# What is GenAI spending by industry expected to be in 2024?

**Worldwide Core IT Spending for Generative AI by Industry will be Over $40 Billion in 2024**



$B

| Industry | Spending ($B) |
|---|---|
| Banking | ~6.5 |
| Retail | ~5.0 |
| Professional Services | ~4.95 |
| Discrete Manufacturing | ~3.3 |
| Process Manufacturing | ~2.7 |
| Insurance | ~2.35 |
| Securities and Investment Services | ~2.2 |
| Telecommunications | ~2.0 |
| Healthcare Provider | ~1.65 |
| Federal/Central Government | ~1.35 |
| Media | ~1.25 |
| Transportation | ~1.25 |
| State/Local Government | ~1.15 |
| Utilities | ~0.95 |
| Wholesale | ~0.85 |
| Personal and Consumer Services | ~0.75 |
| Education | ~0.7 |
| Resource Industries | ~0.65 |
| Construction | ~0.4 |

# Generative AI blunders (just some)

**Google Pauses Gemini AI Model… (due to inaccurate historical images)**
source: Forbes February 2024

**Air Canada ordered to pay customer who was misled by airline's chatbot**
source: The Guardian February 2024

**ChatGPT falsely accused a law professor of harassing a student**
source: The Washington Post, April 2023

**Proprietary data was leaked by Samsung employees to ChatGPT**
source: CIO Dive, April 2023

**ChatGPT bug in mid-March exposed client conversations to other users**
source: Vilius Petkauskas, Cybernews, March 2023
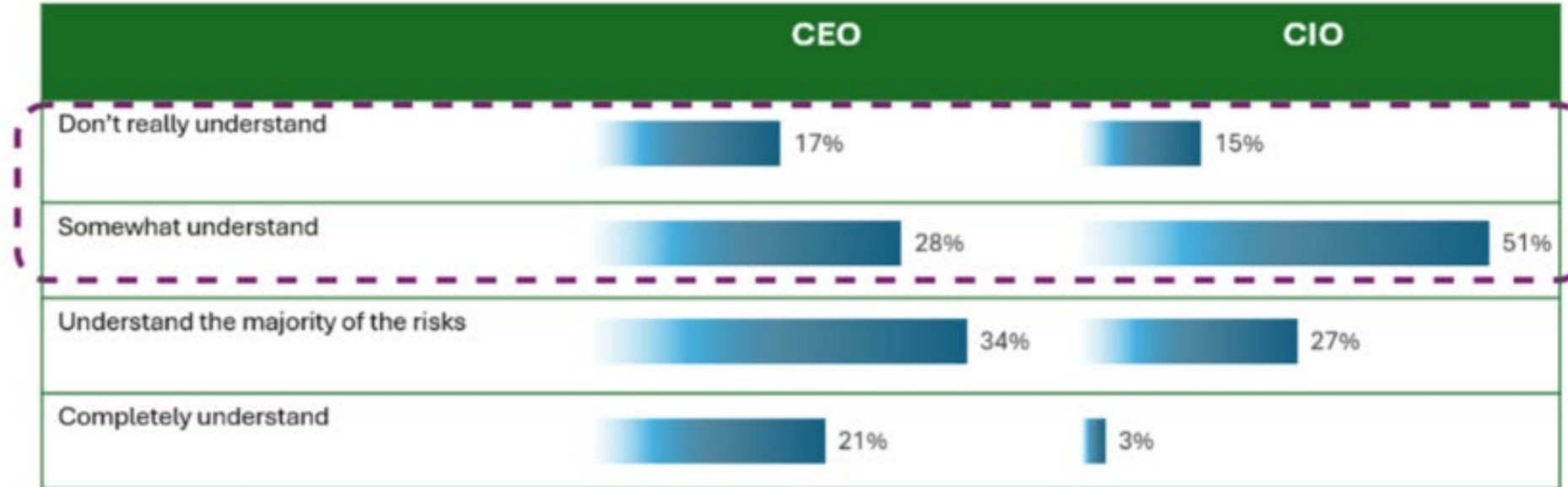
# What risks are we navigating?

**This is not a complete list...**

- AI poisoning
- Bandwidth
- Bias
- Brand threat
- Copyright infringement
- Cost
- Data poisoning
- Data spill
- Environmental impacts
- Evolving regulatory landscape
- Governance
- Implementation complexity
- Integration

- Interoperability
- Lack of common standards
- Limited explainability
- Litigation
- Lying/confabulation
- Performance
- Data quality
- ROI
- Security
- Selection
- Sub-optimization
- Vendor lock-in
- ...

