

**STIFEL | IRIS**

INTELLIGENCE • RESEARCH • INSIGHTS • SERVICE

# DATA CENTRES BEYOND HYPERSCALERS

**Thriving in the shadow of the cloud and AI giants**

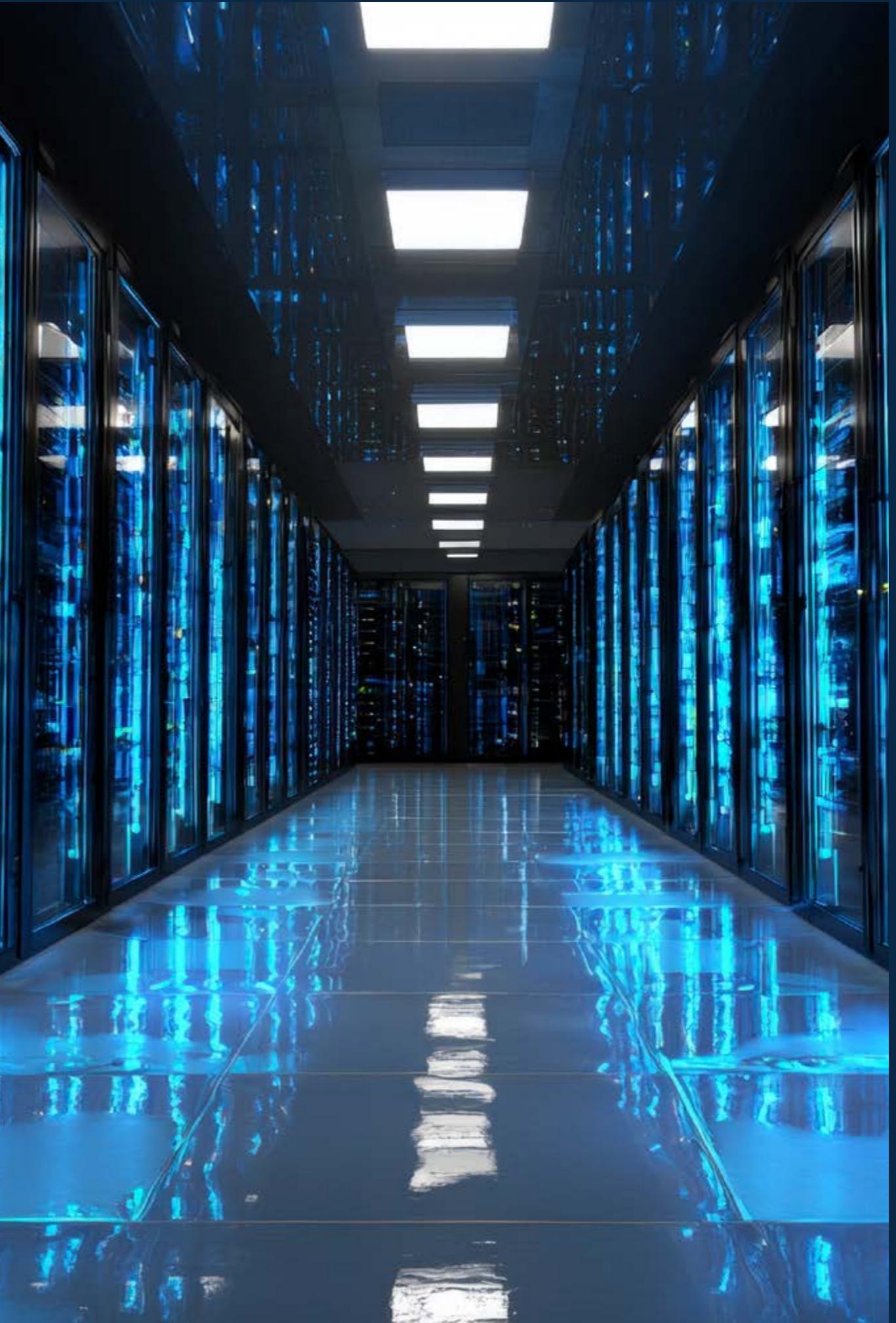
JANUARY 2026



# TABLE OF CONTENT

---

<b>1. The age of giants in AI infrastructure</b>	<b>7</b>
The great compute buildout	8
How tech behemoths dominate an era of AI maximalism	17
Colliding with the power wall	24
<b>2. Carving out value in the data centre ecosystem</b>	<b>31</b>
Unbundling the digital infrastructure stack	32
The AI-native cloud arises	40
The playbook is not to compete, but to bypass	47
<b>3. A European deep dive</b>	<b>55</b>
Understanding the European data centre landscape	56
Identifying European champions	63





## EXECUTIVE SUMMARY

---

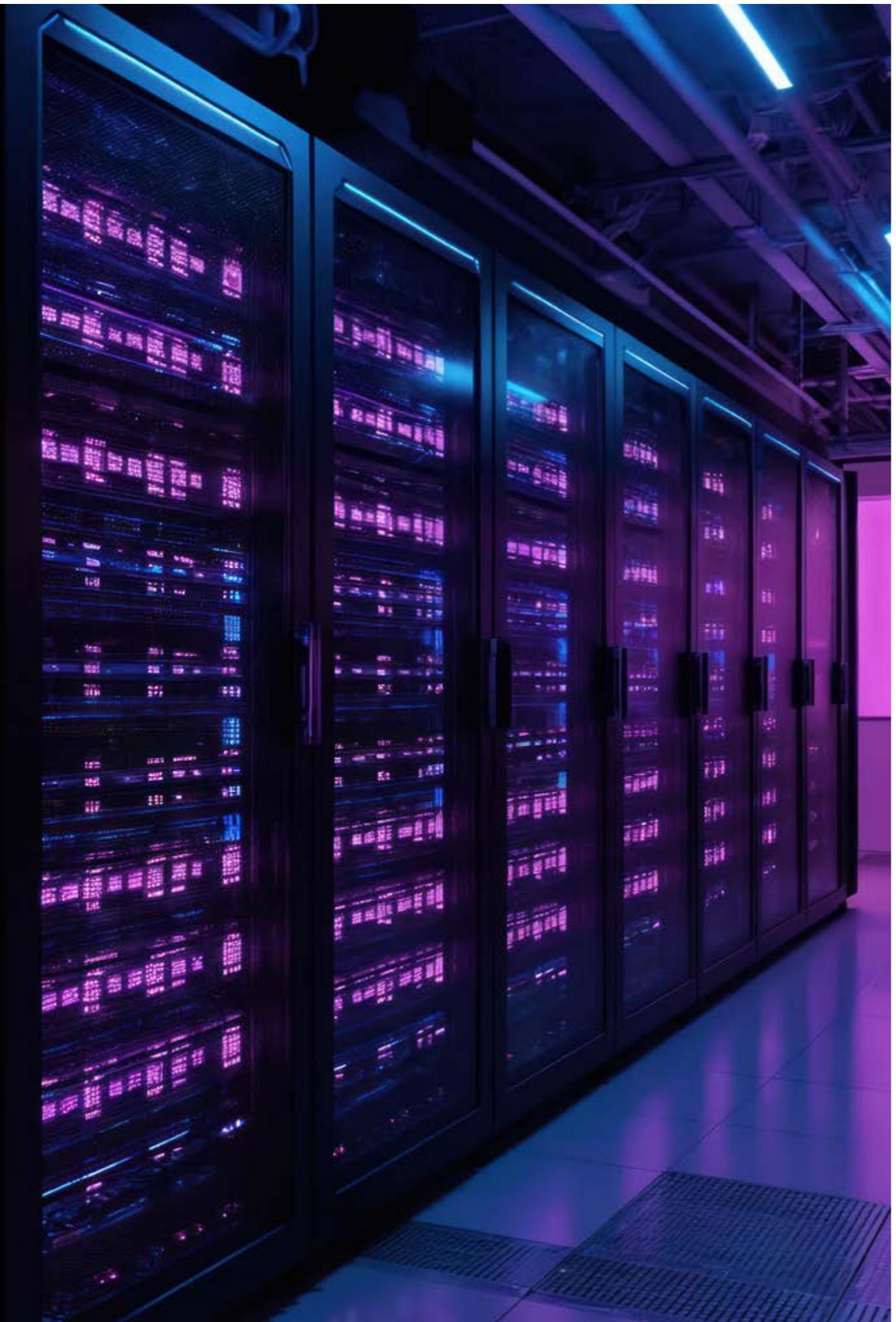
The generative AI boom has triggered a capital supercycle in data centres, defined by a brute-force compute scaling race among tech incumbents. Increasingly, the competitive parameters are such that only a few tech giants look set to dominate the core AI layers.

This narrative, however, overlooks significant market disruption at the periphery, with the emergence of a new and diverse data infrastructure ecosystem. The sheer scale of the AI market means that even niche competitors can capture valuable segments. But as the market attempts to price this trillion-dollar outlay and the bubble debate intensifies, in this paper we aim to shift investor focus from the core to the periphery and ask: while the AI giants build the biggest engines of intelligence, where can agility and specialisation unearth hidden champions?

---

This paper offers a framework for analysing the emergence of specialist firms within the data centre value chain, examining how AI forges new competitive entry points. With a focus on the European landscape, it dissects the strategies firms can use to bypass data centre and AI technology giants.

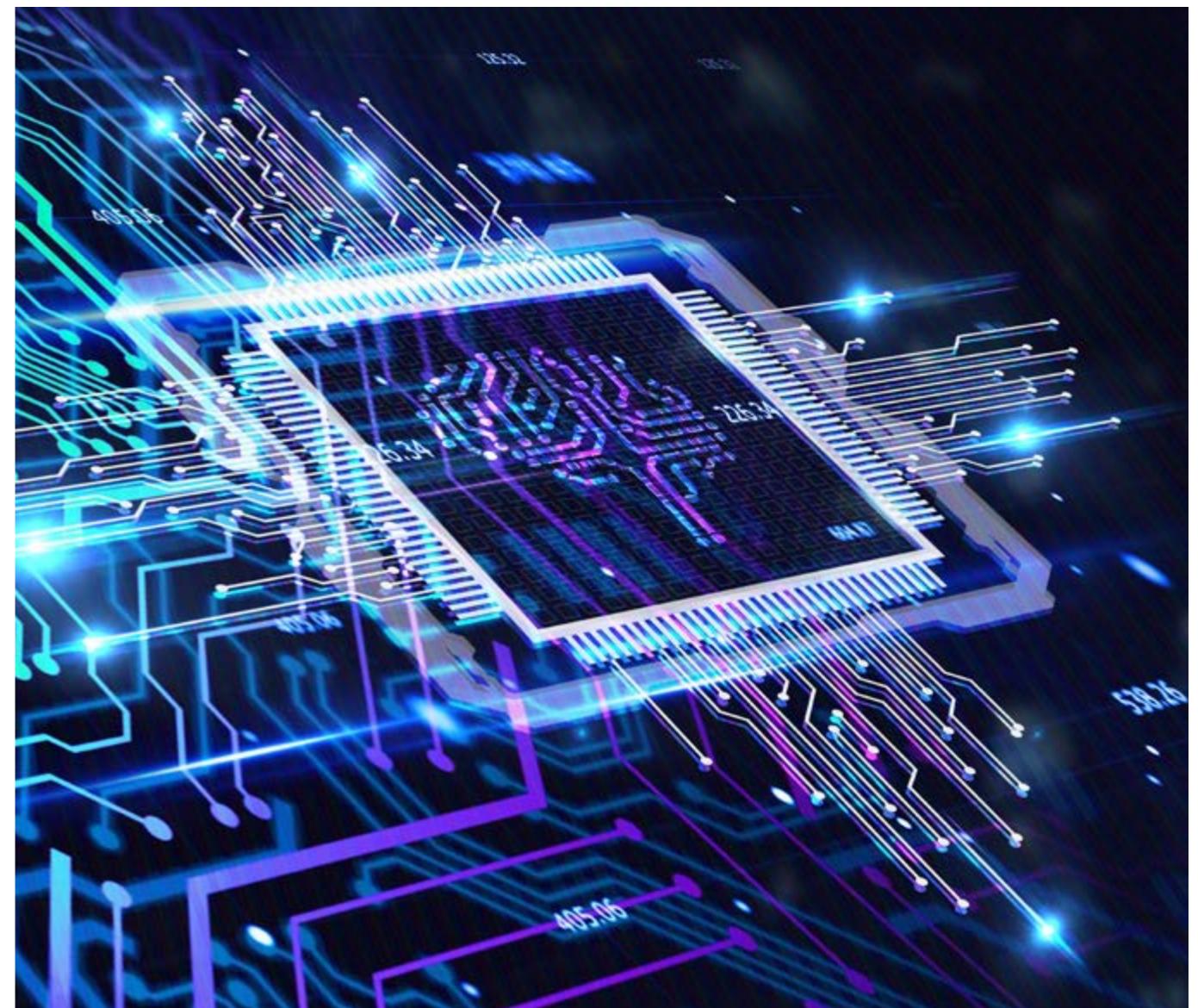
---



# THE AGE OF GIANTS IN AI INFRASTRUCTURE

## SECTION 1

The AI infrastructure race is being forged by its giants. Unprecedented data centre capex, set to escalate further in 2026, is led by hyperscalers whose capital intensity is rising sharply. This surge is concentrating the entire AI stack, from silicon to foundation models, into a few interconnected oligopolies. This dynamic solidifies their dominance and fundamentally questions the scope for competition in the data centre ecosystem.



# 1.1 THE GREAT COMPUTE BUILDOUT

The AI buildout is a story of structural tensions. The digital race is colliding with physical realities like power access. Financially, sizeable capex grapples with an unclear ROI, torn between nascent AI revenues and the need to protect tech giants' core franchises. Meanwhile, the narrative is pivoting from an obsession with brute-force AI training scaling to a increasing focus on inference efficiency.

## DATA CENTRE INFRA MOMENTUM ACCELERATES INTO THE BOTTLENECKS

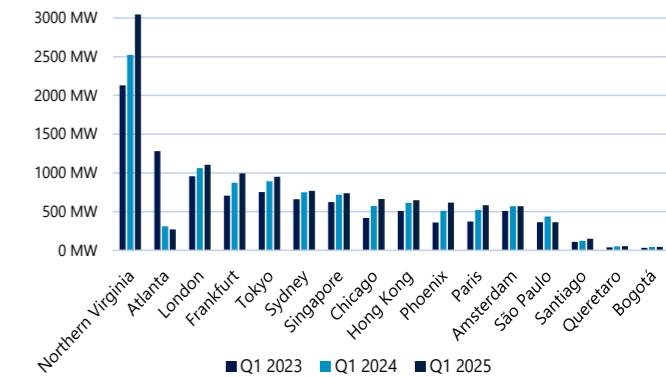
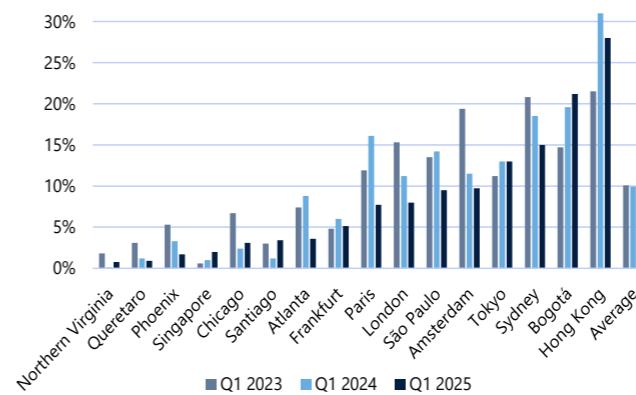
The data centre capex cycle continues to exhibit robust, AI-driven momentum despite mounting complexities. Drawing from our tracker, which aggregates capex consensus (Bloomberg) across approximately fifty data-centre-focused companies, we estimate that data centre capex will grow by 68% in CY25, following a 75% expansion in CY24, marking a 2.9x increase from 2023 levels. This surge is primarily fuelled by demand for hyperscale computing and AI workloads. Yet, the cycle is expected to start maturing from 2026 onward, characterised by decelerating capex growth rates as the industry faces entrenched supply constraints.

Recent results from the major AI firms confirmed a higher-for-longer capex trajectory, with further upward revisions again surprising the market. While this does not defy the mathematical certainty of a deceleration, as the ~70% growth seen in CY24-25 is most likely unsustainable, it significantly lifts the entire forecast curve. Our capex tracker projects this will result in another year of strong, supra-normal growth in CY26 at 30%. Only by CY27 is the cycle expected to enter a true maturation phase, slowing to 12% as it faces deep-rooted supply constraints. Most

indicators signal that these bottlenecks, particularly in power and infrastructure, will govern the market's ceiling, tightening capacity expansion until at least 2029.

This dynamic shifts investor focus. The current, high-demand environment is increasingly seen as a given for the next few years, a reality confirmed by hard colocation data: globally, absorption has consistently outpaced new deliveries, compressing vacancy rates to historic lows. This is most acute in North America, where an availability crisis has pushed vacancy in key markets like Northern Virginia to sub-2% levels. While forecasting data centre capex remains inherently challenging, as stakeholders continually recalibrate commitments amid shifting demand signals, this volatility now applies more to the scale of the peak than the direction of the short-term trend. The combination of significant hyperscaler commitments and a structural supply crisis has largely locked in this near-term trajectory. The focus is thus shifting to the post-2030 horizon, centred on the principal debate: is this a sustainable, permanently high plateau, or is the market inflating a bubble that will result in a painful 'hangover'?

FIGURES 1&2: (LEFT) DATA CENTRE VACANCY RATE BY MARKET (RIGHT) DATA CENTRE INVENTORY BY MARKET

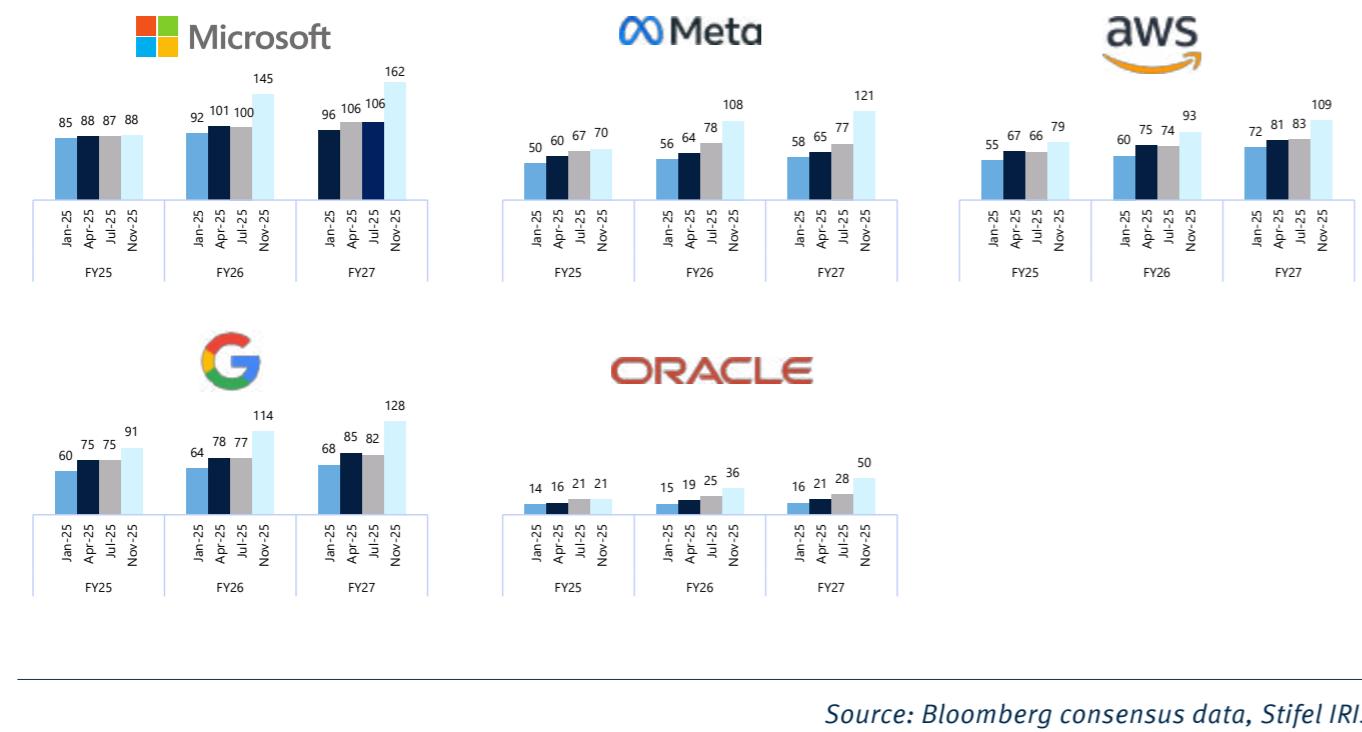


Source: CBRE, Stifel IRIS

The primary AI buildout, focused on significant-scale training data centres for foundational models, is consolidating around the hyperscalers and a handful of specialised neocloud providers. We estimate Microsoft, AWS, Google, Meta and Oracle are driving about 60% of global data centre investment, though the scope for data centre capex depends heavily on how data centres are defined. Evidence from recent months indicates spending projections have been further revised upward, driven by AI infrastructure buildouts, with no immediate peak in sight. Hyperscalers' aggregate capex surged 73% year-over-year in CY24 to \$229bn and is projected to increase by another 70% in CY25 to \$389bn, exceeding \$500bn in CY27, according to Bloomberg consensus data. Signs of flow to smaller competitors are emerging, with neocloud providers such as CoreWeave on track to exceed \$50bn in aggregate capex over CY25-27. Yet, their share remains marginal and barely erodes the hyperscalers' dominance and often relies on partnerships with the giants themselves.

The Q3 2025 earnings season confirmed this accelerating investment thesis, with nearly every hyperscaler signalling upward capex revisions. The results were underpinned by robust cloud performance, validating the demand driving the spend: Microsoft's Azure grew 40% year-on-year, AWS re-accelerated to over 20% growth on a \$132bn run-rate, and Google Cloud surged 34%. In response, the hyperscalers explicitly raised spending forecasts. Microsoft's capex reached \$34.9bn for the quarter, with its CFO admitting they are not catching up to demand, implying sustained or higher future spending. Oracle, whose Cloud Infrastructure (OCI) revenue grew 54%, raised its full fiscal year capex forecast to ~\$35bn, noting that demand "continues to dramatically outstrip supply". Similarly, Alphabet lifted its full-year 2025 capex guide to ~\$92bn and warned of a "significant increase" in 2026, while Meta projected its capex dollar growth will be "notably larger in 2026 than 2025". The commentary was unanimous: demand for AI compute is outstripping supply, forcing upward revisions to historic capex plans and solidifying a multi-year investment cycle.

