation | 25% |
| Extensive vendor support with the AI environment | 25% |
| Performance for AI model training | 25% |
| Specific Acceleration type (GPU, FPGA, ASIC) | 24% |
| Ease of integration with data repositories | 24% |
| Security framework around the stack | 23% |
| Full AI stack included (libraries, SDKs, orchestration, AI tools) | 22% |
| AI app developer recommendation | 19% |
| Long term costs | 19% |
| Vendor recommendation | 19% |
| Short term costs | 17% |

# A typical Software Stack



| Storage File System | | | | | |
|---|---|---|---|---|---|
| BPFS | CephFS | CHFS | DAOS | Dell Poweredge | DDN GPFS/Lustre |
| Dell Powerstore | FlashFS | GekkoFS | IME | IBM Spectrum Scale | Kapok |
| MadFS | NetApp | OceanFS | Panasas | Pure Storage | Quobyte |
| Vast | WekaIO | Qumulo | Accelero | | |

| Schedulers, Workload Management & Orchestration | | | | | |
|---|---|---|---|---|---|
| Adaptive Moab Cluster Suite | Altair PBS Pro | Altair Grid Engine (SGE) | Containers: Docker, Singularity | Oracle Grid Engine (SGE) | Runai |
| Spectrum LSF | SLURM | | | | |

| System Management | | | | | |
|---|---|---|---|---|---|
| Ansible | Apache Mesos | Aspen Cluster Management | Bright Cluster Manager® | Cycle Computing | Collectd |
| Ganglia | Grafana | HPE Performance Cluster Manager | OpenHPC | OneSIS | Puppet Enterprise |
| Prometheus | Scyld Clusterware | Xcat | Zabbix | | |

| Remote Visualization | | | |
|---|---|---|---|
| NICE DCV, EnginFrame | Infovision VCollab | EXceed | |

| Network Fabric | | | | | |
|---|---|---|---|---|---|
| Mellanox® Unified Fabric Manager™ | Ethernet 10/20/40/100GigE | HPE Slingshot Ethernet | Nvidia NVlink | Rockport Network | OSU INAM |

| Storage Benchmark | | | | | |
|---|---|---|---|---|---|
| FIO | Iozone | IOR | Iometer | Cross Platform Disk | Vdbench |

| Operating System | | | | |
|---|---|---|---|---|
| RedHat Enterprise | Centos/Rocky Linux | FreeBSD | Ubuntu Enterprise | Suse |

■ HPE Apollo, Cray, SGI    ■ NVIDIA    ■ ARM    ■ Groq    ■ Cerebras    ■ Graphcore

# A typical AI Applications Stack

| | | | | | |
|---|---|---|---|---|---|
| **Application Software Development Ecosystem** | **Programming Environment** | AOCL | Boost | CMake | Cuda | DDT, MAP Performance Report | Deal.ii |
| | | FFTW | GNU Compiler Collection | Gurobi | Java, Perl. Python (numpy, scipy) | Julia | MKL |
| | | Octave | OpenBLAS | OpenCV | OpenBLAS | R | Tensorflow |
| | | Tensorflow | Theano | oneAPI | VTune | PGI | NAG |

| **Parallel Programming Environment** | | | | | |
|---|---|---|---|---|---|
| MVAPICH2 | MVAPICH2-Azure | MVAPICH2-X/AWS | MVAPICH2-GDR | MVAPICH2-Virt | MVAPICH2-EA |
| OMB | OpenMPI | Platform MPI | Intel MPI | HPE MPI | Cray MPI |

## Data Management

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Aspera | HDF5/PHDF5 | Md5deep Hashdeep | Ncdu | OnDemand | Rclone | Spark | Visidata | Zoltan |

## Visualization

| | | | | | |
|---|---|---|---|---|---|
| Ascinema | Circos | OpenCV | Exceed | Infovision Vcollab | Lowchart |
| Nice DCV Enginframe | Xfig Fig2dev | Visit | | | |

## Productivity

| | | | | | |
|---|---|---|---|---|---|
| Git | GNU Parallel | Jupyter Notebook | Lazygit | OnDemand | Tmux | OEMT |