# CEO, CIO opinion on technology vendor's understanding of the potential risk of AI.

**45% of CEOs and 66% of CIOs feel technology vendors do not completely understand the risk potential of AI**

### Opinion on technology vendors' understanding of the potential downside risks of AI

| | CEO | CIO |
|---|---|---|
| Don't really understand | 17% | 15% |
| Somewhat understand | 28% | 51% |
| Understand the majority of the risks | 34% | 27% |
| Completely understand | 21% | 3% |

Source: IDC Worldwide CEO Survey 2024; IDC CIO Quick Poll (January 2024)

# "Implementation" of GenAI 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

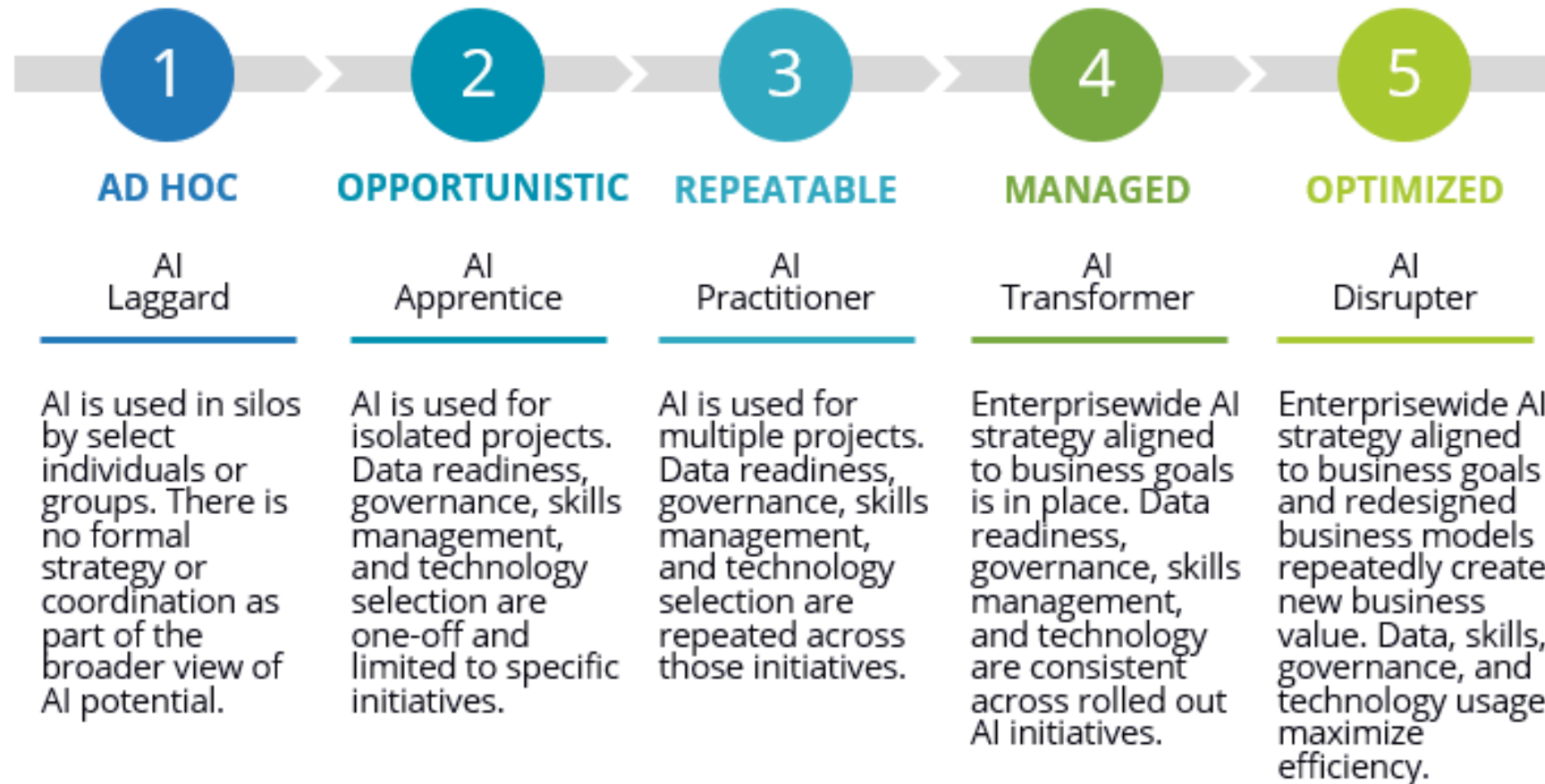**What type of transformation do you expect for your ■ AI software services and ■ AI lifecycle applications?**

| Category | 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% |

# Five components of GenAI implementation
## Checklist for managing risk

**Technology**
- Systems and platforms
- Transparency and interpretability
- Bias mitigation, security, and privacy
- Policies and risks remediations