FIGURE 3: HYPERSCALER CAPEX FORECASTS: A STORY OF CONTINUOUS UPWARD REVISIONS IN 2025 (US\$BN)



Since ChatGPT's November 2022 launch ignited the generative AI boom, the financial narrative has gone through some distinct phases, oscillating between transformative optimism and concern. The AI rally began not at a market peak, but from the ashes of the 2022 tech downturn, with the launch of ChatGPT as the primary catalyst. Early 2023 saw explosive adoption, with ChatGPT crossing the 100m weekly user thresholds in about a year and many organisations integrating generative AI into operations. The market narrative was straightforward: AI was a paradigm shift, and the picks & shovels providers would be the first to benefit. This thesis was very much validated during 2023 when Nvidia's earnings exceeded expectations, proving that demand for AI infrastructure was not just hype but a tidal wave of corporate capex.

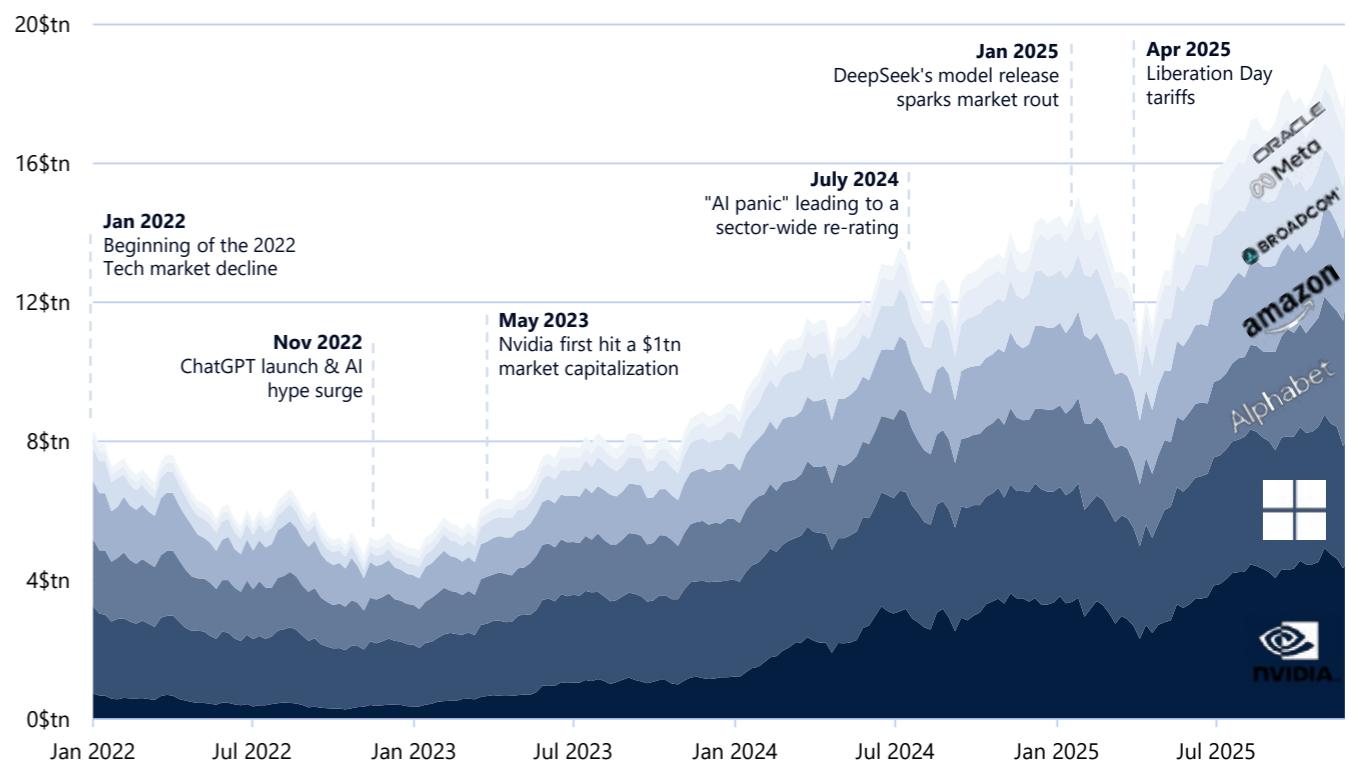
By mid-2024, the story became more complex as relative euphoria gave way to pragmatic concerns over ROI and physical constraints. This tension culminated

in early 2025. Just weeks after the announcement of Stargate (a \$500bn Microsoft/OpenAI supercomputer plan that crystallised the market's 'scale-at-all-costs' optimism) the surprise release of a powerful model from Chinese startup DeepSeek triggered a market sell-off. The event introduced the risk of disruption from lower-cost, highly efficient competitors. However, this panic quickly subsided as the market analysed the underlying trade-offs. While DeepSeek's models gained traction, traffic on its own low-cost service eroded after the initial spike, as the company served its public models with high latency and limited features to conserve scarce compute resources. This injected a new realism: a model's low price does not guarantee market share if the user experience is poor, and DeepSeek's service suffered precisely because it lacked the significant compute infrastructure to serve its popular model effectively at scale.

Today, the financial narrative is more nuanced than ever. The market continues to grapple with the stark contrast between unprecedented AI capex and the still-nascent revenue directly attributable to it. Yet, the rally remains remarkably resilient. The aggregated market capitalisation of the seven largest US AI incumbents has surged from a low of approximately \$5.8tn in early 2023 to over \$18.1tn as of November 24, 2025.

This valuation stands just ~5% shy of the collective all-time high established in early November, defying recent waves of market anxiety surrounding elevated AI valuations. This staggering accumulation of value demonstrates that while the path forward is complex and fraught with operational and competitive risks, the market's conviction in the long-term, transformative power of the AI supercycle is, for now, largely intact.

FIGURE 4: THE \$13 TRILLION AI RALLY: MARKET CAP GROWTH OF LEADING AI TECH COMPANIES (2022-25)



*Source: Refinitiv, Stifel IRIS*



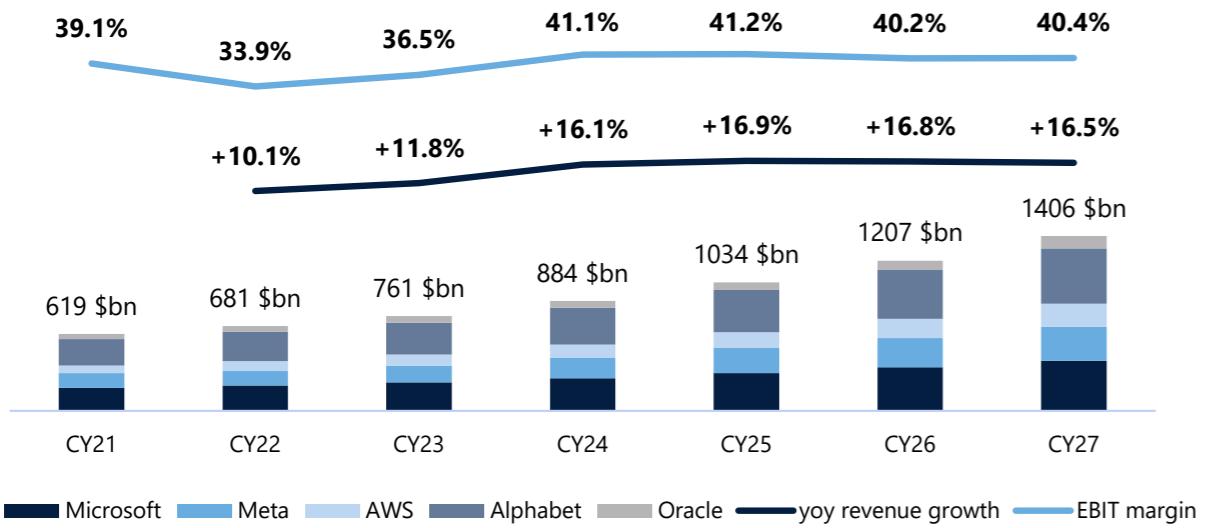
## NEW FRONTIERS OR SHIELDED FRANCHISES?

The immense AI capital outlay from Big Tech sparks debate over its core motivation: is it a defensive fortification of existing franchises or an offensive conquest of emerging AI markets? The financial scale underscores the stakes. The five largest hyperscale investors (Microsoft, Meta, AWS, Alphabet, Oracle) generated \$884bn in aggregated CY24 revenue at a 41.1% combined EBIT margin, excluding Amazon's non-AWS divisions. According to Bloomberg consensus data, this revenue base is projected to expand by nearly \$500bn to \$1,406bn by CY27, implying an accelerating CAGR of 16.7% over CY24-27, from 12.5% over CY21-24. Furthermore, this growth is expected to occur without margin compression as EBIT margins are forecast to average 40.7% over CY24-27, versus 36.5% during CY21-23.

We lean toward interpreting these outlays as a form of compelled self-reinvention. Incumbents are front-loading a race to dominance not merely to capture budding AI markets, but fundamentally to shield their

core, high-margin revenue streams, like Google's Search or Meta's social networks, from AI-driven disruption. The seemingly irrational ROI of this capex surge becomes justifiable only when framed as a need to protect these existing trillion-dollar franchises. This dynamic presents a chicken and the egg dilemma for hyperscalers: underinvestment risks ceding compute supremacy to rivals, yet collective overbuilding risks precipitating a cycle of brutal capital inefficiency and rapid asset depreciation. This is the key difference from the 90s telecom bust: the risk then was building "dumb pipes" that others failed to fill. Today's vertically integrated hyperscalers are building infrastructure for their own services. The risk is not bankruptcy via empty capacity, but a structural erosion of profitability, in our view. This defensive posture is enabled by their deep capital pools and low cost of capital, allowing them to sustain significant investments to maintain their market dominance, irrespective of the immediate profitability of new applications.

FIGURE 5: US HYPERSCALERS: AGGREGATE REVENUE & EBIT MARGIN (CY21-CY27E)



Source: Bloomberg consensus data, Stifel IRIS

The generative AI investment supercycle raises a critical financial question: how will this unprecedented capex be monetised to justify the industry's trillion-dollar outlay? While scepticism over immediate ROI remains, we note that a multi-layered ecosystem of revenue opportunities is already forming.

The most visible pathway is the mass-market web platform. The August 2025 launch of GPT-5 illustrated this pivot: while perceived by many power users as a performance disappointment, its core innovation was a strategic shift toward monetising its significant free user base. The key mechanism introduced was a "router", a system designed to classify user intent. While OpenAI's initial attempt to force this automatic routing on paid subscribers failed after user backlash, the underlying technology remains central to its strategy. This system is unique to OpenAI and appears designed not just for cost optimisation, but to identify commercial intent. This technology lays the groundwork for an agentic superapp model. The strategy is to shift from subscriptions or disruptive ads to a transactional take-rate. In this model, the platform could leverage its vast partnerships (which include Shopify, Booking.com, Stripe, and Instacart) to facilitate end-to-end transactions, taking a small affiliate commission on purchases made via the agent. This move aims to convert user engagement directly into revenue without relying on traditional ad models.

Beyond these mass-market platforms, we believe the most durable, high-margin enterprise value will be created in specialised, functional domains. These niches mark a crucial shift from general-purpose models to inference-driven applications that solve specific, high-value problems. We see three areas as particularly noteworthy:

- **Coding & software development:** Already one of the largest generative AI markets after chatbots, AI-assisted coding has rapidly matured from

demos to integrated developer tools. Platforms like GitHub Copilot and Cursor are reshaping software workflows by moving developers beyond rote tasks (automating code generation, debugging, and optimisation) and embedding a collaborative AI dynamic into daily practice.

- **Agentic AI:** This segment represents a fundamental departure from passive, predictive models. It focuses on proactive agents that chain actions to solve multi-step problems autonomously. By emphasising inference-time compute and reinforcement learning, these systems adapt in real-time to dynamic environments. This is where AI moves from being a simple tool to an autonomous collaborator, finding product-market fit in complex orchestration and decision-making.
- **Physical AI:** This frontier represents the application of modern AI to machines that interact with the physical world, creating a fundamental break from traditional "scripted" robotics. Historically, automation was a rigid, capital-intensive engineering problem, that required perfect, unchanging environments and multi-year integration projects. This model was inflexible and too expensive for most of the economy. AI should fundamentally redefine this paradigm. Instead of being manually programmed, new systems use end-to-end AI, allowing them to learn from data and adapt to new, complex situations rather than just following rigid, pre-programmed rules. This leap is driven by modern foundation models, specifically Vision-Language Models (VLMs), which give the machine advanced perception and "common sense" reasoning, and Vision-Language-Action (VLAs) models, which translate this understanding into generalised physical manipulation. This opens vast new markets that were previously impossible to automate.

FIGURE 6: A FEW EMERGING MARKETS IN AI



## THE LOOMING TRANSITION FROM TRAINING TO INFERENCE

The AI workload landscape is defined by a debate over the transition from training to inference.

- AI Training (the computationally intensive process of creating a foundational model) involves a significant, one-time, power-intensive compute cluster to create the foundational model. Like building a factory, the economic assumption was that this immense upfront investment would be amortised over billions of user queries.
- AI Inference, in contrast, was understood as the process of using that trained model to generate outputs. This was seen as the subsequent opex-driven phase, where each query represents a small, latency-sensitive workload, presumably handled by more distributed, less-dense infrastructure.

This perspective, separating the high fixed cost of creation from a low variable cost of deployment, has been fundamentally challenged as a result of two successive strategic shifts.

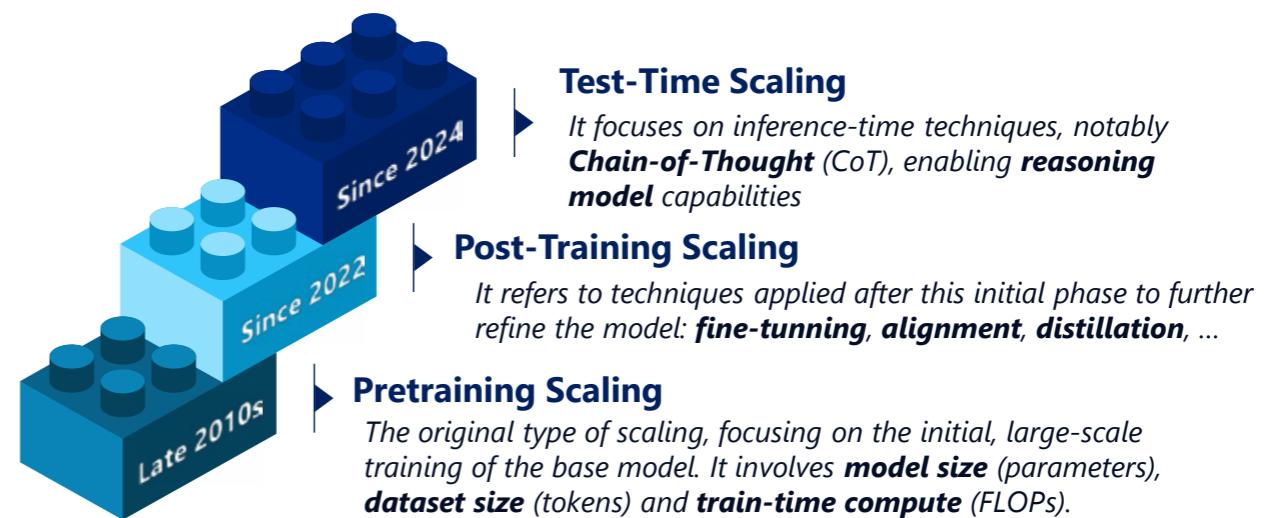
The first shift occurred when brute-force pre-training scaling laws (the belief that LLM performance scales predictably with compute, data, and parameters) hit a wall of diminishing returns. The industry pivoted, discovering that superior reasoning could be achieved not just by building a bigger model, but by using more

compute at the moment of inference via test-time scaling (such as chain-of-thought). This development was the first economic shock to the model: it revealed the variable cost of a high-quality query was not minimal, but in fact, needed to increase to unlock better performance.

The strategic reality of 2025 is now defined by a second transition: the move from the pre-training era to the reinforcement learning era. This new paradigm has irrevocably blurred the line between training and inference. Reinforcement learning is an iterative training process used to refine a model's reasoning capabilities. Crucially, it is very inference-heavy. The process requires the model to generate millions of potential answers (rollouts) which are then scored by a reward model (often another powerful LLM). This entire R&D pipeline (generating rollouts, running LLM judges, generating synthetic data) is fundamentally an inference workload.

Consequently, labs must now build and finance two colossal cluster configurations: a traditional, dense cluster for the initial pre-training, and a second, sprawling, inference-optimised fleet that runs continuously just to power the RL-based training and alignment process. The true barrier to entry is no longer the one-time pre-training cost; it is the capacity to sustain this permanent, inference-heavy R&D cycle.

FIGURE 7: THREE STAGES OF SCALING LAWS



Source: Stifel IRIS

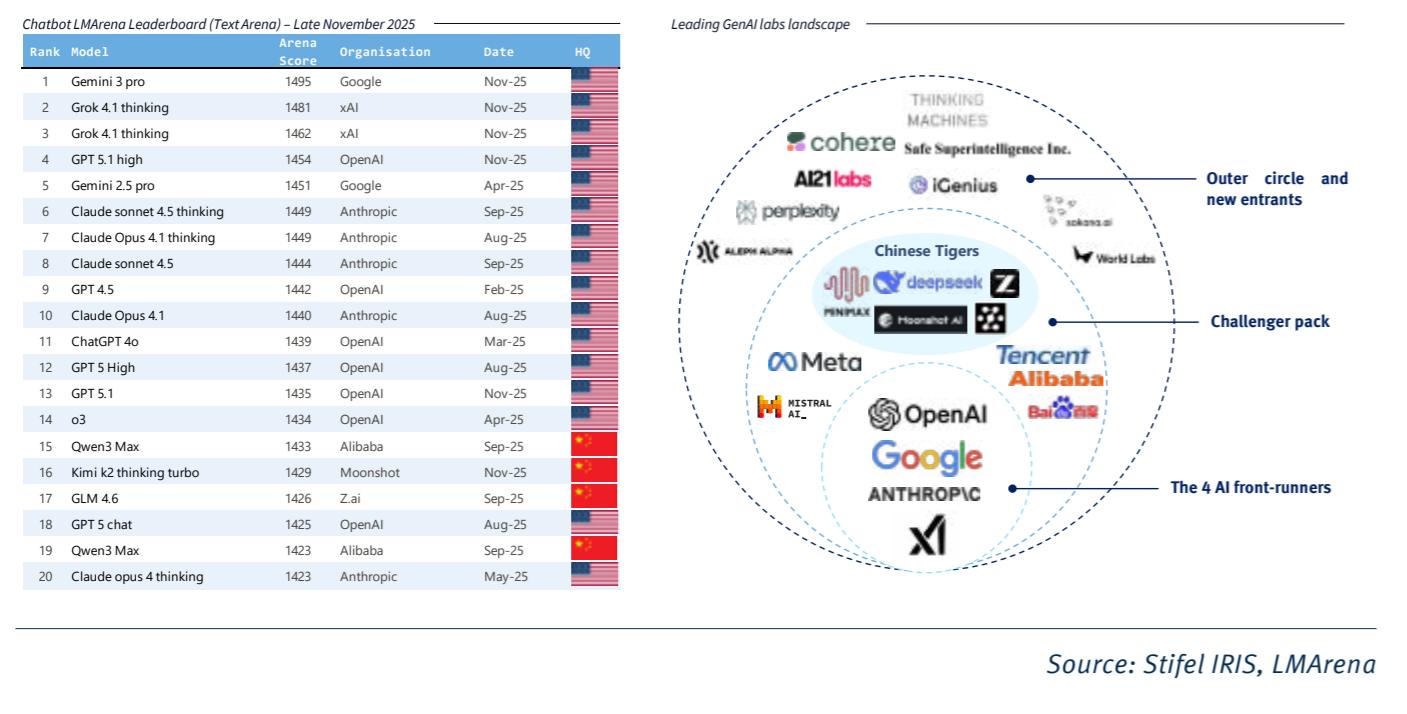
The AI landscape recently experienced a period of perceived commoditisation. Model performance on key benchmarks converged rapidly, and the leaderboard became highly volatile, with state-of-the-art status often appearing to be a function of which lab had the most recent release.

However, we believe this view is evolving, and that a more durable hierarchy is emerging. The top of the leaderboard is increasingly dominated by a quartet of major US AI labs: OpenAI, Google, Anthropic, and xAI. These labs possess, or have access to, the unique capital required to fund both the record-breaking pre-training clusters and the sprawling, inference-optimised fleets needed to scale Reinforcement Learning. This new reality confirms that compute supremacy (for both pre-training and the significant, ongoing RL-inference cycle) is the deciding factor. Meta represents a distinct case. Despite possessing compute power on par with this top tier, its models have not consistently led on closed-model performance benchmarks.

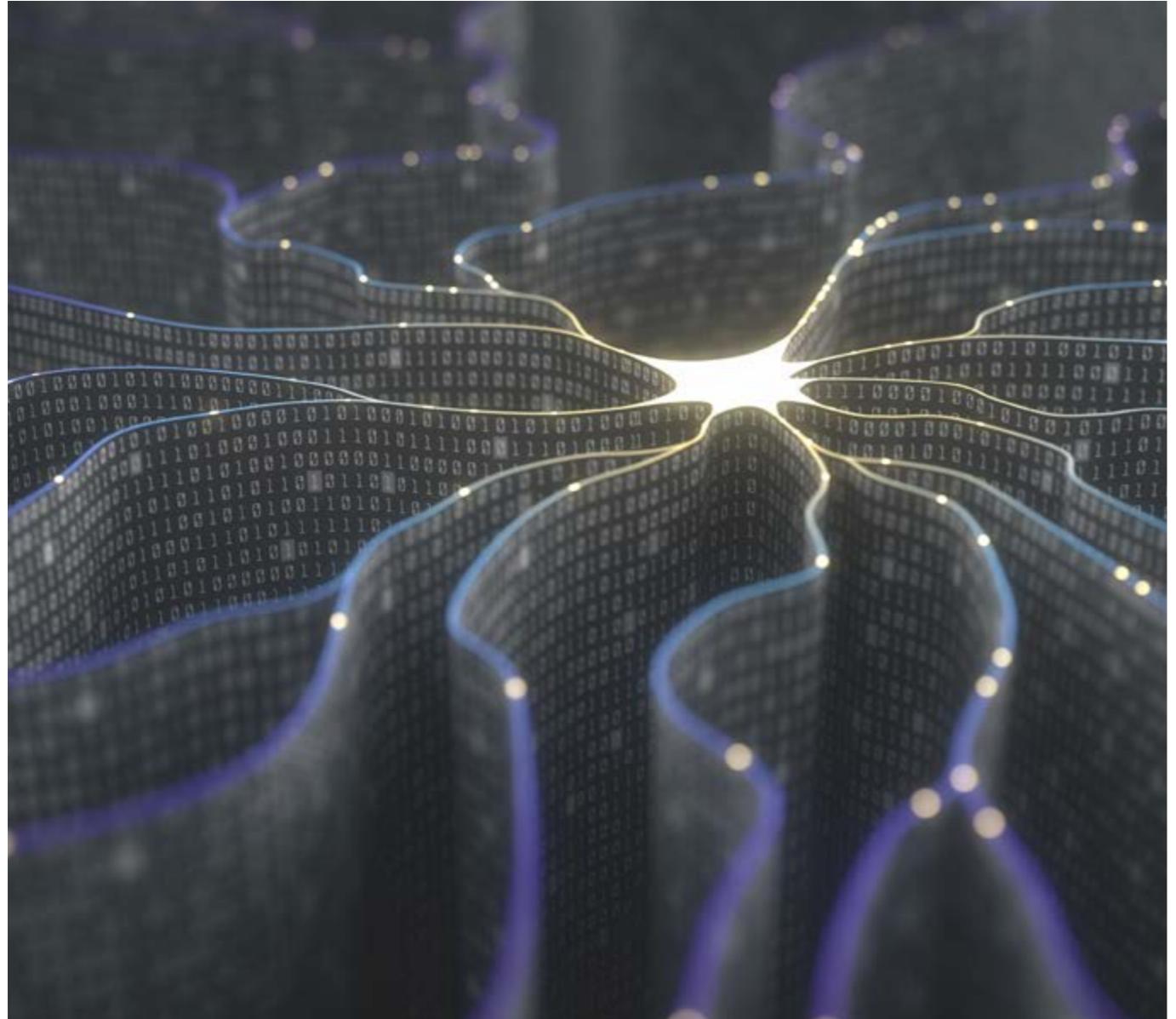
A second tier of challengers operates with less compute, forcing them to pursue alternative strategies. These include Chinese technology giants (Alibaba, Baidu and Tencent) and a group of Chinese smaller contenders often known as the AI tigers. Lacking the compute scale of the major US AI labs, their path relies on differentiation through efficiency or by targeting specific regional or open-source niches.