## AI/Deep Learning

| | | | | | |
|---|---|---|---|---|---|
| Anaconda | Jupyter | Caffe | CUDA | CuDNN | DeepGraph |
| Deepstream | Gensynth | Keras | MXNet | PaddlePaddle | Pytorch |
| Runai | Singularity | Scikit-learn | TensorRT | Tensorflow | Theano |

■ HPE Apollo, Cray, SGI  ■ NVIDIA  ■ Intel  ■ AMD  ■ ARM

# Partnering for success

## Urgency to respond quickly to business disruption at a corporate level is influencing partner selection

**Strategic Generative AI technology partners in the next 18 months**

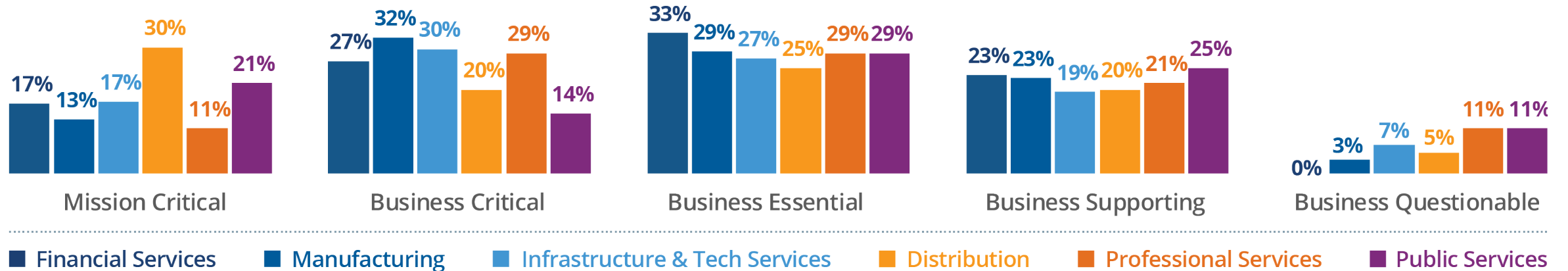# Investing in the right partner when building an AI infrastructure stack

*For areas on which IDC publishes market share data, the top 3–5 market share leaders are represented. For areas on which IDC does not publish market share data, vendor selection is up to analyst discretion.*

# Consider using AI software services

Pre-trained standalone services that provide capabilities based on machine learning, deep learning, and other AI/ML technologies for applications and workflows to help improve business outcomes.
Used to build AI-powered applications.

How would you describe AI Software Services in relation to the application's impact on your business today?



**Mission Critical**
- Financial Services: 17%
- Manufacturing: 13%
- Infrastructure & Tech Services: 17%
- Distribution: 30%
- Professional Services: 11%
- Public Services: 21%

**Business Critical**
- Financial Services: 27%
- Manufacturing: 32%
- Infrastructure & Tech Services: 30%
- Distribution: 20%
- Professional Services: 29%
- Public Services: 14%

**Business Essential**
- Financial Services: 33%
- Manufacturing: 29%
- Infrastructure & Tech Services: 27%
- Distribution: 25%
- Professional Services: 29%
- Public Services: 29%

**Business Supporting**
- Financial Services: 23%
- Manufacturing: 23%
- Infrastructure & Tech Services: 19%
- Distribution: 20%
- Professional Services: 21%
- Public Services: 25%

**Business Questionable**
- Financial Services: 0%
- Manufacturing: 3%
- Infrastructure & Tech Services: 7%
- Distribution: 5%
- Professional Services: 11%
- Public Services: 11%

Legend: ■ Financial Services ■ Manufacturing ■ Infrastructure & Tech Services ■ Distribution ■ Professional Services ■ Public Services

# In Closing...
## Consider an AI Center for Excellence to accelerate maturity of AI adoption



| Category | Percentage |
|---|---|
| Aligned to business goals, as well as redesigned business models repeatedly creating business value | 29% |
| Enterprise-wide artificial intelligence strategy aligned to business goals in place | 27% |
| Multiple projects with data readiness, governance, skills mgmt, & technology selection replicated | 20% |
| Use AI technologies for isolated projects with some level of coordination between them | 14% |
| Silos via select individuals or groups without any formal strategy or coordination as part of a broader vision | 12% |

Thank you!

Ashish Nadkarni
ANadkarni@idc.com