**Processes**
- Data-centric methods
- Model-centric methods
- Product management
- Platform

**Data**

**Governance**
- Oversight and accountabilities
- Trade-offs
- Data governance
- Model governance
- Project criteria

**Talent**
- Transparency and ethical expertise
- Technical and data skills
- Model skills
- Training and communication

IDC

# Start by assessing your maturity*

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| **AD HOC** | **OPPORTUNISTIC** | **REPEATABLE** | **MANAGED** | **OPTIMIZED** |
| AI Laggard | AI Apprentice | AI Practitioner | AI Transformer | AI Disrupter |
| AI is used in silos by select individuals or groups. There is no formal strategy or coordination as part of the broader view of AI potential. | AI is used for isolated projects. Data readiness, governance, skills management, and technology selection are one-off and limited to specific initiatives. | AI is used for multiple projects. Data readiness, governance, skills management, and technology selection are repeated across those initiatives. | Enterprisewide AI strategy aligned to business goals is in place. Data readiness, governance, skills management, and technology are consistent across rolled out AI initiatives. | Enterprisewide AI strategy aligned to business goals and redesigned business models repeatedly create new business value. Data, skills, governance, and technology usage maximize efficiency. |

*While this maturity model was designed for traditional AI, it is consistent with GenAI.

# Typical use cases
## Three categories of GenAI use cases deliver increasing business value with correlated maturity requirements and risk

**Use case categories range from productivity to functional to industry. They provide increasing differentiation while necessitating increased levels of control over model architecture, security, data privacy, and governance**