The shift from training-centric to inference-driven workloads is becoming increasingly visible as usage of LLMs becomes more common. This transition is creating new, acute inference bottlenecks, proving that infrastructure constraints are not limited to pre-training. This is evident even in top-tier labs. Anthropic, for example, saw its API inference speeds degrade as its popular Claude models surged in usage. This performance drop revealed a structural inference constraint, forcing them to batch users at higher rates to manage load. The bottleneck is even more pronounced for Chinese labs like DeepSeek, which were directly impacted by US export controls targeting critical inference chips.

FIGURE 8: WHO'S WINNING THE LLM RACE?



Source: Stifel IRIS, LMArena



## 1.2 HOW TECH BEHEMOTHS DOMINATE AN ERA OF AI MAXIMALISM

In AI's maximalist era, tech giants solidify dominance through concentrated market shares across the stack, from Nvidia in GPUs and TSMC in leading-edge chip manufacturing, to hyperscalers in cloud infrastructure and leading labs in foundational models. Surging capital intensity and vertiginous infrastructure outlays raise entry thresholds and encourage partnerships, from trillion-dollar-scale commitments and cross-investments, to manage counterparties' risks and support a connected value chain.

### THE HYPERSCALERS' INFRASTRUCTURE PIVOT

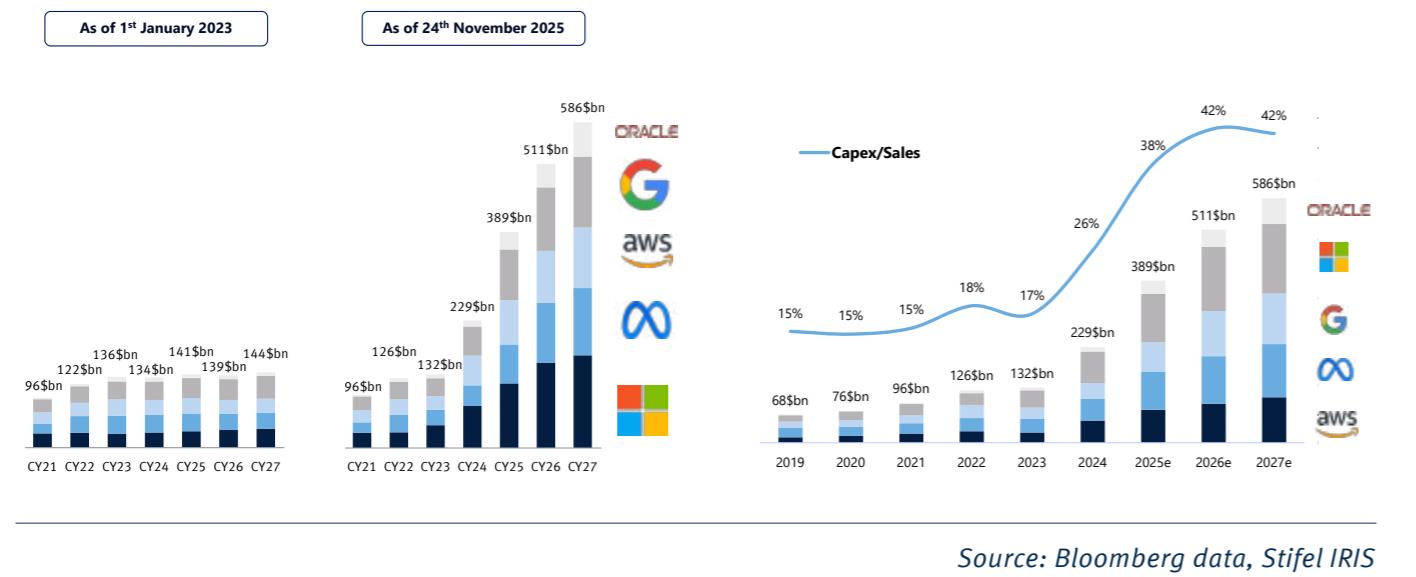
In the AI race for supremacy, software-first giants are reallocating capital towards physical infrastructure at scale. This marks a sharp reversal of recent expectations. After the 2020-21 pandemic-era investment boom, the 2022-23 narrative was dominated by the “year of efficiency”, widespread layoffs, and a focus on cost discipline. The prevailing consensus was that hyperscalers’ capex had peaked and that their capital intensity would decline, unlocking operating leverage.

The generative AI boom has upended this thesis. We observe a structural step-change in capital intensity, with US hyperscalers’ capex-to-sales ratios set to surge from a pre-2024 average of 16% to 26% in 2024, followed by a sustained 37% to 43% plateau. This shift represents a 128% increase, and a c.\$1.2 trillion in incremental infrastructure spending over the 2024-27 period compared to consensus forecasts just three years ago. In CY25, US hyperscalers’ capex is projected to exceed that of all global telcos combined, confirming a digital capex supercycle now driven by compute rather than networks. The centre of gravity is shifting from platform dominance in the network era to data centre control in the compute era.

This has ignited a debate on whether the surge is a temporary “AI tax” or a permanent strategic realignment. We believe many still underestimate the durability of this new capital-intensive model. A similar, though more constrained, trend is visible among Chinese tech giants, underscoring this is a global pivot. ByteDance alone is reportedly targeting a capex of nearly \$20bn in 2025. The data centre industry is fundamentally driven by economies of scale; this significant spending creates a formidable capital barrier and unlocks a virtuous cycle of cost advantages (procurement, power, custom silicon), effectively entrenching their market dominance.

In our view, this infrastructure-first doctrine appears to be financed by a stark optimisation of labour costs. The era of significant workforce expansion is over. Despite strong double-digit revenue growth, tech giants are now in a phase of recurring, “surgical” layoffs, such as Google’s cuts in its core engineering and cloud teams in early 2025 and Microsoft’s ongoing reductions within its Azure division. This contrast suggests a core strategic trade-off: margins are being protected by optimising human capital to fund an unprecedented ramp up in compute capital.

FIGURE 9: THE GREAT REVISION: A GENERATIONAL STEP-UP IN HYPERSCALE CAPITAL INTENSITY

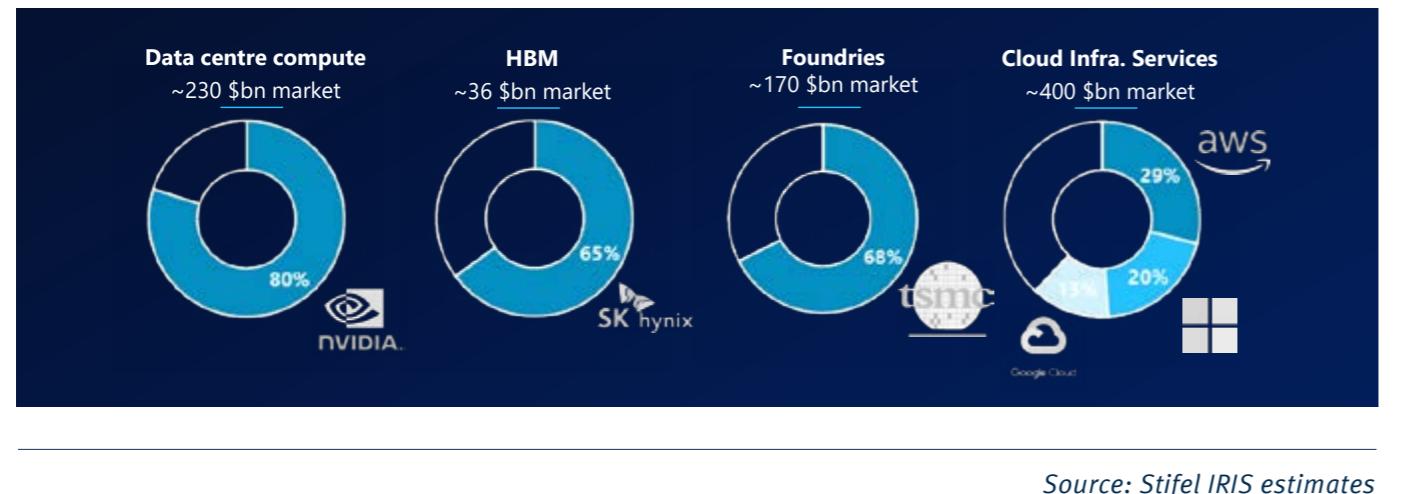


## THE AI STACK IS LOCKED IN OLIGOPOLIES

Our analysis of the AI value chain reveals a market structure defined not by competition, but by a series of interconnected oligopolies. From the silicon foundation to the cloud platforms, a handful of incumbents

capture a disproportionate share of the value, creating formidable barriers to entry and reinforcing a “rich-get-richer” dynamic.

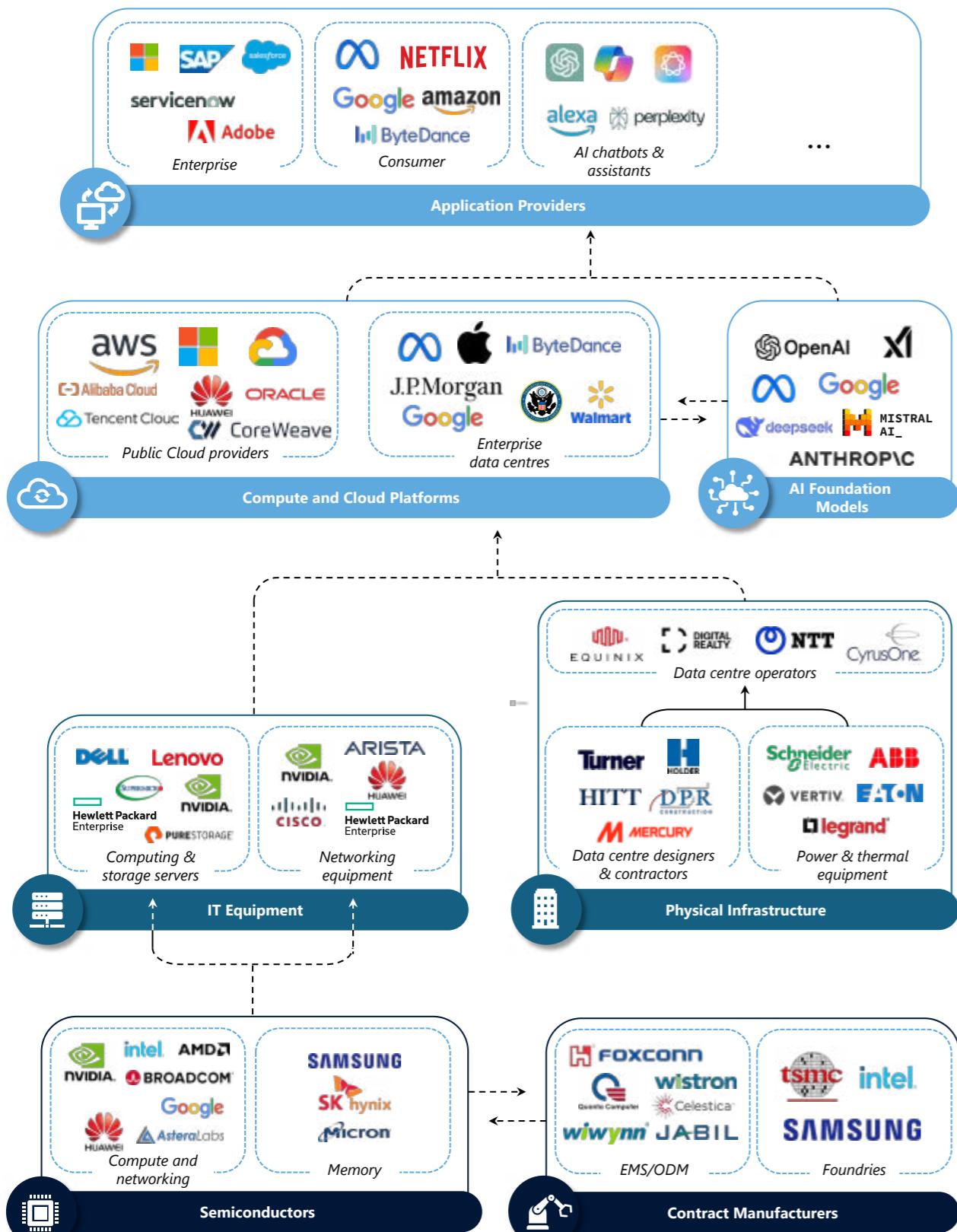
FIGURE 10: 2025E MARKET SHARE IN CRITICAL DATA CENTRE TECHNOLOGY LAYERS



This concentration is evident at many critical layers of the stack:

- **Semiconductors:** This foundational layer provides the raw processing power. The AI accelerator market is dominated by NVIDIA, whose Hopper and Blackwell GPU architectures capture an overwhelming share of revenues. It faces two key challengers: AMD in the merchant market and Broadcom as the established leader in the custom ASIC (Application-Specific Integrated Circuit) space, notably developing bespoke chips for Google's TPU and Meta's AI programs. Most AI hardware startups have struggled to penetrate the market at scale. This concentration extends to critical components. HBM (High Bandwidth Memory, stacked DRAM essential for feeding AI processors) is a tight oligopoly controlled by the established memory giants: SK Hynix, Samsung, and Micron. Crucially, the advanced manufacturing layer is a dual chokepoint. TSMC holds a near-monopoly on fabricating the advanced sub-5nm logic chips and providing the critical 2.5D/3D advanced packaging (like its CoWoS technology) required to integrate these logic chips with HBM.
- **IT and networking equipment:** This layer bundles the components into functional supercomputers. This includes the high-density servers, a market led by ODMs like Supermicro and enterprise giants Dell and HPE, who maintain the closest partnerships with NVIDIA. It also includes the high-performance, low-latency networking fabrics essential for large-scale clusters. This networking segment is itself a duopoly, where we note NVIDIA's proprietary InfiniBand is widely considered the gold standard for large-scale training, challenged primarily by high-speed Ethernet solutions from Arista.
- **Physical infrastructure:** This is the physical shell that provides space, power, and cooling. As AI drives extreme rack densities, the electrical and thermal management systems become critical bottlenecks. The core global market for this mission-critical power and cooling equipment is dominated by a “Big 3” (Schneider Electric, Vertiv, and Eaton), which holds roughly 50% market share. However, the overall physical infrastructure market remains fragmented due to immense product diversity (racks, cabling, busways) and geographical fragmentation.
- **Cloud infrastructure:** The IaaS/PaaS layer, where most enterprises access AI, remains a mature oligopoly. The Big 3 – Amazon (AWS), Microsoft (Azure), and Google (GCP) – account for a formidable 67% of the public cloud market. Oracle has been placing a lot on AI and joined them. However, their dominance in high-performance AI workloads is now being challenged by the emergence of specialised, AI-native “neoclouds” (as discussed in Section 2.2) that compete directly on architecture and performance.
- **Foundation models:** At the very top, the AI frontier itself is not an open ecosystem. As previously discussed, we believe it is defined by the domination of major US AI labs: OpenAI, Google, Anthropic, and xAI. We believe their dominance is to a large extent a direct result of their preferential access to the significant, proprietary compute clusters that the underlying oligopolies provide.

FIGURE 11: FROM SILICON TO SERVICES: MAPPING THE COMPLETE AI & DATA CENTRE VALUE CHAIN



Source: Stifel IRIS

## A TIGHTLY KNIT ECOSYSTEM

The AI infrastructure landscape is not just dominated by a few giants; it is an increasingly tightly-knit and interconnected ecosystem. This environment is characterised by a dense web of cross-investments, vendor financing, and strategic partnerships that blur the lines between customers, suppliers, and competitors. This interconnectedness has led some commentators to describe the relationship as increasingly insular in the AI economy, and sparked debates about the risk of it becoming circular.

Yet, this dynamic serves a dual purpose: it aims to reassure counterparties and mitigate risks in a highly capital-intensive market. It also creates systemic vulnerabilities, as a significant setback at one major AI firm could amplify shockwaves throughout the entire industry.

The system relies on large-scale financial and operational linkages between the infrastructure owners and the AI model developers:

- Nvidia's strategic financing:** Nvidia is a focal point of this integration, with its client base highly

concentrated among hyperscalers (Microsoft, Amazon, Meta, Google) who account for over half of its data centre sales. This interdependence is formalised through sizeable deals and strategic equity stakes. For instance, NVIDIA frequently uses its financial strength to strike large cloud-computing deals with key partners and clients in which it holds an equity position, such as its deal with CoreWeave. This practice effectively links vendor sales volumes to capital investment flows, ensuring long-term hardware offtake.

- The Hyperscaler-AI lab alignment:** Foundational model developers are structurally tied to the major cloud platforms. Microsoft's investment in OpenAI secured its position as its primary cloud provider and user of its technology. Amazon has invested billions in rival AI model developer Anthropic, ensuring that the model's consumption feeds directly into its respective cloud arms. Google and Meta favour an in-house approach, though Google has also invested more than \$8bn in Anthropic to date.

