

## 7-layer Al Stack framework & Al Market Maps

April 2025

## **Executive summary**

The AI ecosystem is becoming increasingly complex, making it more challenging for organizations to navigate how systems are built, scaled, and governed. To bring structure to this landscape, **AI Insider has developed a 7-layer AI Stack** framework that clarifies the interdependencies between infrastructure, models, workflows, and governance.

Across these seven layers – Hardware & Datacenter Enablers, Data Structure & Processing, Model Development & Deployment, Inference, Orchestration, Tooling, and Security & Governance – we have mapped over 450 solutions to provide a detailed and comprehensive view of the AI market.

The framework is intended as a baseline for evaluating AI strategies, priorities and understand the AI ecosystem from its foundational infrastructure to services and applications, building a structured reference for understanding how the market is evolving. Organizations adopting a full-stack perspective may be better positioned to manage complexity, optimize resources, and align AI initiatives with broader enterprise objectives.

This deck presents a curated selection of the full mapping for illustrative purposes. To access the complete maps or discuss the framework in more detail, please contact us at *hello@resonance.holdings*.



## The 7 Layers of the AI Stack A strategic framework for market understanding

#### 6 Tooling

The top layer empowers builders. It includes developer environments, agent builders, app creation tools, and MLOps – making it easier to prototype, ship, and scale AI-powered applications.

#### 5 Orchestration

Al systems often use multiple models and tools in sequence. This layer decides how components interact – routing prompts, managing agent behavior, and evaluating outputs across workflows.

#### Inference

Once models are live, this layer handles real-time execution. Inference systems, APIs, and microservices ensure models can respond at scale, in milliseconds, to user inputs or system gueries.

#### 3 Model development & deployment

With compute and structured data ready, this is where models are trained, fine-tuned, and deployed. It includes foundation models, deep learning frameworks, and interfaces to production environments.

#### **Data Structure & Processing**

Raw data must be cleaned, normalized, and organized before AI can learn. This layer structures data into usable formats and enriches it with embeddings and vectorization for learning and retrieval.

#### Hardware & AI Datacentre Enablers

The stack begins with infrastructure. This layer provides the computational foundation – from GPUs and AI chips to power, cooling, and storage. Without it, no training or inference is possible.



7

Security & Governance



## **The 7 Layers of the Al Stack** 41 building blocks to map the Al Stack across the 7 layers



AI INSIDER

### Market Map: Core Al Stack April 2025 – A work in progress



## Market Map: Hardware & AI Datacenter enablers Workflow

Core AI stack layers

#### Power

1. The flow begins with Power, which delivers a continuous and stable energy supply to the entire Al infrastructure. This includes redundant power distribution systems, battery backups, and (in cutting-edge cases) small modular nuclear reactors. Power is critical, because even brief interruptions can crash long training runs or interrupt 24/7 inference systems.

#### **AI Chips**

4. Al chips may be generalist (like TPUs) or specialist (e.g., vision-only or edge inference chips) and provide higher efficiency, speed, and energy savings compared to traditional GPUs in targeted Al scenarios. They are especially useful for companies that need scalable, repeatable Al performance without the overhead of general-purpose compute.

#### Storage infrastructure

7. Alongside compute, storage infrastructure keeps datasets, logs, checkpoints, and model outputs safely stored. This may include HDDs, SSDs, or flash arrays – all optimized for throughput and durability to support both training and inference workflows.

#### Racks

2. With power in place, Racks physically house the hardware — such as GPUs, AI chips, and storage systems — in dense, modular configurations optimized for airflow and cable management. These racks form the scalable physical core of compute clusters.

#### Cooling

5. These processors generate a massive amount of heat, so Cooling systems are used to stabilize thermal conditions. Whether air-based, liquid-cooled, or immersed, cooling keeps the cluster operational and prevents performance throttling or failure.

#### **Physical monitoring**

8. Physical Monitoring tools continuously track temperature, power draw, uptime, and failures across the rack and server infrastructure. This ensures reliability and early problem detection, feeding telemetry into datacentre management systems.

#### **Processors (GPUs)**

3. Inside the racks, GPUs and AI chips handle the heavy lifting. GPUs specialize in parallel matrix operations critical for deep learning, unlike general-purpose CPUs. Their massive parallelism accelerates AI tasks: training, fine-tuning, and inference.

#### **Networking infrastructure**

6. While computation happens, Networking Infrastructure connects all nodes across the datacentre – enabling distributed training, data movement from storage, and real-time communication between inference servers. High-bandwidth, low-latency networking (e.g. InfiniBand, 400G Ethernet) is essential to scale.

#### Virtualization (Graphics & DC)

9. Finally, Virtualisation (Graphics & Datacentre) abstracts the underlying compute into shareable, scalable resources. Through virtualization layers, GPU capacity is assigned to users and jobs, allowing teams to run workloads in isolated environments (e.g., vGPU, Docker, or K8s).