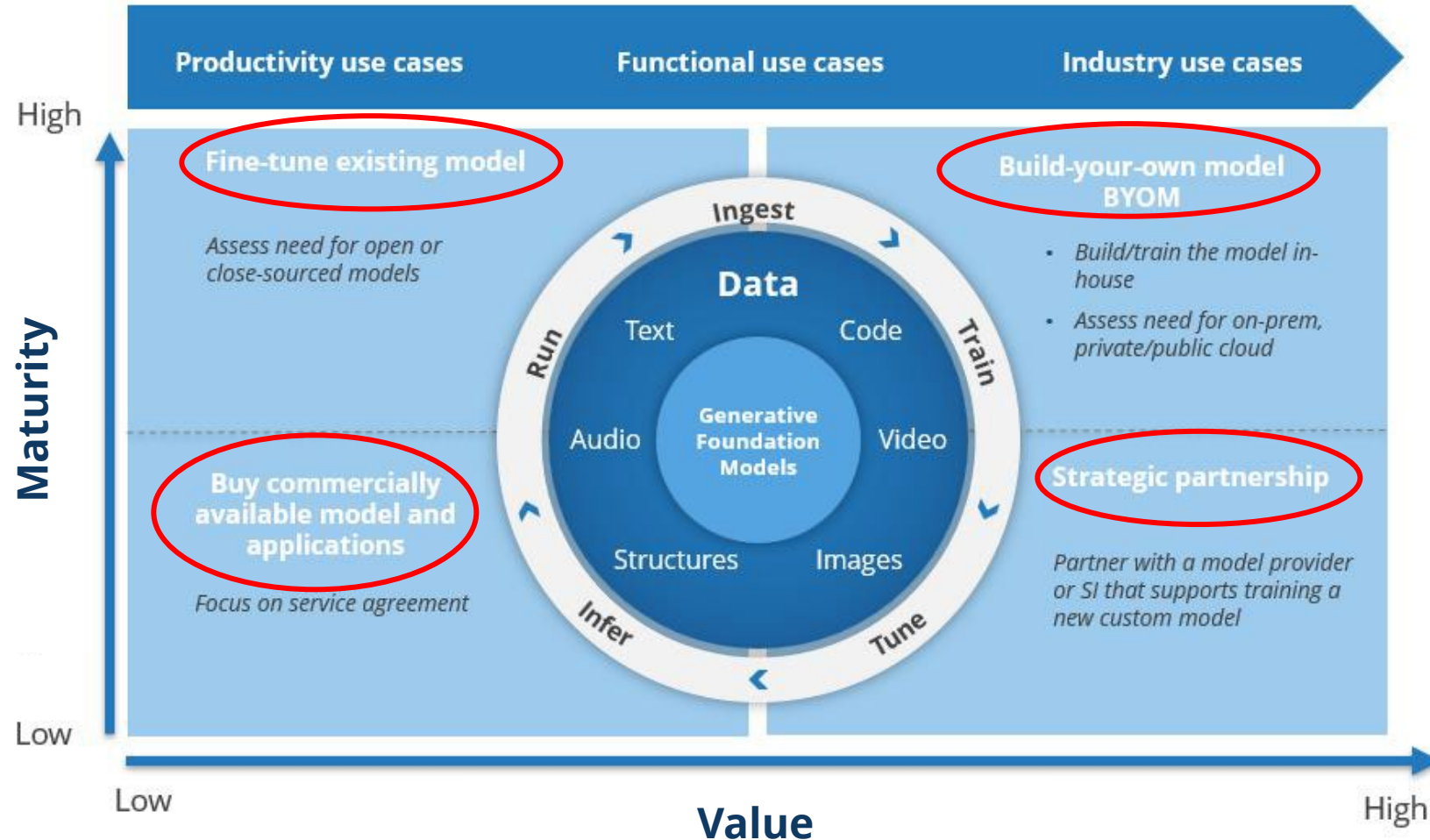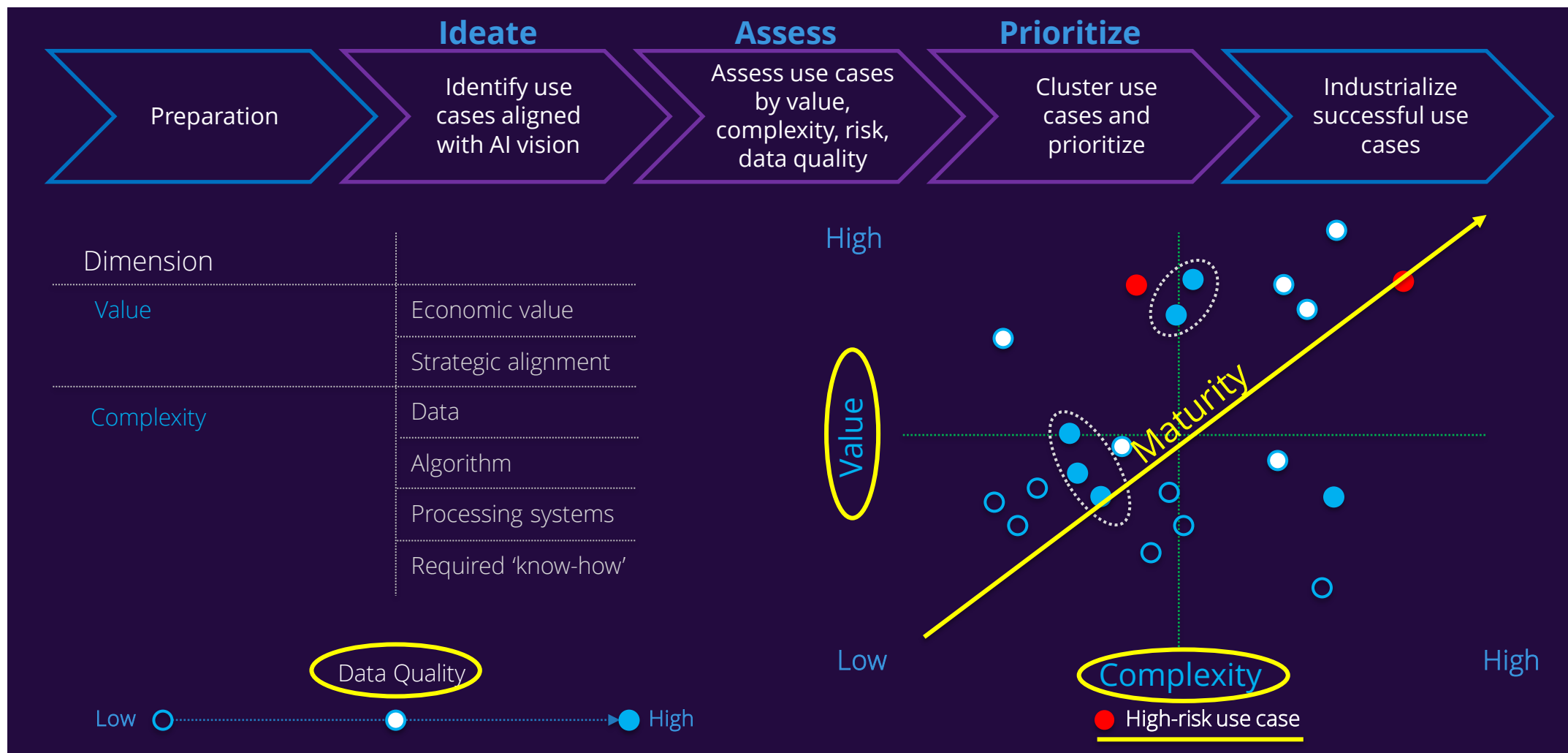| Use Cases Categories | Business Impact | Drivers | Possible Implementation Approach | Use Case Example | |
|---|---|---|---|---|---|
| **Productivity or Efficiency** | • Task **productivity**<br>• Operational **efficiency** | • **Limited talent** in-house<br>• Limited **budget**<br>• **Low risk** appetite<br>• Early adoption<br>• Limited/**poor data** | • Commercial applications with **embedded GenAI**<br>• Native GenAI **standalone applications** (e.g., Microsoft Copilot, Jasper AI, etc.)<br>• **Commercial** models | • Summarizing a report<br>• **RFP** creation<br>• **Code** production<br>• Use case development | **Increasing Maturity, Risk, and Benefit** |
| **Functional** | • Increased functional **effectiveness**<br>• **Contextualized experiences** | • **Good data**<br>• Available **talent in-house**<br>• **Budget** available<br>• **Medium risk** appetite | • Fine-**tuning open-source** models<br>• Fine-tuning models available from model hubs and AI platforms<br>• Retrieval-augmented generation **(RAG)** | • Hyper-personalized sales and marketing<br>• Hyper-**personalized wealth and investments** management<br>• Generative product design and prototyping | |
| **Industry or Transformational** | • **New digital business models**, products, and services<br>• Competitive **moats** | • **Strategic** differentiator<br>• Talent in-house or partner<br>• **Quality** and quantity of **institutional data** | • Fine-**tuning third-party or industry** models<br>• Custom-built models **(BYOM)**<br>• **Strategic partnering** | • Generative **drug discovery**<br>• Generative **material design** | |