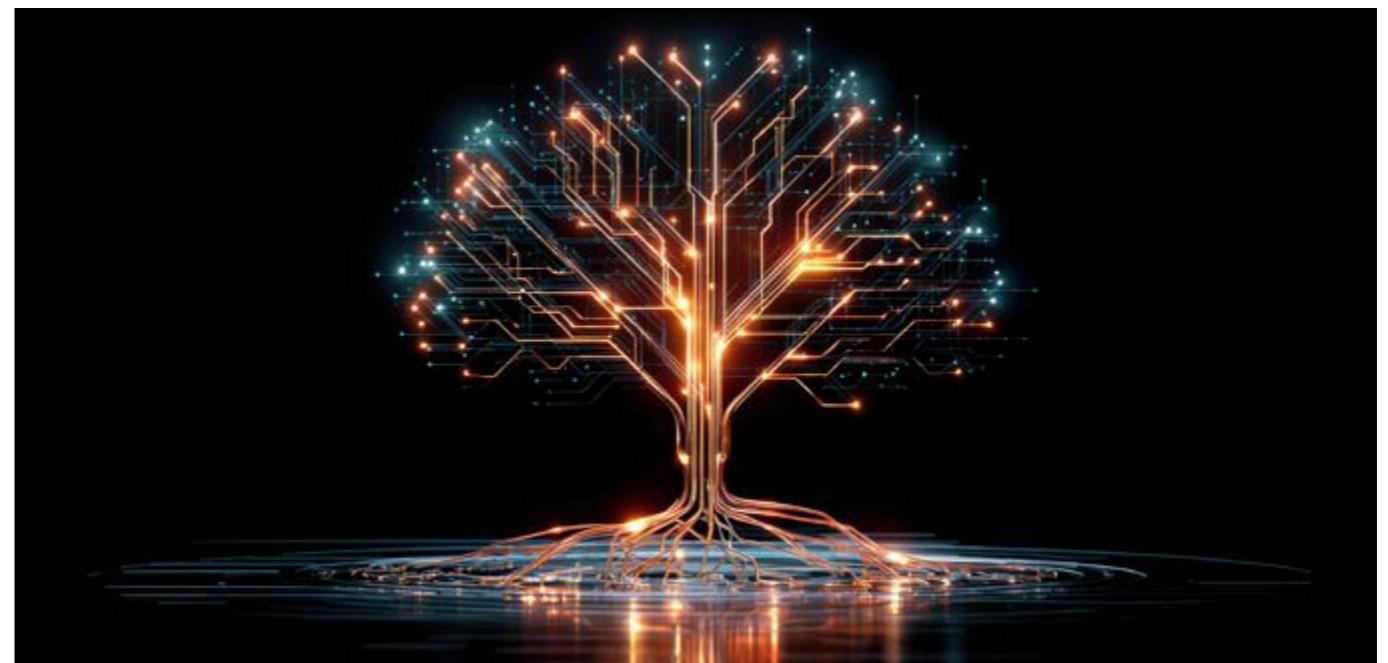
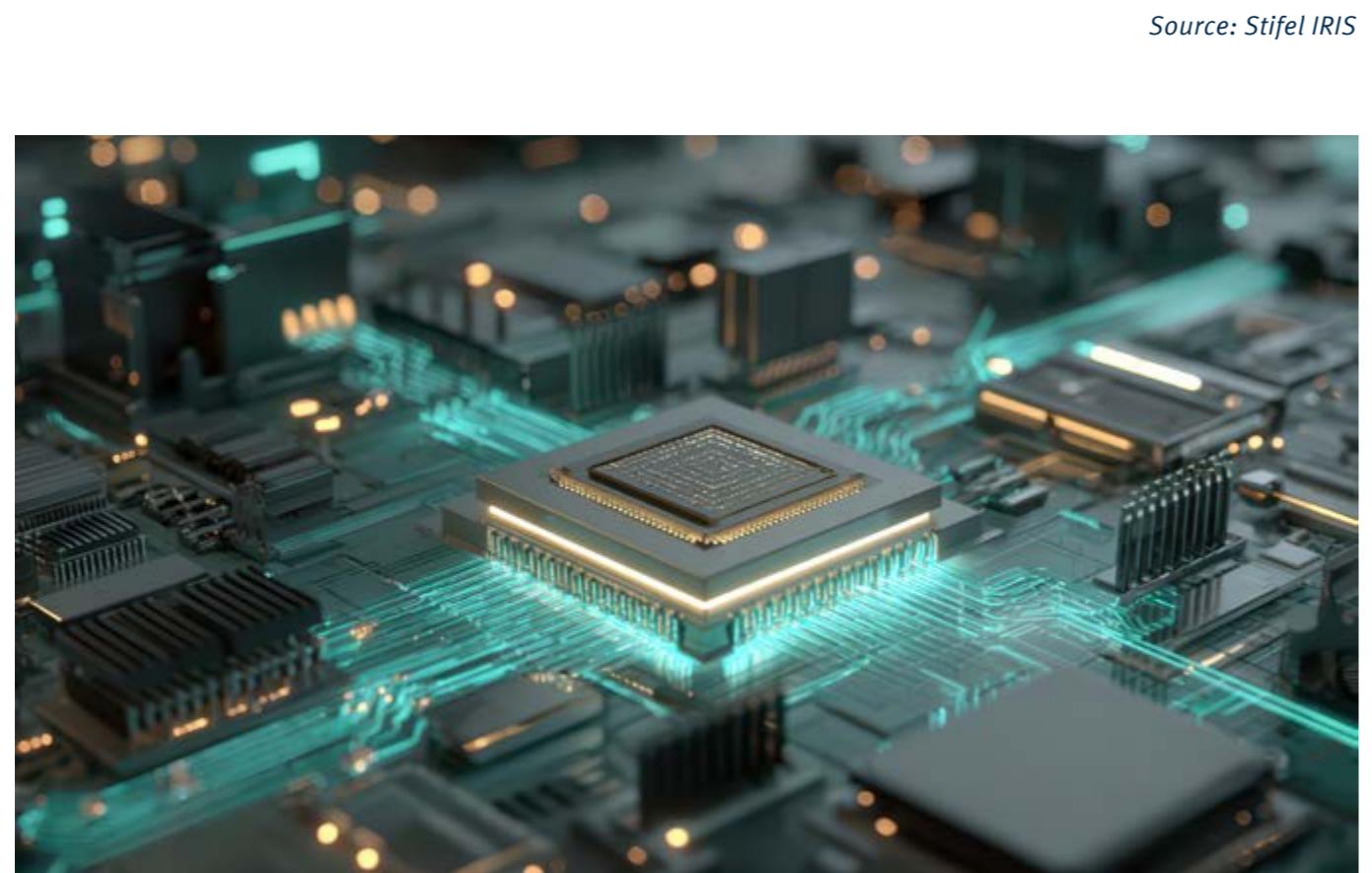
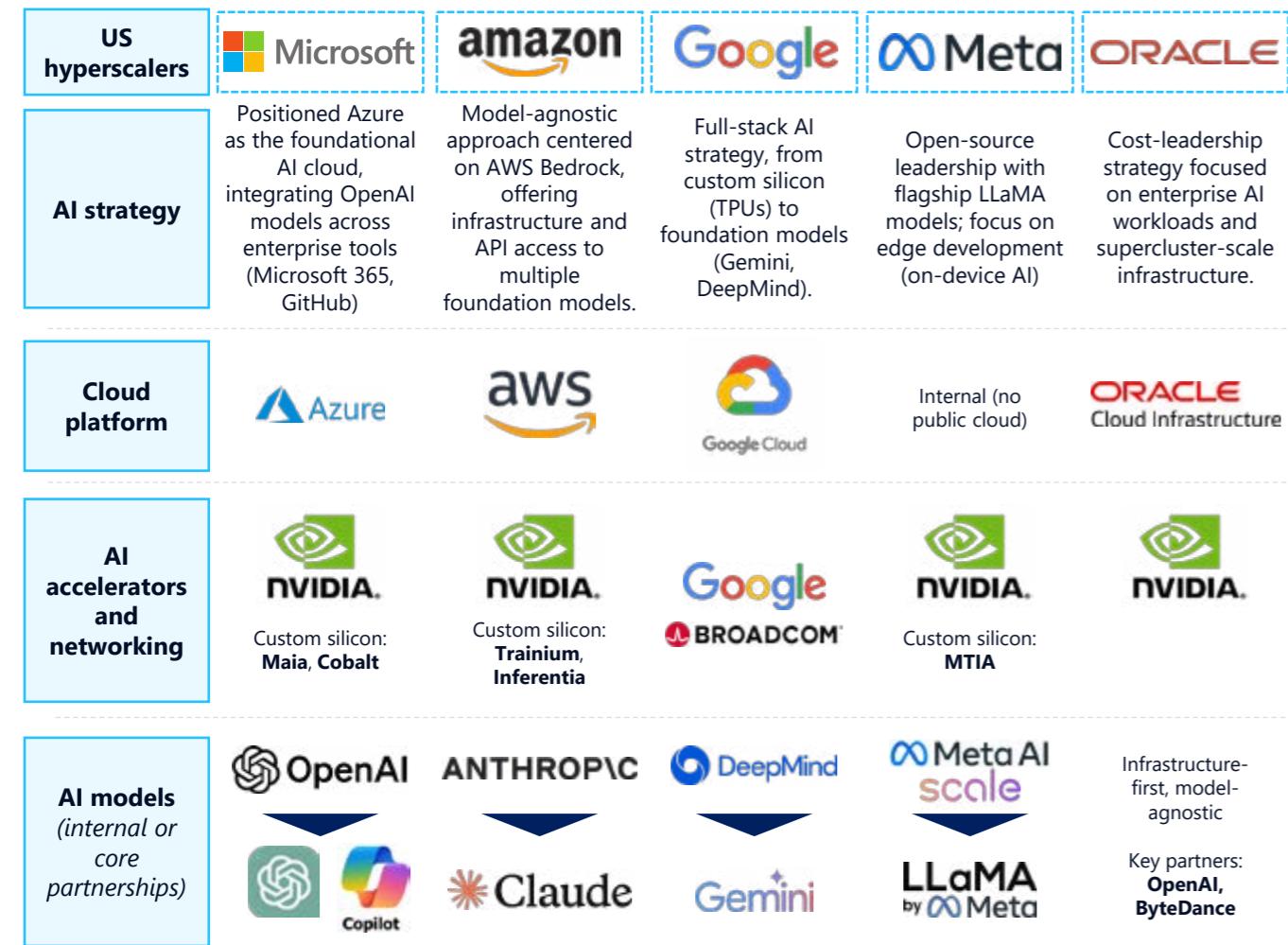


FIGURE 12: US TECH CONGLOMERATES AND THEIR VERTICAL AI STACK STRATEGIES

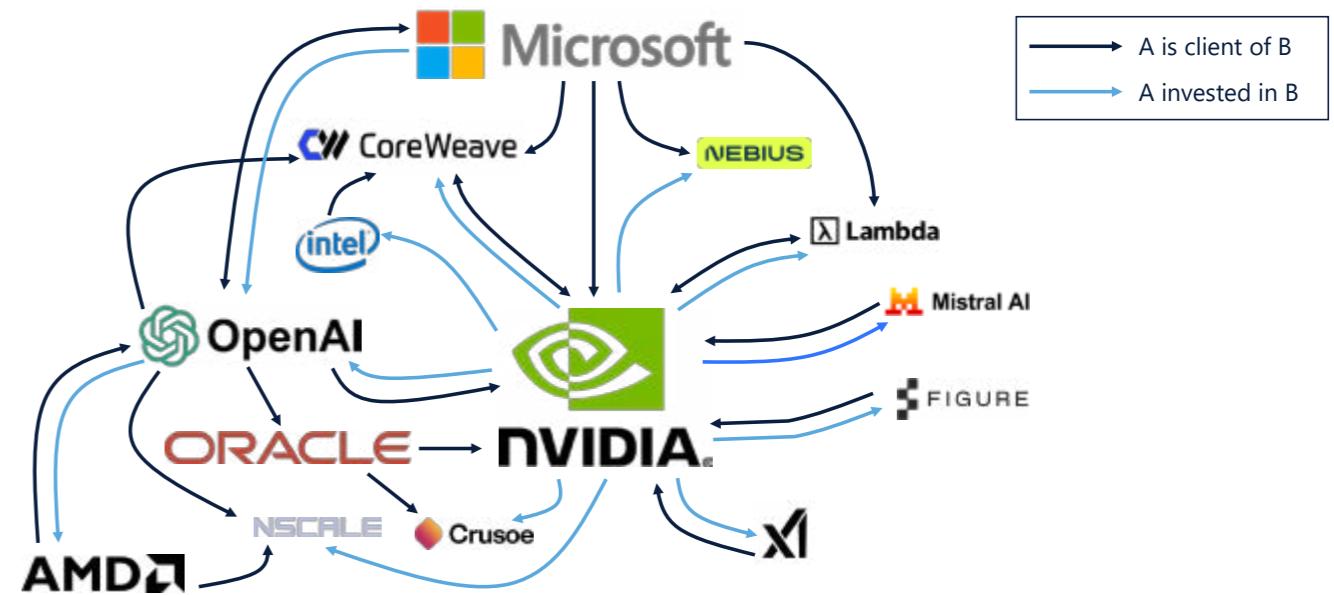


The most acute illustration of this financial architecture is seen in OpenAI's ambitious infrastructure planning. To secure the capacity required for its research roadmap, the company has had to engage in a series of unprecedented financial commitments with nearly every major technology firm (Nvidia, AMD, Oracle, Broadcom, Google, Samsung). OpenAI's recent compute contracts, which commit the company to securing future capacity as it is developed, with the total close to \$1.4tn as of early November 2025. These sizeable commitments are structured as incremental, capacity-contingent payments, effectively hedging on future revenue growth to meet infrastructure supply targets. The recent \$38bn deal with Amazon Web Services to access their infrastructure and NVIDIA chips lessens OpenAI's dependence on its primary backer, Microsoft, highlighting a strategy to secure resource redundancy across the key cloud oligopolies.

This accelerated build-out is introducing high levels of financialisation into the infrastructure market, most

notably through increased debt financing. Initially, a key argument differentiating the AI boom from the telecom bubble was the balance sheet strength of hyperscalers and other major AI participants. However, we are now seeing a trend towards greater debt, even among the largest firms. The substantial capital required for AI infrastructure is driving a debt surge across the sector. For instance, Alphabet recently sought to raise approximately \$22bn through debt, following a \$30bn bond sale by Meta, as companies increasingly rely on capital markets rather than FCF to finance capex. To date in 2025, bond issuance from AI giants has already exceeded \$200bn. With AI-related debt representing a growing share of net supply, we note that this perceived re-leveraging is sparking more debate in credit markets. This growing dependence on debt to finance expansion raises concerns. If AI demand fails to meet expectations, the associated financial risks could spread beyond the tech sector and impact the wider financial system.

FIGURE 13: MAPPING STRATEGIC INVESTMENTS AND CLIENT DEPENDENCIES IN THE AI ECOSYSTEM



Source: Stifel IRIS

# 1.3 COLLIDING WITH THE POWER WALL

The AI buildout is on a collision course with physical limits. As facility-level energy savings plateau, global data centre capacity is set to double by 2030, triggering a structural power bottleneck. Simultaneously, rising AI rack densities render traditional air cooling obsolete. This mandates a non-negotiable, industry-wide shift to liquid cooling as the new baseline for performance.

## THE END OF EASY EFFICIENCY

To understand the data centre energy dilemma, one must first understand its core metrics. Today, the industry's physical scale is no longer measured in square footage, but in power capacity (Megawatts, or MW), as this dictates the amount of IT equipment a facility can support. Convention focuses on IT power (or critical power), which is the electricity allocated only to compute and storage hardware, distinct from the power used for cooling and other infrastructure. The ratio between these two is the Power Usage Effectiveness (PUE), the sector's key efficiency benchmark. A PUE of 2.0 means that for every watt used for IT, another watt is consumed by cooling and overhead. A perfect PUE of 1.0 means zero overhead.

For the past 15 years, the data centre industry achieved a decoupling of capacity growth from energy consumption, driven by two engines of efficiency:

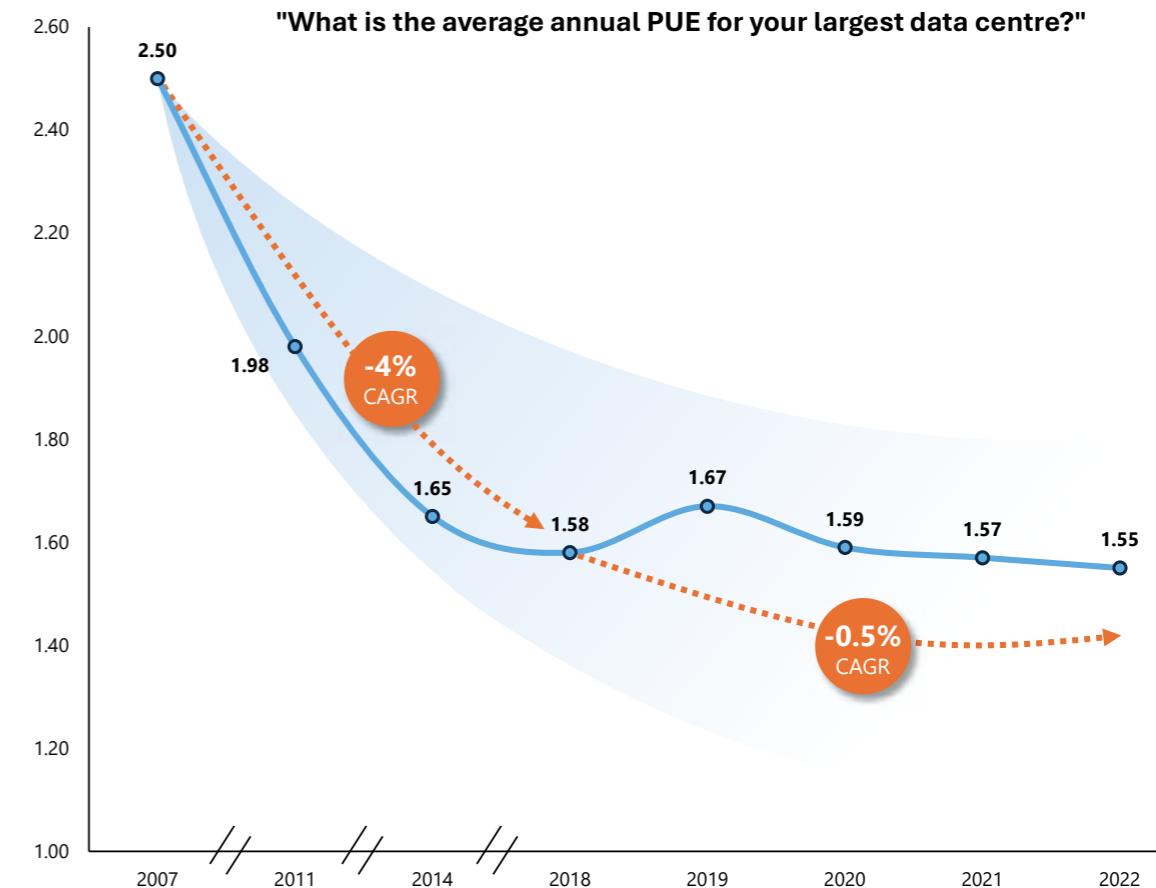
- Continuous innovation in facility-level infrastructure. More energy-efficient electrical systems (UPS, switchgear) and vastly improved air-cooling technologies drove the industry-wide average PUE down from a highly inefficient 2.5 in 2007 to a 1.55, according to Uptime Institute data.
- The second, and perhaps more powerful, engine was the significant migration of workloads from legacy,

high-PUE (>2.0) on-premise servers to state-of-the-art hyperscale cloud facilities. These hyperscale centres are engineering marvels, operating at PUEs far below the industry average, often sub-1.4, with leaders like Google reporting figures as low as 1.10.

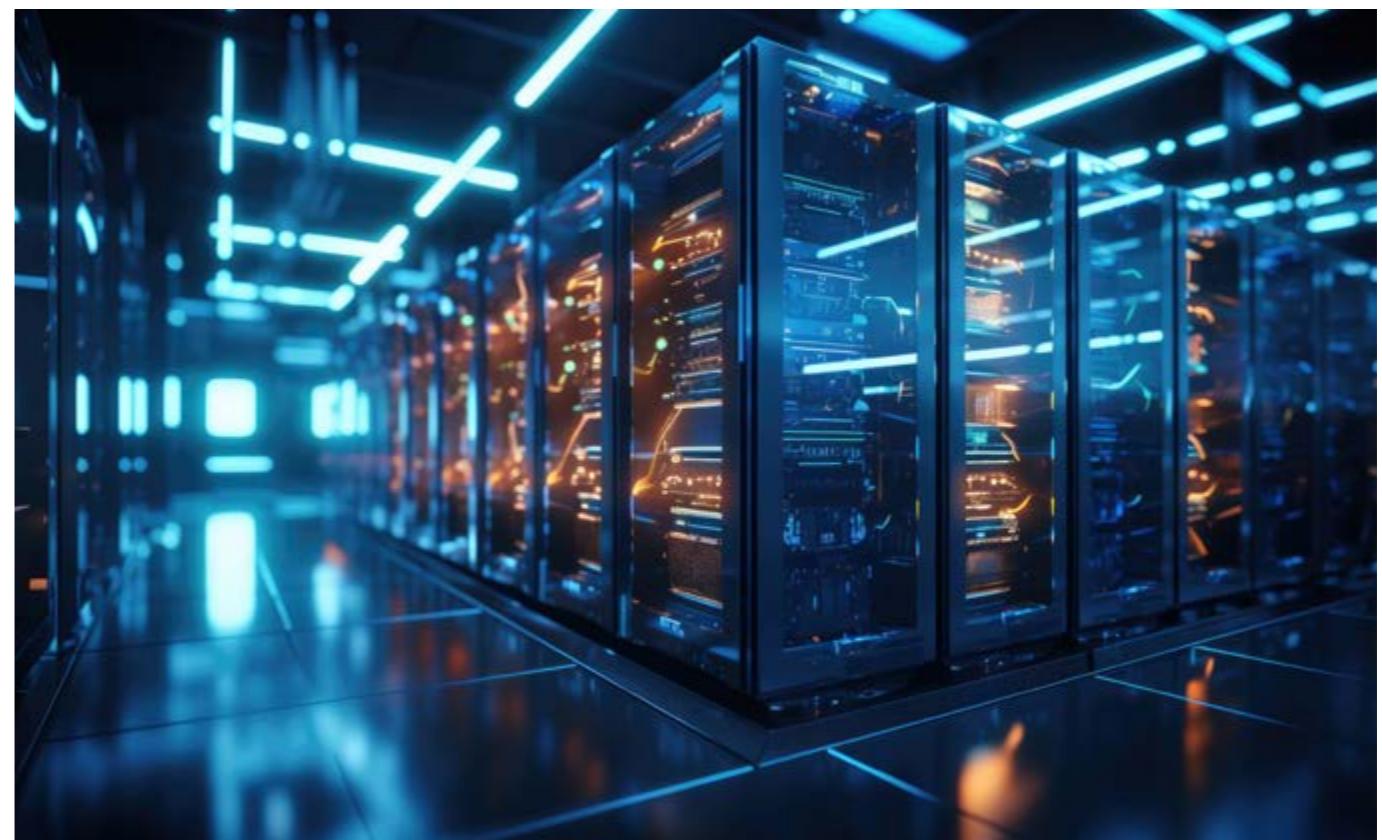
Together, this dual movement allowed compute capacity to explode while curbing total power consumption. That era appears to be over. Both engines of efficiency are now stalling. The structural gains from cloud migration are maturing. More critically, facility-level PUE gains have hit a wall of diminishing returns. The industry's last major efficiency lever, the shift to liquid cooling, is now being pulled. While its primary role is to enable the extreme compute densities required by AI (a topic we explore next), its direct impact on efficiency represents the final step in facility optimisation. It offers a last, incremental PUE reduction (estimated at 0.2-0.3) on an already highly-optimised hyperscale design.

With no other major facility-wide efficiency levers remaining, the implication is stark: the “free lunch” is over. The industry is entering a new era where, for the first time, growth in data centre power consumption (TWh) will track the growth of physical IT capacity (MW) on an almost one-to-one basis.

FIGURE 14: DATA CENTRE AVERAGE ANNUAL PUE: PROGRESS IS STALLING



Source: Uptime Institute, Stifel IRIS



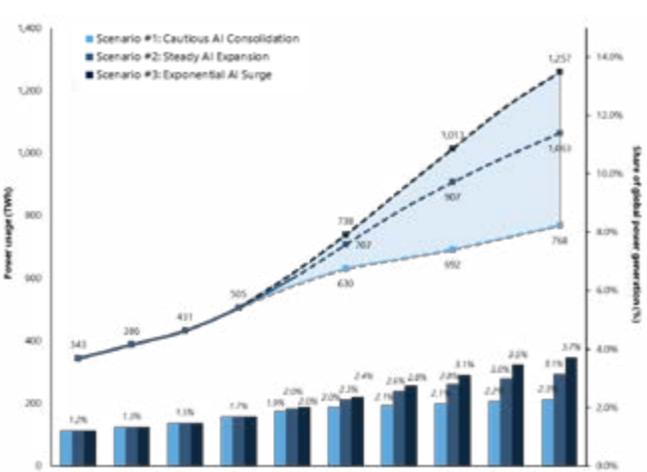
This new paradigm where consumption tracks capacity collides directly with an unprecedented expansion of physical infrastructure. As a baseline, installed data centre capacity stood at approximately 57GW in 2023, with AI workloads accounting for only 4.5GW, or less than 8% of the total. This figure effectively represents the starting line before the full force of the AI-driven capex cycle. Estimates of global capacity vary by ±10GW, reflecting both methodological differences and the inherent opacity of data centre markets. However, multiple forecasts converge on a doubling, or even tripling, of capacity by 2030. Omdia projects an additional ~90GW of installed capacity between 2023 and 2030, exceeding 160GW.

This direct link between capacity and consumption reframes the energy forecasts. Based on our modelling (from our previous report published in February 2025), we project global data centre power demand to surge by ~150% from 2023 to 2030, reaching 1,063 TWh in our base case scenario. This trajectory would see the sector's share of global electricity demand double from ~1.5% to 3.1% by the end of the decade. However,

given recent trends, we believe the sector may now be tracking closer to our high scenario, in which its share could rise to 3.7%

While these absolute figures often fuel alarmist narratives, we find the primary challenge is not the total global consumption; a 3.1% share, while significant, is manageable on a planetary scale. The true crisis stems from two structural mismatches: concentration and speed. This demand is not evenly distributed; it is intensely focused on specific regions, particularly in the US, where data centre consumption is projected to surge from less than 5% in 2023 to nearly 11% of the total national grid by 2030. Furthermore, data centres are built at venture speed (18-24 months), but the energy infrastructure required to power them operates on five-to-ten-year timelines. This asynchronism is the structural source of the current scarcity. It has fundamentally inverted the logic of site selection: where land and latency were once key decision factors, utility power availability is now the primary constraint dictating where, and if, new data centres can be built at all.