## Market Map: Hardware & AI Datacenter enablers

April 2025 – A work in progress

## Core AI stack layers Hardware & AI Datacenter enablers

尔线程

ELECTRONICS

AOORE THREADS

SAMSUNG

#### Non exhaustive





130 +

neider

Electric

**GABYTE**<sup>®</sup>

( Room Alert

SIEMENS

🤼 Sunbird®





### Market Map: Data Structure & Processing Workflow

Core Al stack layers

Orchestration Inference odel development & deployme Data Structure & Processing

ecurity &

In many environments, this logical storage layer is also connected to Storage Infrastructure (Layer 1) to enable caching or fast local access during runtime (e.g., training or inference). However, the core focus here is to provide a platform-agnostic, structured foundation for all downstream Al tasks – including pre-processing, embedding, retrieval, and labeling.

#### **Cloud & Logical Data Storage**

1. After the infrastructure is provisioned and compute is enabled at the hardware level (Layer 1), AI workflows begin by organizing and storing raw data inside Logical Data Storage systems. At this stage, files such as text corpora, image datasets, or structured exports are ingested into cloud object stores, structured data lakes, or dataset catalogs. These systems are abstracted from physical hardware and are accessed via APIs and pipelines — enabling reproducible, versioned, and collaborative data workflows.

#### **Data integration**

2. Because datasets often originate from fragment systems or formats, the next step is Data Internet Pipelines ingest content from databases warehouses, or SaaS platforms of consistent schemas and the to ensure that all the coherent set



# Market Map: Data Structure & Processing April 2025 – A work in progress



Non exhaustive



# Market Map: Model dev. & deployment

Tooling Orchestration Inference Model development & deployment Data Structure & Processing Hardware & AI Datacenter enablers

#### **Compute & accessibility providers**

1. With data structured and embeddings generated, the system calls on Compute & Accessibility Providers to provision infrastructure — GPUs, TPUs, or Al-optimized chips, either via the cloud or on-prem clusters. This is the raw horsepower required to run large-scale training and inference jobs.

#### **Foundation models**

4. With data and frameworks ready, teams may choose to start from scratch — or use a Foundation Model. These are large, pre-trained models like GPT Claude, or LLaMA that already capture generation knowledge. They're retrieved from models and serve as a strong base for

#### ML deep learning framework

2. Once compute is available, developers use ML deep learning frameworks to define model architectures, loss functions, optimization routines, and training workflows. These frameworks are where the training logic is programmed and run.

#### Data labelling

Simultaneously

e Data L

## Market Map: Model dev. & deployment

### April 2025 – A work in progress

Model development & deployment
Hardware & AI Datacenter enablers

#### Non exhaustive





### **Market Map: Inference** Workflow

**AI** INSIDER



Inference

Core AI stack layers

# Market Map: Inference April 2025 – A work in progress

Non exhaustive





Al inference Anyscale aws 40+ **Solutions** mapped For more details, please email hello@resonance.holdings



### Market Map: Orchestration Workflow

Labom Core AI stack

Orchestration

#### **Prompt & Model Orchestration**

1. Orchestration begins with Prompt & Model Orchestration, where developers define the logic that sequences prompts, tools, and model outputs. This is the "brain wiring" of an LLM system — it decides how to query models, chain tasks, and retrieve external data.

#### Model gateway & routing

 As multiple models or tools are used, Mode & Routing systems dynamically decide which endpoint to use. They enable A/B tech load-based routing — choose Mistral for speed, for the system.

#### For more details, please email hello@resonance.holdings

#### Agent frameworks

2. Once logic is defined, Agent Framework autonomy to the system. These free models into stateful agentee maintain memory of respond to us orche

## Market Map: Orchestration April 2025 – A work in progress

#### Non exhaustive

Al21 Jobs 😤 Langfuse APromotLaver 🗿 Haystack	
S stonebranch HoneyHive PREFECT S Flyte	
	50+
	50+ Solutions
	50+ Solutions mapped



Core AI stack layers

# Market Map: Tooling

roomig
Model development & deployment
Hardware & AI Datacenter enablers

#### **Dev environments**

Cloud-native or local environments that enable developers to build, test, and debug AI code with preconfigured dependencies and GPU access. These platforms streamline reproducibility, CI/CD workflows, and agent/model prototyping.

#### End-to-end AI app builder

Platforms to build, integrate, and deploy full Alpowered applications — combining logic, models, APIs, UI, and workflows. Targets builders looking ship LLM apps, dashboards, and copilor managing infrastructure.

#### AI Ops (LLM & ML)

Operational tooling that monitors, evaluates, and improves AI systems in production. Includes LLMOps for prompt tracking, hallucination detection, and token usage, and MLOps for retraining, model versioning, and pipeline observability.

#### Compute & cost on t

Core AI stack layers

Platforms that impresent infrastructures acrossed

## Market Map: Tooling April 2025 – A work in progress

#### Non exhaustive



Tooling

Core AI stack

layers

### Market Map: Security & governance Workflow

Core AI stack layers

#### Al security & risk mitigation

This category covers tools and practices that protect AI systems from emerging threats, vulnerabilities, and misuse hacking, prompt injections, model theft, hallucinations, or unsafe outputs. The main focus is to ensure the reliability robustness, and ethical integrity of models throughout their the training to inference

#### **Data Governance, privacy & complia**

Data Governance, prive platforms mane control control

# Market Map: Security & governance April 2025 – A work in progress

Core AI stack layers

Security & governance

Non exhaustive





# Thank you