# Build versus buy
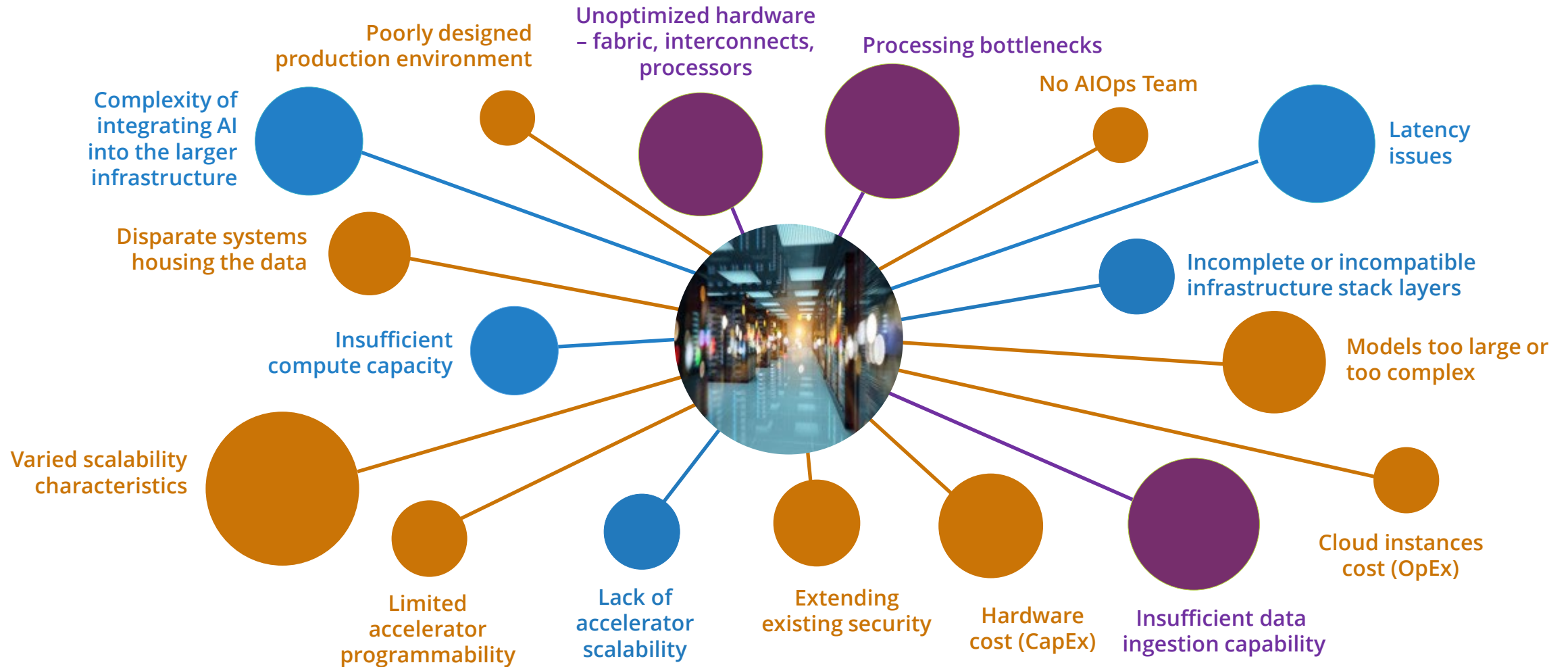## A false dichotomy – really build **and** buy

**Organizational GenAI goals will be influenced by their individual business and technology circumstances. This will result in a mix of build and buy approaches across different use cases**

# Balancing risk, value, complexity, and data quality

# Many AI projects fail because of improper attention to infrastructure



Complexity of integrating AI into the larger infrastructure

Poorly designed production environment

Unoptimized hardware – fabric, interconnects, processors

Processing bottlenecks

No AIOps Team

Latency issues

Disparate systems housing the data

Insufficient compute capacity

Incomplete or incompatible infrastructure stack layers

Models too large or too complex

Varied scalability characteristics

Limited accelerator programmability

Lack of accelerator scalability

Extending existing security

Hardware cost (CapEx)

Insufficient data ingestion capability

Cloud instances cost (OpEx)

# Selecting the right infrastructure stack
## Part 1 – What and how

**Like any blanket recommendation: YMMV. These are guides, not rules**

| | Lean on-premises or collocation (Private Cloud) if | Lean off-premises (Public Cloud) if |
|---|---|---|
| **AI Initiatives** | On-going AI, consistent initiatives | Intermittent AI initiatives |
| **System utilization** | Keeping utilization rates high (keep expensive processors busy) | Limited/inconsistent utilization rates |
| **IT Skills** | In-house skills for complex AI deployments | Limited IT skills for AI deployments |
| **Facilities** | Sufficient 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 |

# Selecting the right infrastructure stack
## Part 2 – Model considerations

| | Lean on-premises or collocation (Private Cloud) if | Lean 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 | Limited customization |
| **Model customization** | Heavily customized | No API changes or customization |
| **Model Performance** | High performance requirements | Lower performance requirements |

# Selecting the right infrastructure stack
## Part 3 – Data considerations

| | **Lean on-premises or Collocation (Private Cloud) if** | **Lean off-premises (Public Cloud) if** |
|---|---|---|
| **Data sensitivity** | Sensitive data, strict data compliance requirements, proprietary data | Data is not proprietary, limited compliance requirements or sanitized |
| **Data isolation** | Model data <u>cannot</u> mix with public data, requires isolation | Model data can safely mix with public data, does not require isolation |

# Selecting the right partner when building an AI infrastructure stack



Source: IDC, 4Q23

*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.*

# Some Takeaways

A realistic understanding of AI maturity is key to implementation decisions

If you don't have or can't build AI maturity, buy it with knowledge transfer

A platform approach with integration to existing systems is key

Know your data with confidence (or at least the data you will use)

"Implementation" risk may arise from existing, licensed products incorporating GenAI

Generative AI is a business transformer that happens to be a technology

Daniel Saroff

Group VP
End User Consulting and Research

dsaroff@idc.com
https://www.linkedin.com/in/daniel-saroff-9301991/

IDC.com

linkedin.com/company/idc

blogs.idc.com

© IDC