FIGURE 15: AI & DATA CENTRE POWER DEMAND SET TO EXCEED 1,000TWH/YR BY 2030



Source: Stifel IRIS, from "Energy-Aware AI", February 2025

## REDEFINING INFRASTRUCTURE FOR EXTREME DENSITY

While IT equipment (the GPUs, servers, and networking gear) understandably captures the spotlight, it is the underlying physical infrastructure that defines a data centre's capacity. For a new-build, this "facility" infrastructure (building, electrical systems, and cooling) typically costs USD 7-10m per MW. This figure varies by region but is most sensitive to scale: large hyperscale facilities benefit from significant economies of scale, making their cost-per-MW often far lower than smaller sites.

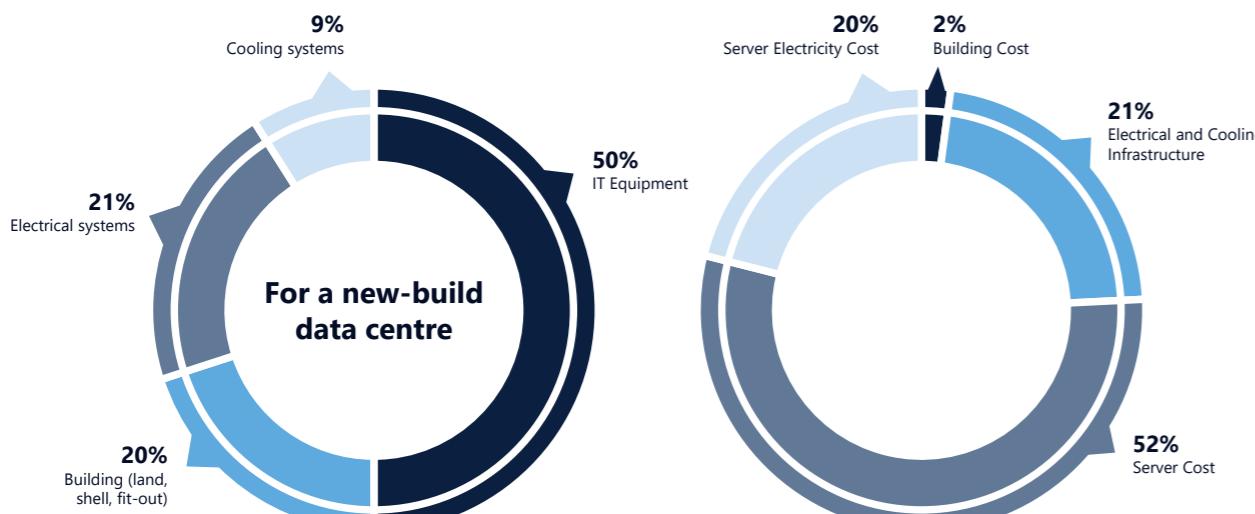
In a traditional data centre, this infrastructure accounts for approximately 35% of total capex, versus 65% for the IT equipment. This share is often lower in AI data centres, dropping to ~20%. This shift is not because the facility is cheaper; rather, it is because the AI-optimised IT hardware is vastly more expensive, roughly doubling the cost-per-MW of the IT load while the underlying electrical and thermal infrastructure cost rises only incrementally. This distinction is critical given their divergent lifecycles: IT hardware is refreshed on a rapid three-to-six-year cycle, whereas the facility infrastructure is a foundational, 15-to-20-year investment. The initial design choices are therefore structural, defining the facility's capabilities for decades, even as retrofits remain a costly possibility.

Data centres hinge on specialised energy and thermal management systems to ensure uninterrupted

operations. These systems, essential for uptime, fall into four primary categories. First, the grey space encompasses the back-end electrical infrastructure, representing 30-40% of total electrical and thermal infrastructure costs. This includes switchgear (which protects downstream components from power surges) and Uninterruptible Power Supplies (UPS), which deliver instantaneous battery backup for 5-15 minutes to bridge the gap until generators engage. Second, the white space accounts for 10-15% of costs and includes the IT-room equipment such as the racks themselves and the Power Distribution Units (PDUs) that deliver power to the servers. Third, the cooling system mitigates the heat generated by IT operations; it represents 30-35% of infrastructure costs, a share that is increasing due to the shift towards liquid cooling. Finally, back-up generators, predominantly diesel-powered, represent ~15% of costs and activate within minutes of a failure to sustain critical loads.

While this system is complex, the value is highly concentrated. We estimate five key component categories (chillers, CRAH, UPS systems, switchgear, and backup generators) account for ~70% of the total facility infrastructure expenditure. This entire foundation of conditioned, resilient power is ultimately delivered to the white space, where PDUs feed the individual racks.

FIGURE 16 & 17: THE COSTS OF BUILDING AND RUNNING A DATA CENTRE



Source: Stifel IRIS

This traditional architecture is not adapted to AI. In 2022, before AI adoption accelerated, average rack densities plateaued at ~10kW, though hyperscalers were already operating at 2-3x that level. For years, the industry's push toward higher densities was fragmented. While some hyperscalers, notably Google, were early pioneers of liquid cooling for their custom TPUs, the broader market lacked a unifying standard.

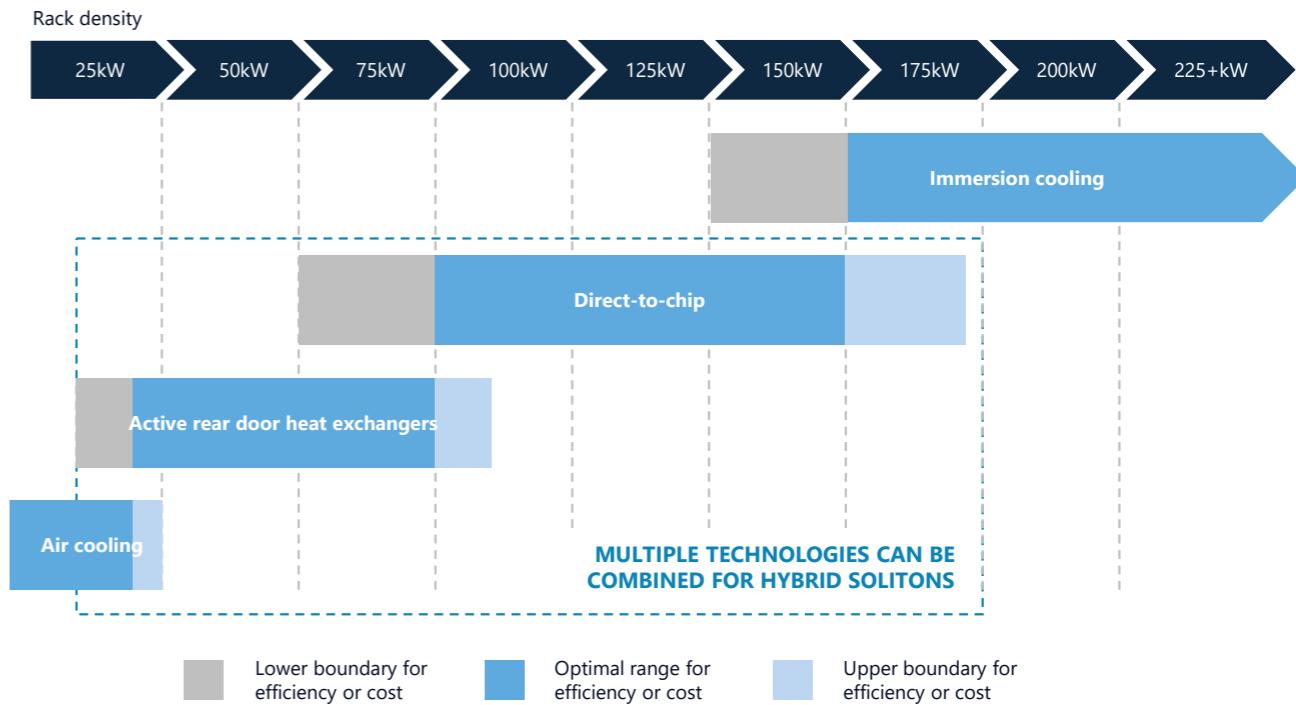
This dynamic was irrevocably altered by the Blackwell platform from Nvidia, the AI market's centre of gravity. For example, Nvidia's GB200 NVL72 rack-scale server, unveiled in March 2024 as part of the Blackwell series,

push total rack density beyond 130kW. It renders traditional air cooling obsolete and makes liquid cooling a non-negotiable baseline for performance.

This shift creates a clear technological line: while legacy technologies like Active Rear Door Heat Exchangers (RDHx) may service mid-density upgrades (often cited up to ~80kW), they are insufficient for this new standard. As the roadmap pushes beyond 150kW, DTC and full-immersion cooling are the only viable paths forward. Blackwell has effectively ended the architectural debate, standardising the high-density, liquid-cooled blueprint for the AI industry.



FIGURE 18: APPLICABLE COOLING TECHNOLOGIES BY RACK DENSITY



Source: Vertiv, Stifel IRIS

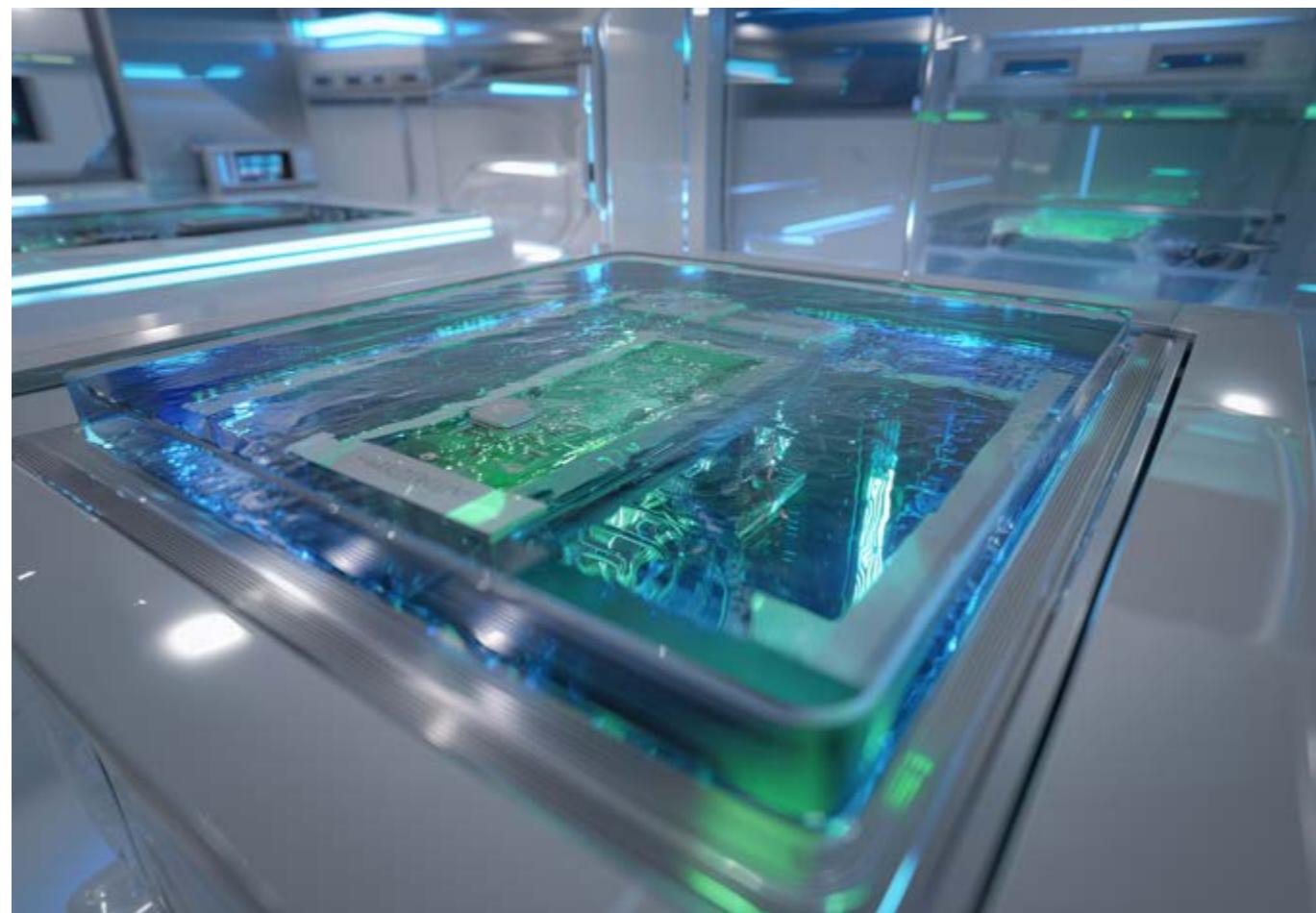
This reality mandates the shift to liquid cooling, which is rapidly moving from a niche solution to an essential enabling technology. The market is standardising around three core methods that manage heat with increasing effectiveness:

- **Direct-to-Chip (DTC):** The current standard for high-performance AI, where liquid coolant circulates through "cold plates" attached directly to the hottest components (GPUs/CPUs) to absorb heat at source.
- **Immersion Cooling:** The most efficient method (PUEs approaching 1.05), where servers are fully submerged in a non-conductive dielectric fluid. It remains a niche solution today but is the primary path for future extreme-density designs.

FIGURE 19: OVERVIEW OF DATA CENTRE COOLING TECHNOLOGIES



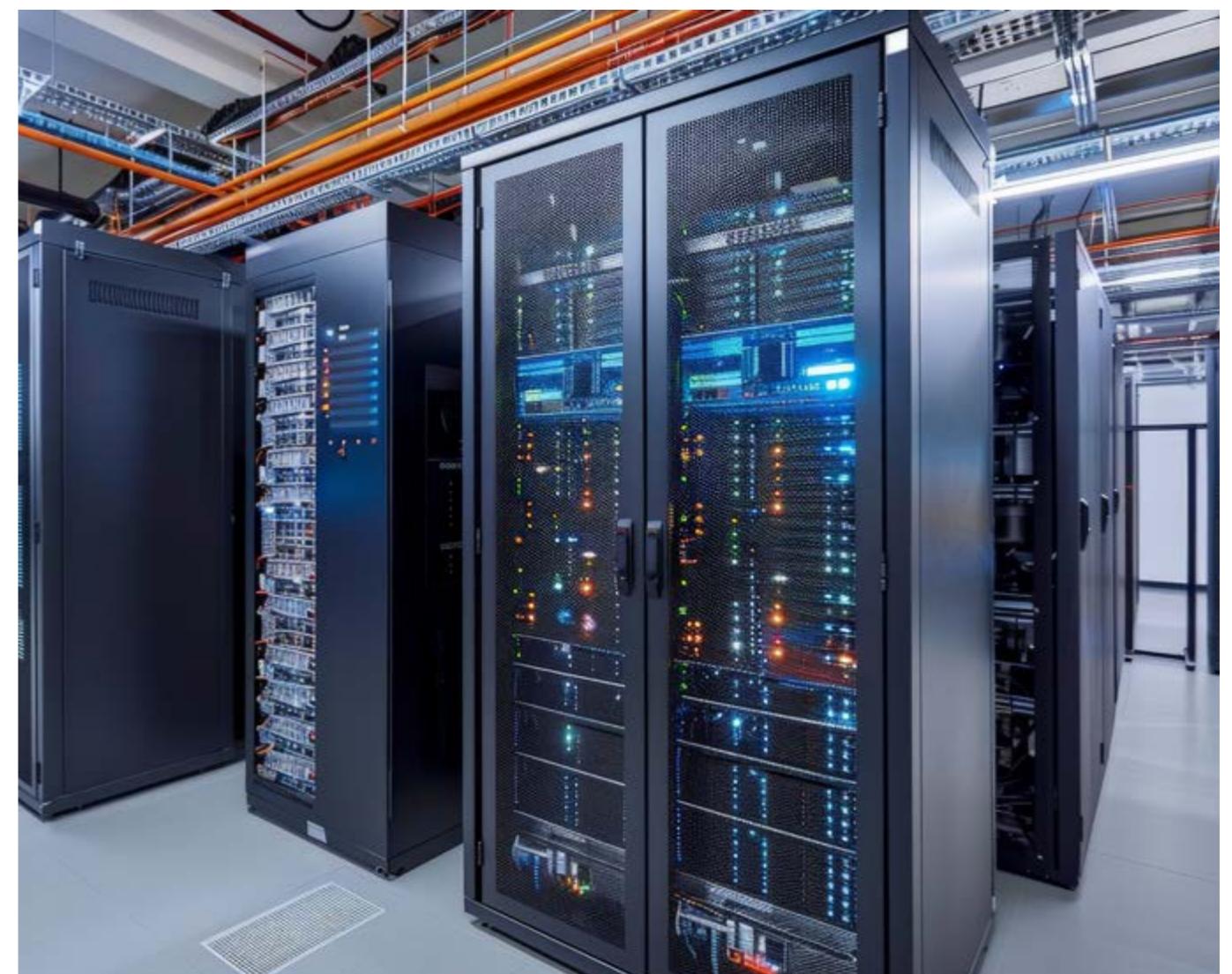
Source: Stifel IRIS



## CARVING OUT VALUE IN THE DATA CENTRE ECOSYSTEM

### SECTION 2

The hyperscalers' dominance on cloud computing, while formidable, is not total. It is a fortress that creates its own periphery. The AI-native stack, demanding operational excellence over mere scale, has created an opening for a new class of specialised neoclouds to thrive. This section unbundles the digital infrastructure stack to map this new ecosystem. It details the five proven bypass playbooks, from low-cost simplicity to sovereign AI, that allow focused competitors to capture value in the fast-growing cloud and data centre market.



## 2.1 UNBUNDLING THE DIGITAL INFRASTRUCTURE STACK

The cloud ecosystem is evolving beyond the simple IaaS/PaaS/SaaS hierarchy. The top 3 hyperscalers consolidate their 67% market share by leveraging a vast range of services, as the \$390bn cloud infrastructure services market re-accelerates, fuelled by AI. The remaining 33% is a complex, multi-layered ecosystem of specialised providers. Identifying value requires deconstructing this stack to find opportunities created by hyperscaler structural gaps.

### THE CLOUD'S BATTLE FOR VALUE CAPTURE

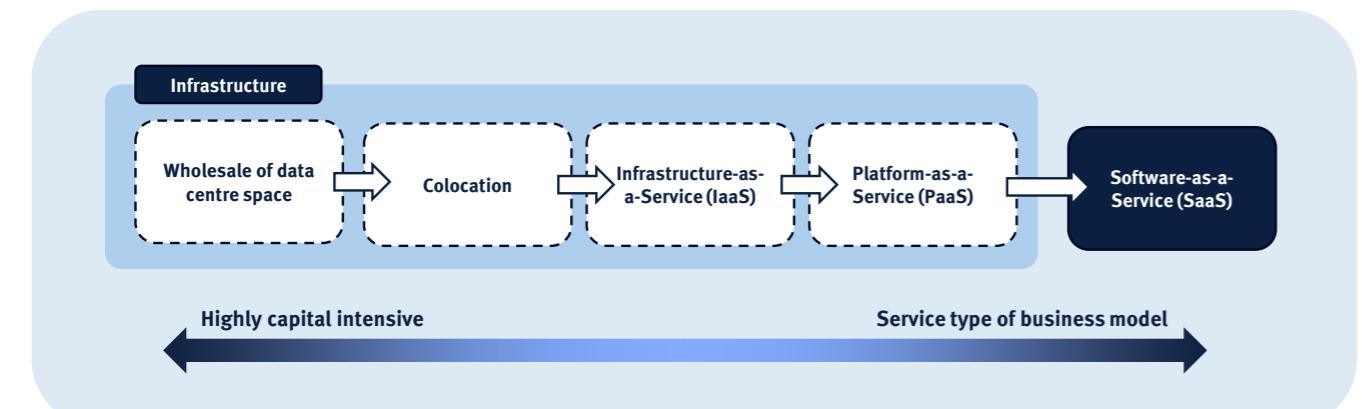
For investors, analysing the cloud market requires moving beyond the static, horizontal definitions of Infrastructure (IaaS), Platform (PaaS), and Software as a Service (SaaS) that defined its first era. While this framework remains a useful descriptive tool, the strategic reality of the 2025 market is one of intense vertical competition and blurring boundaries. The primary battle is no longer within these layers, but across them for ultimate control of developer workflows and enterprise data.

The US hyperscalers AWS, Azure and GCP (Google Cloud Platform) have become IaaS giants, and are aggressively pushing up the stack. They leverage their foundational dominance in compute and storage to offer an ever-expanding portfolio of high-margin, proprietary PaaS

services such as managed databases and AI platforms. The strategic goal is technical and clear: increase value capture per customer and create deep, systemic vendor lock-in at the platform level.

Simultaneously, a new generation of specialised PaaS and tooling providers, such as Databricks and Snowflake, is building on top of the IaaS layer. These platforms effectively abstract the underlying infrastructure, intentionally commoditising the hyperscalers' core product. By winning developer loyalty through superior, focused workflows, they intercept the customer relationship and can arbitrage IaaS providers for cost and performance.

FIGURE 20: THE CONVENTIONAL CLOUD VALUE CHAIN



Source: Stifel IRIS

This strategic conflict has fragmented the monolithic "Cloud" into the multi-dimensional, complex market we see today. The simple public cloud model, once seen as the final destination, is now just one of several core deployment models. We are witnessing a rise in private cloud deployments, as enterprises with predictable, high-volume workloads find it more economical to use dedicated hardware rather than paying the variable premium of the public cloud. To bridge these different environments, the hybrid cloud model has emerged, providing an essential orchestration layer through platforms like Red Hat OpenShift and VMware Cloud Foundation, which allow applications and data to move fluidly between public, private, and edge domains. Building on this complexity, the multi-cloud strategy has evolved from a simple redundancy tactic into an offensive financial decision. It allows firms to strategically select

"best-of-breed" tools for specific jobs, such as using cost-effective, specialised clouds for GPU-intensive training while relying on AWS for core data warehousing. This strategic diversity extends even further, fragmenting the very way services are purchased and secured. In consumption models, traditional virtual machines are now complemented by serverless computing. Platforms like AWS Lambda or Azure Functions radically redefine the opex model, abstracting away the server entirely and enabling billing based on millisecond-level execution time, which is ideal for spiky, event-driven applications.

This multi-axis fragmentation, driven by complex technical, financial, and strategic imperatives, is precisely what creates the structural openings for the specialised competitors, despite the strength of hyperscalers.

FIGURE 21: THE MODERN CLOUD BUSINESS MODEL MATRIX

Physical infrastructure layer		Service abstraction layer		Orchestration & consumption layer	
Data centre format	Hyperscale / Wholesale / Retail	Service model	IaaS / PaaS / SaaS	Consumption mode	On-demand / Reserved / Spot / Subscription
Deployment model	Self-built / Colocation / Leaseback	Cloud model	Public / Private / Hybrid / Multi-cloud	Control model	Fully managed / Customer-managed / Embedded in software stack
Integration	Vertically integrated / Bare metal / Overlay	Offering scope	Compute / Storage / Networking / Sovereign / AI-specialised	Differentiation axis	Sovereignty / Energy source / Latency / Interconnectivity
Geography	Core (FLAP, N. Virginia, Singapore) / Tier-2 / Edge	Consumption granularity	VM / Container / Serverless / Function / Model-level / GPU-by-the-minute	Billing model	Per resource / Per usage unit / Per outcome (e.g. per token, per query)
Power envelope	Standard / High-density / Liquid-cooled	Resilience tiering	Standard / HA / Geo-redundant / Zone-isolated	Platform exposure	API-first / GUI-driven / Abstracted into SDKs
Ownership model	Owner-operator / REIT / JV / Asset-light				

Source: Stifel IRIS

The cloud infrastructure services market is demonstrating robust health, entering a phase of sharp re-acceleration driven unequivocally by the generative AI investment cycle.

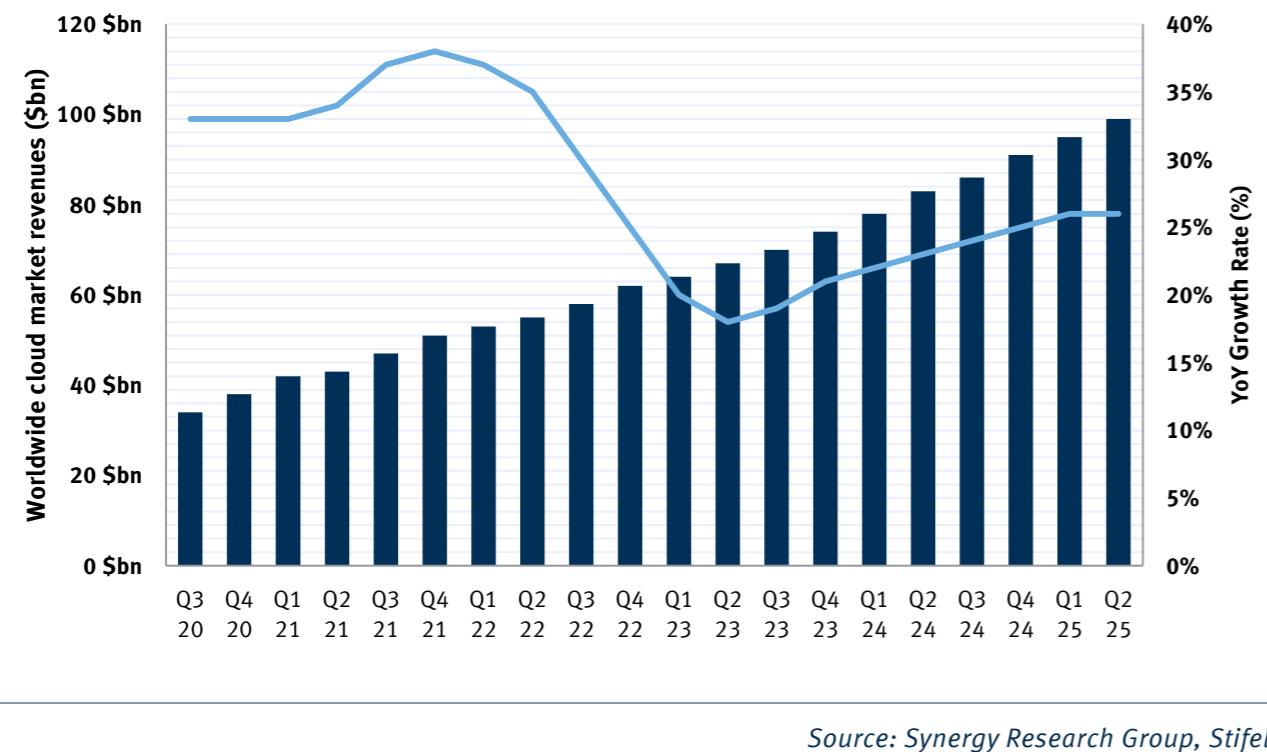
After a multi-year golden age of expansion (pre-2022) where growth rates consistently hovered in the 34-40% range, beginning in the second half of 2022 and proliferating throughout 2023, this momentum cooled significantly. During this period, market growth rates compressed, dipping to an average of 19% for 2023. This slowdown sparked a significant debate: was this merely a cyclical “cloud cost optimisation” phase driven by macro-headwinds, or was it the first sign of structural maturation in the core cloud market?

This trend has now reversed. Before this question could be fully answered, the generative AI investment cycle

arrived as a powerful new catalyst. Beginning in late 2023, the market began a strong rebound, fuelled by significant new spending on AI infrastructure. Q3 2025 marked the fourth successive quarter of accelerating year-on-year growth rates (adjusted for currency) to 28%.

This renewed dynamism highlights the market’s sheer scale. According to Synergy Research Group, quarterly enterprise spending on cloud infrastructure (defined as IaaS, PaaS, and hosted private cloud) reached nearly \$107bn in Q3 2025. This brings the market to a consequential trailing twelve-month revenue of \$390bn. Generative AI has thus acted as a powerful new engine, effectively supersising and reigniting a market that was otherwise beginning to face difficult questions about its long-term growth profile.

FIGURE 22: CLOUD INFRASTRUCTURE SERVICES MARKET GROWTH - IAAS, PAAS, HOSTED PRIVATE CLOUD



## THE CLOUD'S PHYSICAL LANDLORDS

Beneath the intangible cloud stack lies the physical layer: the “data centre landlords” providing the critical foundation of space, power, and cooling. The industry’s roots trace back to the late 1990s and early 2000s, when enterprises began outsourcing server management to avoid the capex and operational burdens of in-house data centres. This birthed the retail colocation model, where businesses rent individual racks, cabinets, or secure cages in multi-tenant facilities. Far from mere space rental, the true value proposition emerged in interconnection: these centres became “digital crossroads” fostering direct peering among carriers, financial institutions, and enterprises, reducing latency and costs while creating network effects. This high-margin, sticky, network-effect-driven business is the fortress upon which giants like Equinix built their global empires.

The 2010s marked a pivotal shift with the cloud computing surge, demanding not bespoke, interconnected racks but significant, undifferentiated scale. This propelled the wholesale leasing model, where a single tenant, often a hyperscaler, leases entire data halls or buildings for raw power and space efficiency at minimal cost. Operators like Digital Realty dominated this space, focusing on cost-optimised delivery rather than interconnection density.

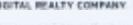
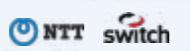
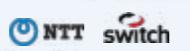
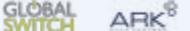
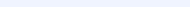
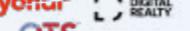
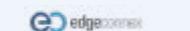
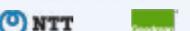
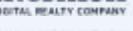
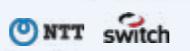
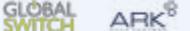
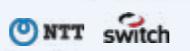
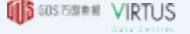
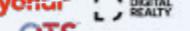
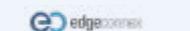
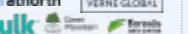
This dual evolution (retail for specialised, multi-tenant needs and wholesale for bulk hyperscale) has created a complex and multi-layered market structure. The top is highly concentrated, with a few global, publicly-traded leaders (namely Equinix, Digital Realty, and NTT) controlling the most critical, interconnection-rich assets and a significant portion of the wholesale market. Beyond this oligopoly, however, the landscape

remains highly fragmented. It is populated by a long tail of regional champions and private operators, many of which are now backed by a wall of capital from private equity and infrastructure funds.

Data centre operators have demonstrated resilience in response to the rise of hyperscalers, largely through strategic collaboration. Hyperscalers do self-build extensively, but reliance on developers for agility is evident in partnerships and financing models. Firms like STACK Infrastructure or EdgeConneX operate as build-to-suit partners for the tech giants. They serve a vital function by arbitraging the hyperscalers’ core strategic trilemma: speed-to-market (especially in an AI race that cannot wait for multi-year build cycles); capital allocation, allowing tech giants to spend their billions on GPUs rather than concrete; specialised expertise, as developers provide essential local knowledge regarding engineering, supply chains, and regulatory compliance. This highlights the central market tension: while hyperscalers can self-build, they also rely on this specialised developer ecosystem to preserve capital and agility.

The AI wave has further bifurcated the data centre market. While incumbents adapt, a new breed of HPC/AI pure-play operators is changing the rules. Many of these new entrants have roots in crypto mining, forged in the hyper-competitive digital asset markets. They leverage a unique DNA: treating low-cost energy as a core competency, controlling the full infrastructure stack from power to GPUs, and using agile financing to scale hardware procurement. This vertically integrated, energy-first model is redefining the economics of AI infrastructure, challenging traditional data centre operators.

FIGURE 23: GLOBAL DATA CENTRE OPERATOR AND DEVELOPER LANDSCAPE

					
Core activity	Interconnection	Scale colocation	Hyperscale developers	AI/HPC	Regional & Edge
Primary customers	Operating carrier-neutral hubs designed for traffic exchange. Core focus on minimising latency through direct physical interconnections.	Operating large multi-tenant data centres offering scalable power and space, from retail cages to wholesale halls.	Securing land and power at scale to deliver rapid, build-to-suit capacity. Arbitraging deployment speed and capital efficiency for hyperscalers.	Delivering AI-optimised infrastructure and/or immediate massive power availability. Focus on high density readiness and time-to-power.	Bridging core hubs to end-users. Operating distributed capacity across Tier-2 metros and proximity nodes to optimize latency and backhaul costs.
Data centre profile	"Carrier Hotels" located in dense city centres. Prioritize fibre density and network redundancy over massive power scale.	Located in Tier-1 and Tier-2 metro areas. Facilities offer a versatile balance of connectivity, security, and power.	Large campuses (100MW+ to GW-scale) located in power-rich zones. Designed to scale from cloud to AI densities.	Purpose-built for high density (>40kW) and liquid cooling. An emerging asset class with evolving technical standards.	Ranging from mid-sized facilities in secondary cities (1-10MW), to micro-modular units (<500kW) at the extreme network edge.
Commercial Model	High-margin model driven by connectivity fees (cross-connects) and high customer retention due to network stickiness.	Flexible model selling private suites or secure cages via medium-term contracts.	Based on long-term, secured leases (10-15 years) with investment-grade tenants.	A nascent, supply-constrained market with pricing reflecting scarcity. Long-term contract models are still maturing.	Retail colocation dominance with high unit pricing based on proximity and connectivity.
Key operators/developers	   	          	       	              	          
	   	          	      	              	          

*Source: Stifel IRIS . Many operators are diversified across segments; classification is illustrative only.*



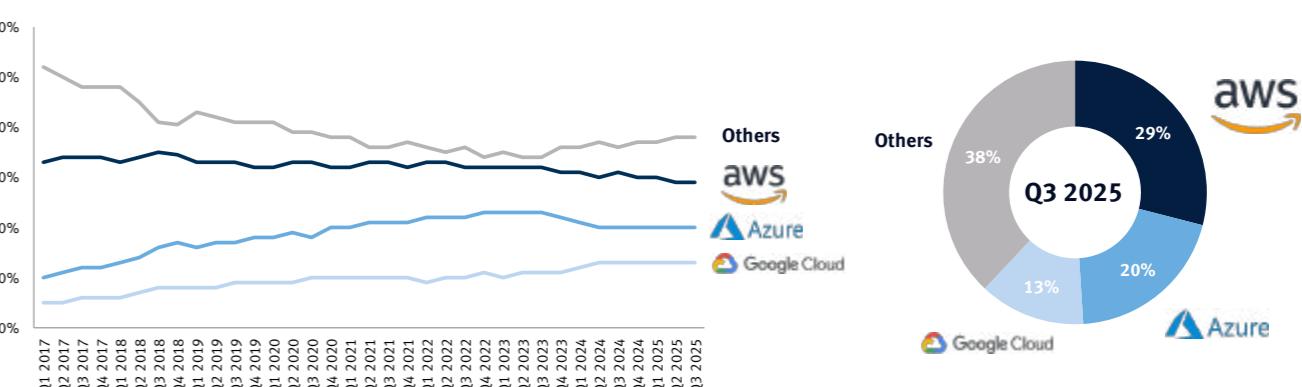
# THE HYPERSCALE DOMINANCE QUANTIFIED

The narrative of hyperscaler dominance is not an abstract thesis, it is a statistical reality. The cloud market remains structurally defined by an entrenched oligopoly. Q2 2025 data from Synergy Research Group confirms this structure. In the total worldwide cloud market (defined as IaaS, PaaS, and hosted private cloud), Amazon, Microsoft, and Google command 29%, 20%, and 13% market share, respectively. This concentration is even more acute in the core public cloud segment, where the top three account for a formidable 67% of the market.

Since early 2017, the consolidation of market share by hyperscalers has been evident, with AWS, Microsoft Azure, and Google Cloud steadily gaining dominance.

However, starting from early 2023, these giants have begun to witness a slight stabilisation and even a decline in market share. Within the hyperscaler dynamic, there has been a notable internal shift, particularly with Microsoft Azure and Google Cloud gaining traction. Azure's share expanded from approximately 10% in early 2017 to 20% by Q3 2025, attributable to its strength in hybrid cloud solutions and extensive enterprise integrations, particularly with Microsoft 365. Google Cloud also improved its share, climbing from about 5% in 2017 to 13% in 2025, driven by its strengths in data analytics and AI offerings.

FIGURE 24: CLOUD PROVIDERS MARKET SHARE TRENDS - IAAS, PAAS, HOSTED PRIVATE CLOUD



*Source: Synergy Research Group, Stifel IRIS*

This 67% market share figure, however, obscures the complex and highly dynamic battle for the remaining 33%. This “other” market is not a homogenous periphery; it is a multi-layered ecosystem. Understanding the market requires deconstructing this segmentation, which we see as comprised of several distinct strategic layers.

At the foundational infrastructure layer, several strategic groups compete for workloads:

- The first layer consists of enterprise challengers, large technology firms like Oracle, Alibaba, Tencent, and IBM. They leverage significant IaaS/PaaS offerings to compete, typically by targeting their existing enterprise customer base or specialising in hybrid cloud and specific geographic strongholds.
- Distinct from these giants are a group of specialised and regional IaaS providers. This diverse group, including European providers like OVHcloud and Hetzner alongside developer-focused platforms like DigitalOcean, deliberately chooses not to compete on the hyperscalers’ terms. Instead, they build successful businesses by excelling in specific niches. Their differentiation is clear and compelling: some offer superior price-performance on raw compute and storage, appealing to cost-sensitive customers. Others, like DigitalOcean, focus on a simplified, developer-centric user experience that removes the complexity inherent in hyperscaler platforms.
- A distinct group is formed by large telecom participants, which often encompass characteristics of two of the previous categories. They leverage their core network infrastructure and deep-rooted enterprise and public sector relationships to compete. Their primary strategy is not to rival hyperscalers on raw IaaS, but to offer managed

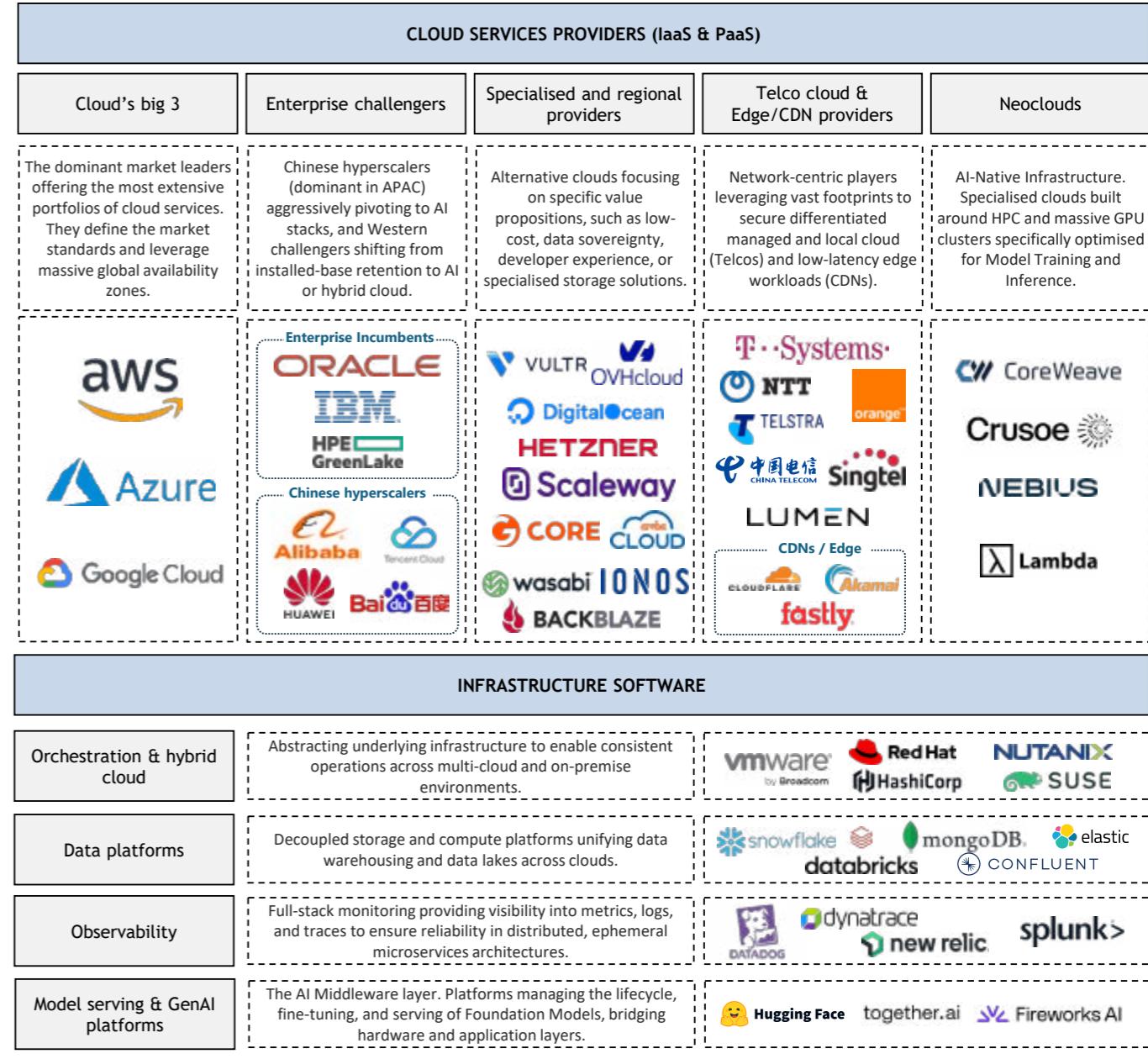
multi-cloud services and, most critically, sovereign cloud platforms that guarantee data residency and compliance within specific national or regional jurisdictions.

- Finally, it is essential to acknowledge the emergence of AI cloud specialists. As we will discuss in greater detail later, these providers are carving out a high-value niche by offering purpose-built, GPU-accelerated infrastructure optimised for AI workloads.

Sitting above this infrastructure is a powerful abstraction layer composed of advanced PaaS and SaaS platforms. Data platforms like Snowflake and Databricks, along with application deployment platforms such as Vercel, create a unified control plane for developers and data scientists. By abstracting away the underlying IaaS, these platforms commoditise the infrastructure providers and mitigate vendor lock-in, shifting the point of strategic control up the technology stack. They enable organisations to operate seamlessly across different cloud environments without being tied to a single provider’s ecosystem.

This entire multi-cloud ecosystem is unified by a horizontal layer of cloud-agnostic management and automation tools. This category, which includes observability platforms like Datadog, security solutions from firms like Palo Alto Networks, and infrastructure-as-code tools like HashiCorp’s Terraform, is indeed an industry standard in practice, if not by a single formal name. These tools provide the essential “picks and shovels”, enabling consistent operations, security, and governance across disparate environments, from hyperscalers to the most specialised providers. They are the connective tissue that make a true multi-cloud strategy viable.

FIGURE 25: CLOUD MARKET SEGMENTATION BEYOND THE HYPERSCALERS



Source: Stifel IRIS

## 2.2 THE AI-NATIVE CLOUD ARISES

The advent of AI clashed with legacy cloud stacks built for flexibility, creating a gap for a new class of specialised “neocloud” providers. While their initial edge was hardware access and aggressive pricing, sparking a commoditisation debate, the market has evolved. Key providers now build deep software moats to guarantee operational reliability, effectively de-commoditising the high-end market.

### HOW AI SHATTERED THE TRADITIONAL CLOUD STACK

Early cloud architectures were designed for general-purpose workloads. Centered on commodity CPU servers, they heavily leveraged virtualisation and multi-tenancy to maximise utilisation by sharing resources among thousands of different customers. This “flexibility-first” model excelled at providing scalable, low-cost resources for common business applications, from websites to databases.

This model, however, is fundamentally misaligned with the demands of AI. In particular, training is not a flexible, multi-tenant workload; it is a single-task, brute-force supercomputing problem. When run on a legacy stack, the layers of virtualisation and network abstraction create a “performance tax” or overhead, siphoning critical compute cycles. This inefficiency didn't just require a new software layer, it forced a complete redesign of the data centre itself.

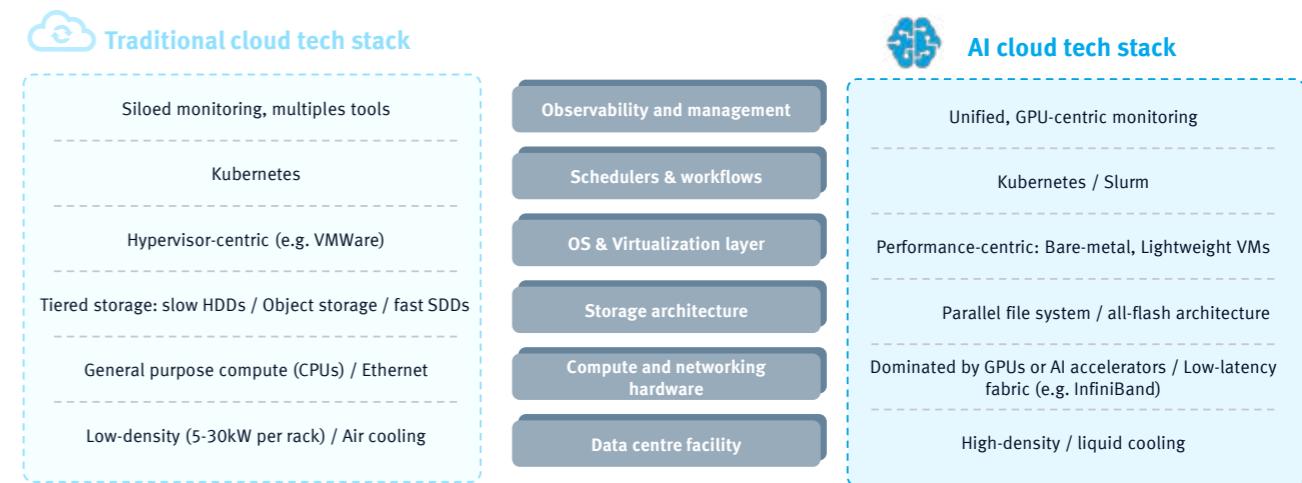
At the physical layer, an AI-optimised data centre is not a collection of servers; it is architected as a single, giant supercomputer. The design revolves around GPU racks with extreme power densities (>100 kW) requiring liquid cooling. Commodity Ethernet is replaced by low-latency, high-bandwidth fabrics like NVIDIA's InfiniBand or NVLink, essential to keep thousands of GPUs “talking” in sync without bottlenecks. The storage hierarchy is similarly overhauled, using parallel file systems (like those from DDN or Vast Data) to feed data at extreme throughput, avoiding the I/O bottlenecks that plague general-purpose clouds.

This shift also transforms the software stack. AI-oriented

infrastructure tends to strip away the virtualisation and multi-tenant abstractions that traditional clouds were built on. The main reasoning here is that hypervisors and host operating systems can siphon a portion of compute used in a workload, which in turn takes away some percentage of energy from the AI workload running. The legacy cloud's “one-size-fits-all” abstraction is replaced by specialised schedulers. The traditional HPC scheduler, Slurm, appears a favourite for training workloads, while Kubernetes has become a standard for inference workloads. Specialists like CoreWeave are deploying proprietary solutions, known as Slurm-on-Kubernetes (SUNK), as a core differentiator. This technology allows them to run both job types on a single GPU cluster, competing for resources in real-time. This hybrid capability is a major efficiency advantage, ensuring expensive GPUs are never idle.

This drive for efficiency also re-ignited the bare-metal vs. virtualisation debate. AI clouds initially championed bare-metal (like CoreWeave or Oracle) to strip away the hypervisor “tax”. The result is an “AI-native” stack stripped to essentials: the hardware and OS do just what's needed for large-scale model training, and nothing that isn't. Now, this debate is evolving. A new generation of lightweight VMs claims to offer better operational benefits like rapid provisioning and cleaner security isolation. Some of these lightweight VM solutions have demonstrated they can achieve bare-metal class performance. The choice is no longer a simple performance trade-off but a deeper, and still unsettled, architectural philosophy.

FIGURE 26: COMPARING TRADITIONAL AND AI CLOUD STACKS



Source: Stifel IRIS

The need for specific cloud services for AI is carving out a new, high-growth market within the broader IaaS landscape. According to IDC, spending on AI IaaS stood at \$19bn in 2024 and is projected to reach \$93bn in 2029 (+37% CAGR). This significant deployment of compute directly catalyses a parallel boom in AI infrastructure software, which IDC projects will grow from \$10bn in 2024 to \$19bn in 2029.

This growth is defined by a critical transition, moving from a market dominated by a handful of foundation model labs to one driven by widespread enterprise adoption. The initial 2023-24 wave was defined by training-centric capex, as major AI labs engaged in a training build up. This phase saw significant, multi-billion dollar contracts for GPU clusters, where training workloads constituted the visible majority of compute spend.

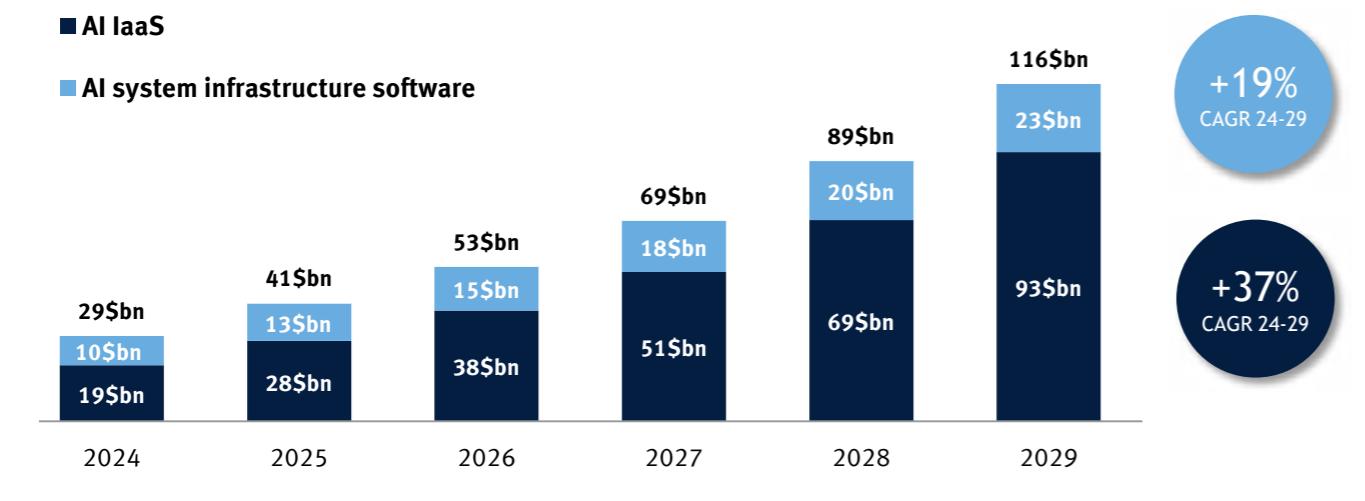
Today, the market is evolving into a more complex, multi-layered structure:

- **The foundation model training:** This segment remains the apex of compute demand, driven by the need to train next-generation models on cutting-edge architectures. This demand is lumpy, capital-intensive, and concentrated among a few AI giants.
- **Enterprise fine-tuning & private AI:** A new, high-growth wave of demand is emerging from enterprises. This involves fine-tuning smaller, specialised models or existing open-source models (like Llama) on proprietary datasets. This workload is more distributed and represents a new, secondary training market.
- **Mass-market inference:** This is arguably the largest and most sustainable driver. As millions of users adopt AI-powered applications, the resulting inference workloads become a significant, continuous, and distributed opex cost. While training is a discrete event, inference is a 24/7 utility, and its share of AI compute is already estimated to exceed 50%.

This evolution resolves the apparent paradox of hardware depreciation. While the useful life of a GPU for peak training is compressed (1-2 years), its economic life is far longer. This creates a tiered infrastructure market:

- **Tier 1 (Bleeding-edge):** Newest hardware (e.g., B200, GB200) serves the high-margin, peak-performance training market.

FIGURE 27: ANNUAL SPENDING ON AI IAAS AND INFRASTRUCTURE SOFTWARE, 2024-2029E



Source: IDC's Worldwide AI and Generative AI Spending Guide, V2 2025 (August 2025), Stifel IRIS



- **Tier 2 (Inference & fine-tuning):** Older hardware (e.g., H100, H200) is pushed down the stack to serve the mass-market inference and fine-tuning segments, which are more cost-sensitive. This allows neoclouds and hyperscalers to fully amortise assets over a three to five year lifespan, creating a long and profitable tail for existing infrastructure.

## THE BIRTH OF NEOCLOUD PROVIDERS

The upheaval in cloud architecture described previously paved the way for a new breed of provider purpose-built for AI: the GPU clouds, now widely known as neoclouds. Their emergence was not an accident, but a direct response to the fundamental architectural mismatch at the heart of the legacy cloud. Traditional hyperscalers were engineered for general-purpose, multi-tenant workloads, prioritising flexibility via virtualisation. Generative AI, however, is a brute-force, single-task supercomputing problem. As AI labs sought to build larger models, they found the hyperscalers' legacy stacks were technically unsuitable and operationally complex.

Neoclouds like CoreWeave, Nebius, and Crusoe rose to prominence by exploiting this gap. Their first advantage was agility. During the 2023-24 H100 supply crisis, these firms moved faster than the giants, providing impressive, AI-native clusters in record time. They offered AI labs what they desperately needed: faster deployment, flexible terms, and immediate access to the correct hardware stack that the hyperscalers were still struggling to deploy at scale.

The strategy has attracted immense investor interest and armed with capital, neoclouds have scaled aggressively. Many large neoclouds have landed marquee contracts with companies such as OpenAI, wins that were once thought unattainable for a startup competing with Microsoft and Amazon. By capitalising on unique opportunities and focusing on what AI customers care about most, these upstarts have cemented themselves as important competitors in the cloud ecosystem.

This agility was not a coincidence; it was a result of their unconventional DNA. Many neocloud leaders share an origin in cryptocurrency mining. This background provided the perfect, albeit accidental, training ground for the AI era. This "Crypto-to-AI Pivot" has taken two primary forms:

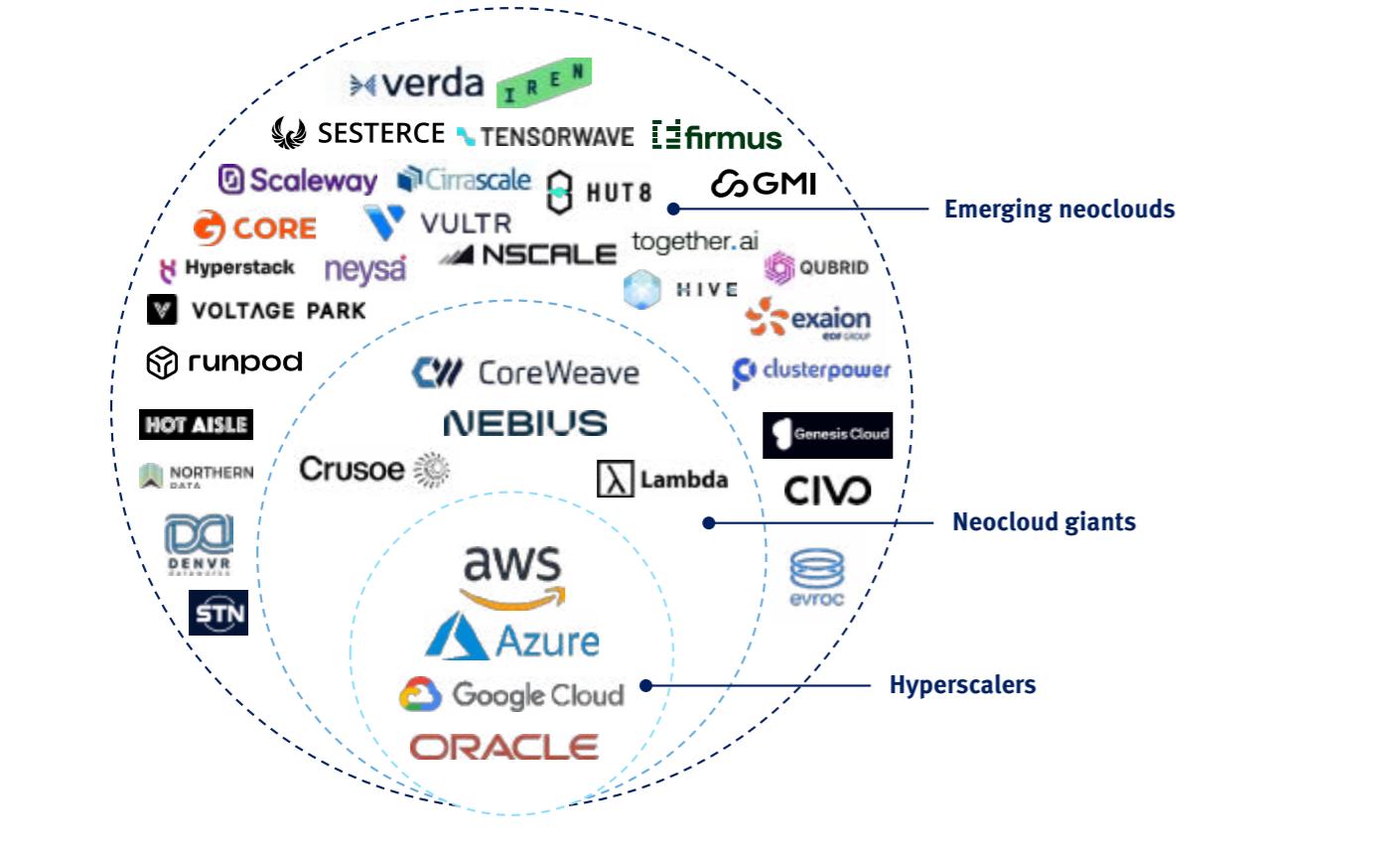
- **The operational pivot:** Companies like CoreWeave (an Ethereum miner since 2017) and Crusoe (mining

with flared gas) repurposed their deep, low-level hardware expertise. They were already specialists in running high-density, low-margin GPU infrastructure 24/7, giving them a significant operational head-start in both supply and engineering talent.

- **The infrastructure pivot:** A second wave, which continues to shape the market, includes publicly traded miners like IREN, Hut 8, and TeraWulf. Their primary asset was not GPUs, but gigawatt-scale, high-density data centre sites with access to low-cost power. These firms are now vertically integrating into wholesale bare-metal providers or, more commonly, leasing their powered shells to other operators. This has also enabled a new hybrid model where specialised operators like Fluidstack partner with these miners to transform raw power and land into elite, managed AI clusters.

To compete and scale in the 2023-2024 market, neoclouds initially paired their technical advantages with aggressive pricing. On price, they often undercut the major cloud providers by a wide margin. This led to a widespread debate among investors about the long-term viability of the model and the risk of GPU-as-a-Service becoming a fully commoditised market, with a race to the bottom on margins. This initial price gap was possible due to leaner cost structures; unlike hyperscalers, neoclouds did not have to support a sprawling portfolio of legacy services. A hyperscaler's cloud offers millions of individual billing items, while a neocloud's focus on a few core hardware configurations allowed those savings to be passed to customers. But as the market matures, price is only half the story. Today, as the GPU shortage eases and hyperscalers launch their own dedicated HPC/AI offerings (often copying the neocloud playbook), the competitive battleground has shifted from hardware access to operational excellence.

FIGURE 28: MAPPING NEOCLOUDS



## BEYOND PRICE PER HOUR

The idea that GPU compute is a commodity, easily compared by the price per GPU-hour, is a fundamental misunderstanding of the market, in our view, and stems from the 2023-24 dynamics when neoclouds used aggressive pricing to gain market share, often undercutting hyperscalers, which led many to predict a margin-destroying race to the bottom.

By late 2025, this narrative is obsolete. The market has bifurcated: the best providers now command a pricing premium for their services, with prices closer to those of the hyperscalers. They no longer compete on raw price-per-hour but on Total Cost of Ownership (TCO). For a multi-week training job, uptime is an irrelevant metric; the only metric that matters is goodput (the volume of successful computation). A job that fails at 90% completion does not cost 10% more; its TCO approaches infinity.

This de-commoditisation is achieved through deep engineering moats in software and operations. Two examples illustrate this operational gap. First is superior software orchestration. The primary technical challenge is integrating HPC schedulers (like Slurm), built for significant jobs, with modern platforms (like Kubernetes). This “Slurm-on-Kubernetes” (SonK) battle is a primary software moat. Best providers leverage mature, battle-tested proprietary solutions that are proven reliable, such as CoreWeave’s SUNK or Nebius’s Soperator. In contrast, second-tier providers often deploy buggy open-source alternatives, which leads to failures and proves this orchestration layer is a critical differentiator.

A second moat is proactive reliability. Elite providers maximise goodput by moving beyond reactive support tickets. They implement automated Passive and Active Health Checks (diagnostics that constantly monitor

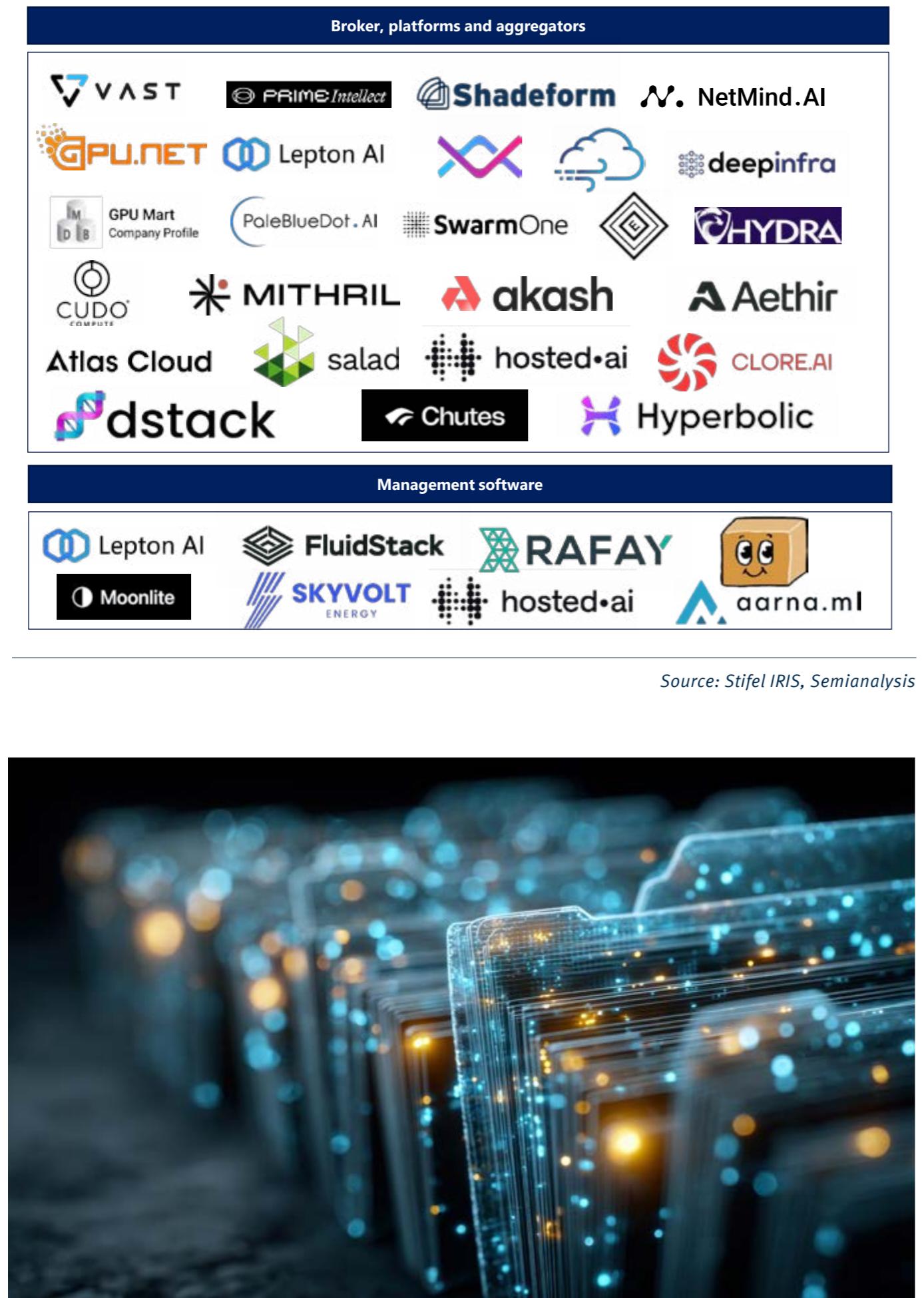
hardware). This system finds and automatically removes failing hardware before it can crash a customer’s critical job, directly reducing TCO.

This software and reliability gap solidifies a complex competition landscape. Hyperscalers, despite their scale, have proven unable or unwilling to fill this operational void. Their own high-performance offerings, such as AWS HyperPod or Azure CycleCloud, are notoriously suffering from complex usability issues, and offload the most critical operational burdens directly onto the customer. This creates a permanent last mile service gap that specialised neoclouds are purpose-built to exploit.

NVIDIA’s own strategic missteps further validate this market structure. After an acquisition spree (including run:ai and Lepton) to build a competing DGX Cloud, the effort failed to produce a coherent or competitive product. NVIDIA was forced to pivot, effectively acknowledging that it could not directly compete with its largest and most important customers: the neoclouds themselves.

A genuine commodity segment for cloud compute exists, primarily served by marketplaces, aggregators, and brokers. These entities serve as intermediaries, pooling GPU supply by connecting buyers with multiple suppliers, often without operating their own GPU infrastructure. While this model enables competitive pricing and broad accessibility, it typically emphasises spot-market offerings and caters to non-critical workloads or developer use cases. As a result, these platforms may not always meet the full range of enterprise-grade requirements around security certifications or end-to-end operational control, which remain the core focus of premium, enterprise-oriented neocloud providers.

FIGURE 29: MAPPING GPU PLATFORMS



## 2.3 THE PLAYBOOK IS NOT TO COMPETE, BUT TO BYPASS

We note that the hyperscaler model, while dominant, is not the only path to success in the data centre industry. Their significant scale and standardised approach inevitably create structural gaps and specialised market needs. This section deconstructs the successful strategies that enable this co-existence, proving that value is captured not only by competing with the giants, but by complementing or bypassing them.

### UNDERSTANDING THE HYPERSCALER FORTRESS

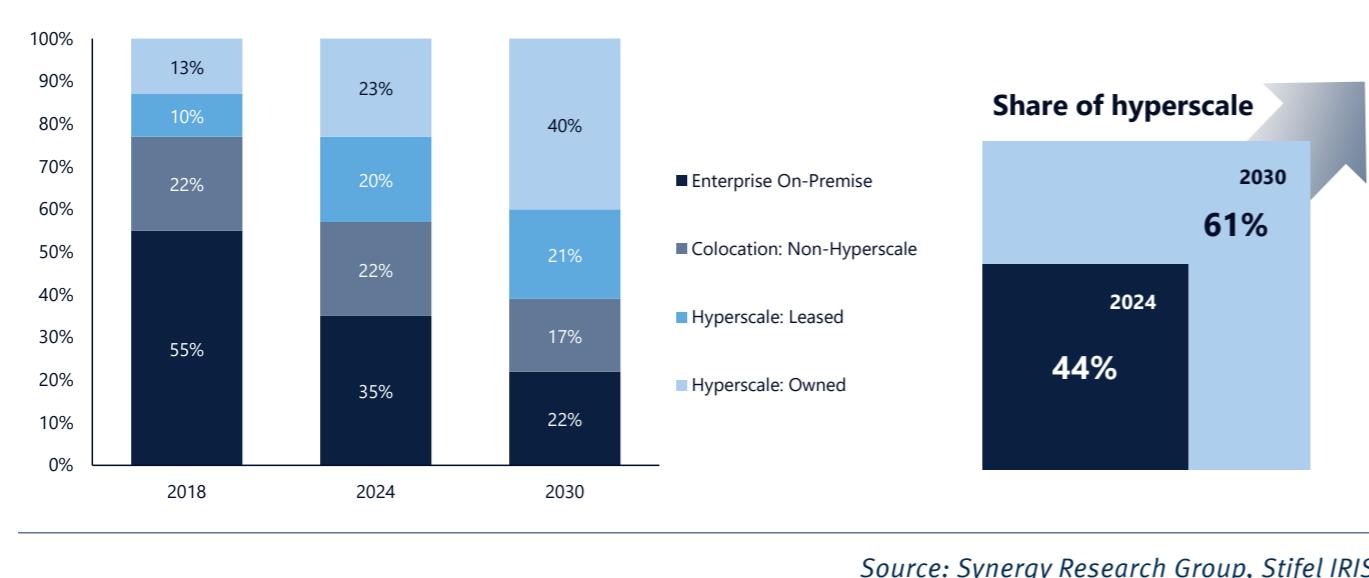
The dominant position of hyperscalers is no accident; it is the structural outcome of a market defined by sizeable, multi-layered economies of scale. This “rich get richer” dynamic suggests the public cloud functions as a quasi-natural oligopoly, as the industry is defined by economies of scale and scope at every layer.

First is the physical infrastructure layer. The business is defined by a virtuous cycle of scale: building significant data centres is far cheaper per megawatt, but this requires prohibitive, front-loaded capital expenditures that lock out new entrants. This scale unlocks compounding cost advantages in procurement, operational automation, and optimised PUE, and is extended globally via proprietary fibre backbones and subsea cables. The most pivotal development in the data centre market over the past decade has thus been the rise of these hyperscale facilities. Synergy Research Group reports that the hyperscale share of global data centre capacity surged from 23% in 2018 to 44% in 2024, and is projected to reach 61% by 2030. These facilities, which typically exceed 40MW of critical power and require more than \$1bn in capex (inclusive of IT equipment). They can form expansive campuses, with future “AI megacentres” potentially reaching multiple GW. This global footprint generates additional network

effects: architectures can be replicated for efficiency, and workloads can be optimised across territories.

The second level of this fortress is deep vertical and horizontal integration. The significant fixed costs of developing and maintaining a vast, proprietary portfolio of IaaS and PaaS services create a high initial software barrier. Horizontally, this allows the same foundational infrastructure to be leveraged across a wide array of distinct services. More importantly, this integration extends vertically. Less visible has been the internalisation of hardware design, as hyperscalers develop their own custom silicon, optimising performance and cost in a way competitors reliant on merchant silicon cannot match. This full-stack strategy now targets every layer: general compute chips like AWS’s Graviton series, Google’s Axion, and Microsoft’s Cobalt are engineered to optimise core IaaS costs. Simultaneously, AI acceleration chips like Google’s TPUs, AWS’s Trainium, and Microsoft’s Maia are designed to slash the cost of running high-margin AI platforms. This custom silicon strategy introduces a distinct capital barrier: with fixed costs for designing a leading-edge chip exploding, now often exceeding hundreds of millions of dollars, this level of deep vertical integration is financially prohibitive for nearly all other market participants.

FIGURE 30: HYPERSCALE BUILDOUTS REDEFINED THE 2010S, AND STILL SHAPE THE 2020S - DATA CENTRE CAPACITY TRENDS: SHARE OF CRITICAL IT LOAD (WORLDWIDE, MW)



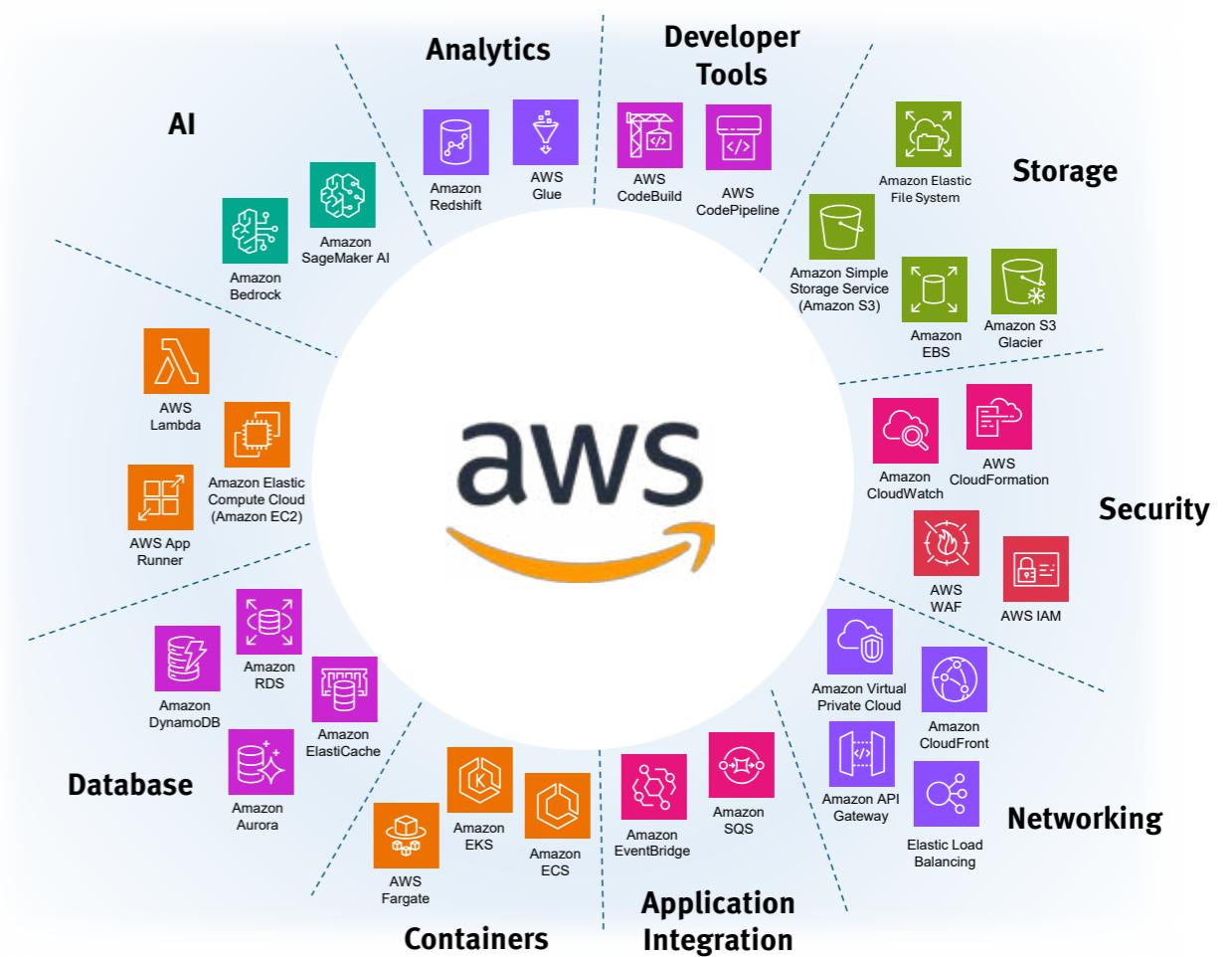
The sheer breadth of the hyperscaler service catalogue functions as a powerful competitive moat. Customers often select a provider based on its ability to cover all current and future needs, creating a one-stop shop dynamic. This advantage is then monetised through a vast portfolio of high-margin, proprietary, and deeply integrated PaaS services. As customers build complex applications intrinsically tied to these unique APIs, migration to a competitor becomes prohibitively expensive and operationally complex.

This technical lock-in is reinforced by deliberate commercial strategies that have drawn significant regulatory scrutiny. Customers often find it difficult to anticipate future costs, a situation compounded by complex offerings and a lack of pricing transparency.

This is amplified by two key practices:

- **Cloud credit programmes:** These are not just free trials but strategic investments targeting high-potential users, like startups, with amounts (e.g., hundreds of thousands of dollars) that smaller providers cannot profitably replicate. This effectively locks in the next generation of high-value customers.
- **Egress fees:** Hyperscalers charge significant fees for transferring data out of their environment, creating the financial lock. This is a major commercial barrier to switching providers or adopting a multi-cloud strategy. The French Autorité de la concurrence noted these fees appear “disconnected from the costs directly borne by providers”, with one hyperscaler admitting the price difference for ingress (free) versus egress (paid) is “essentially commercial”.

FIGURE 31: BEYOND INFRASTRUCTURE: THE STRATEGIC DEPTH OF THE AWS PLATFORM

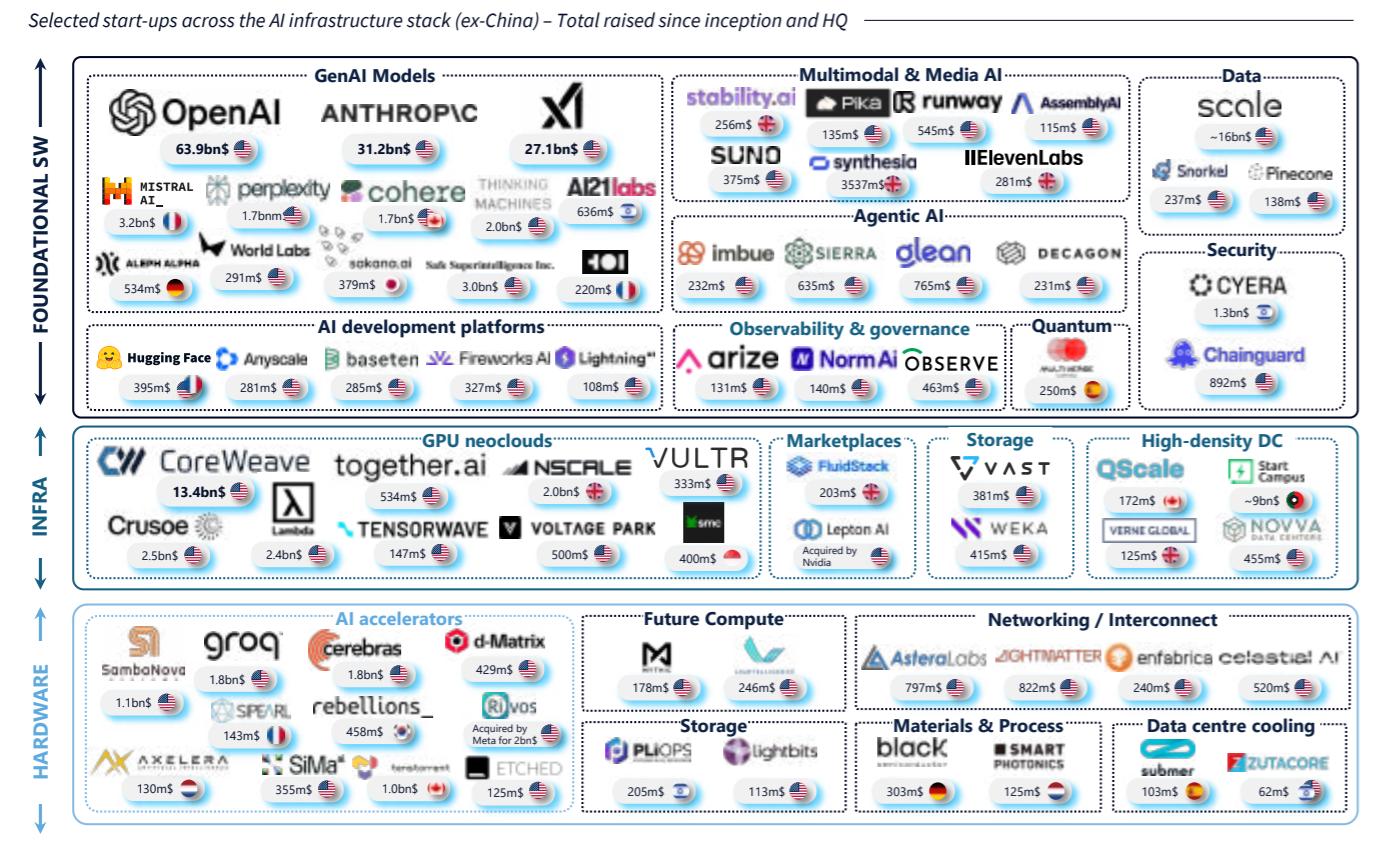


Source: AWS, Stifel IRIS

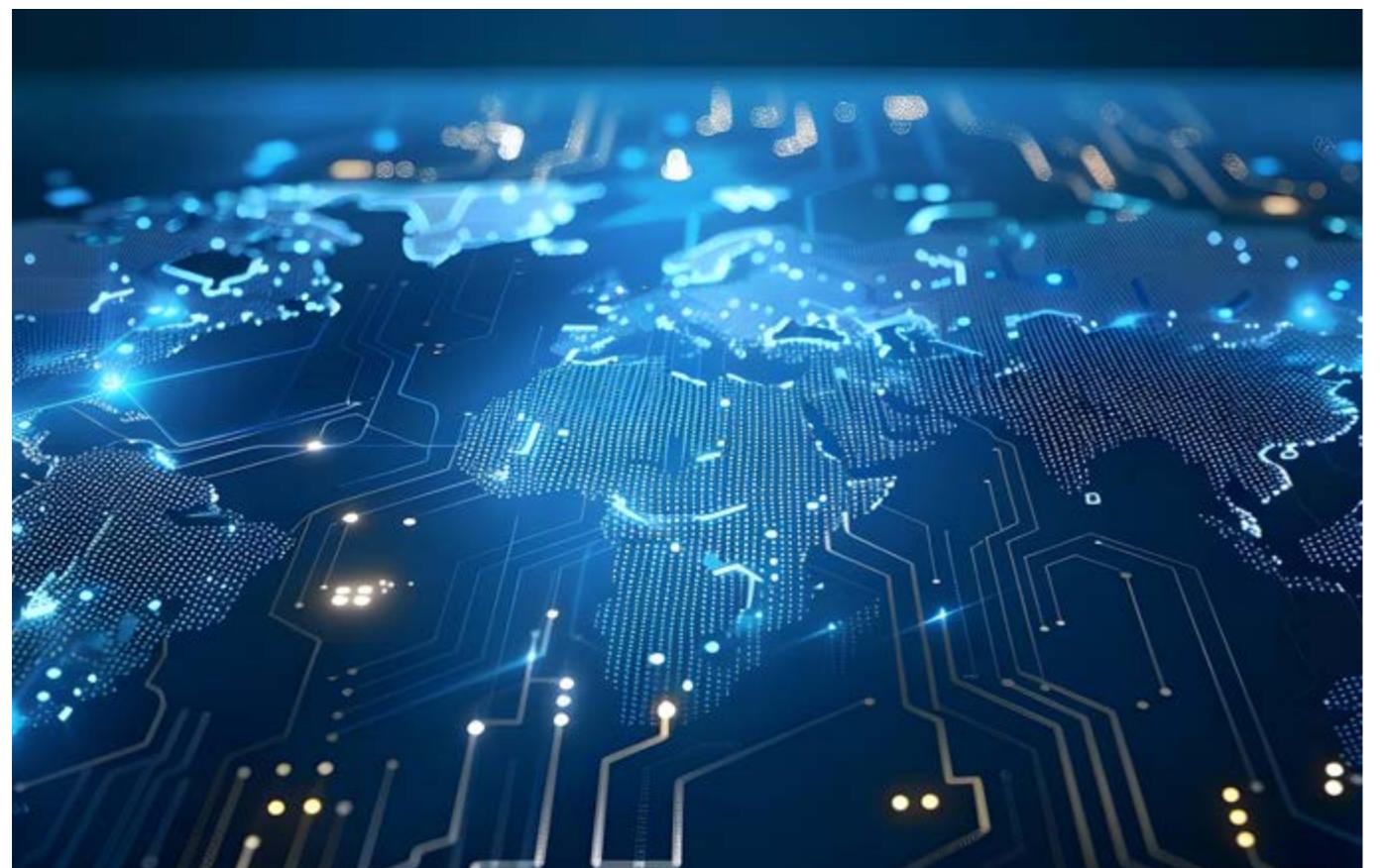
Yet, this narrative of an unassailable oligopoly is not the complete picture. The market remains highly dynamic, and we identified dozens of startups successfully raising over \$100m to tackle specific bottlenecks. The sheer velocity and scale of the AI-driven capital cycle

has created structural openings across the entire value chain. This new wave of innovation is visible at every layer of the stack, from hardware and infrastructure to the foundational software layer itself.

FIGURE 32: MAPPING THE AI INFRA & FOUNDATIONAL START-UP ECOSYSTEM



Source: Institute cyclos-HTP and Stifel IRIS



## THE BYPASS PLAYBOOK

The dominant hyperscaler model is not without structural vulnerabilities. We believe the playbook for emerging champions is not to compete head-on with a trillion-dollar incumbent; it is to bypass. This section deconstructs five proven “bypass strategies” employed by specialised providers, both established and new. These playbooks demonstrate how differentiation and focus can capture value.”

### Bypass Strategy #1: Low-cost and simplicity play

This bypass strategy counters the hyperscalers’ integrated complexity with a focus on radical simplicity and superior price-performance. It is built on the thesis that a significant market segment values predictable billing and legible configurations over access to a sprawling, integrated service catalogue.

This is most evident in the mature compute market. European operators like Hetzner and OVHcloud, along with DigitalOcean, have built highly successful businesses by not emulating the hyperscaler playbook. They target developers and SMBs (customer segments often over-served by the hyperscalers’ sprawling service catalogues) with a “no-frills” offer: simple, legible VM and bare-metal configurations with predictable, transparent billing. For this audience, predictable opex and ease-of-use are more valuable than access to 200 integrated services.

This strategy of specialised price-performance was recently replicated in the AI market. The neoclouds’ initial success was a classic bypass manoeuvre. They identified a new, high-value workload (GPU compute) where hyperscalers were initially uncompetitive, offering faster access to specialised hardware at aggressive, transparent price points, thereby capturing a significant share of the nascent training market.

A third vector attacks the hyperscalers’ most extractive commercial practice: data egress fees. Hyperscaler billing is notoriously opaque, with high, unpredictable

egress charges acting as a powerful driver of vendor lock-in. While hyperscalers like AWS and Google have recently waived some fees for full exits amid competition and regulations, this practice is being unbundled further. Storage specialists like Cloudflare (R2), Wasabi, and Backblaze (B2) offer S3-compatible object storage with zero egress fees, a move that dramatically reduces TCO for data-intensive applications.

### Bypass Strategy #2: Sovereignty

This strategy transforms a geopolitical liability into a core commercial product. In an era of rising data nationalism, US-based hyperscalers are a systemic risk for foreign governments and regulated industries. Legislation like the US CLOUD Act grants US authorities jurisdiction over data regardless of its storage location, creating a fundamental conflict with regional frameworks like the EU’s GDPR and Data Act.

Historically, however, sovereignty alone was a weak commercial argument. The trade-off was un compelling: a theoretical legal risk (CLOUD Act) versus tangible operational risks. Local champions often lacked hyperscaler resilience, a fact highlighted by high-profile failures like the 2021 OVHcloud fire and the 2025 NIRS data centre fire in South Korea, which shut down 647 government services. This dynamic made operational failure a more immediate threat than foreign jurisdiction. China remains the exception, where state protectionism created a domestic market dominated by local giants (Alibaba, Tencent, Huawei).

In open markets like Europe, hyperscalers neutralised the regulatory threat. Top-down standardisation attempts like Gaia-X largely failed, bogged down in bureaucracy and joined by the hyperscalers themselves. In contrast, potent national certifications like France’s SecNumCloud forced hyperscaler adaptation via local JVs (e.g., Bleu, S3NS), which neutralised the sovereignty argument by ceding operational control.

We note this strategy is regaining critical relevance due to two factors. First, escalating geopolitical friction elevates the CLOUD Act from a legal theory to a non-negotiable risk. Second, the battleground has shifted from data sovereignty (protecting storage) to sovereign AI (controlling model development). Dependency on foreign foundational models is now a strategic threat,

creating urgent demand for trusted infrastructure to train national AI models.

The ultimate challenge for this playbook remains: moving sovereign cloud from a niche product for government and defence to a commercially competitive offering for large enterprises, who are still reluctant to accept a performance trade-off.

FIGURE 33: MAJOR SOVEREIGN CLOUD INITIATIVES WORLDWIDE



### Bypass Strategy #3: Winning the edge market

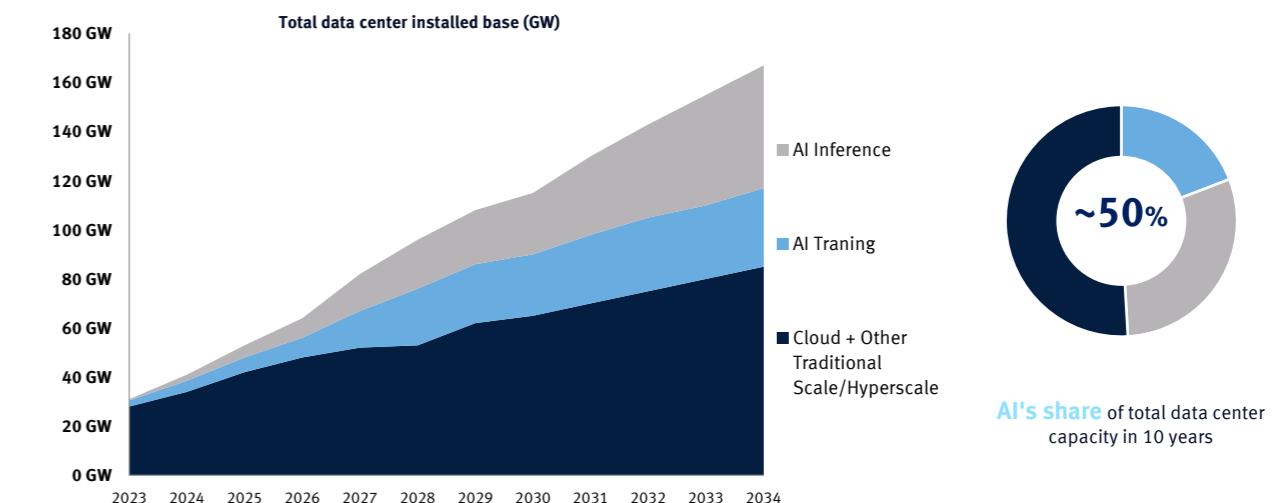
This bypass strategy operates on a different axis, trading the hyperscalers' significant scale for two specific advantages: low latency and data confidentiality. It exploits the fact that the giants' centralised, training-optimised architecture is structurally ill-suited for critical, real-time AI workloads.

Crucially, this must be distinguished from the mass-market inference. While latency requirements will force hyperscalers to build more localised large-scale data centres, we believe the bulk of AI inference will remain a centralised, high-volume business. The “true edge” bypass is different: it targets workloads that must be processed significantly more locally, driven either by sovereignty mandates or mission-critical, ultra-low latency needs.

This “true edge” market has seen slower-than-hyped growth, as previous use cases in 5G and Industrial IoT did not materialise as quickly as anticipated. The AI era, however, introduces a new, more potent driver: confidentiality. For sensitive applications in government, healthcare, defence, or critical industries, running inference on a US-based hyperscaler is often not an option, regardless of cost.

This creates a high-margin, defensible niche. While the total addressable market is an order of magnitude smaller than the mass-market inference opportunity, its strategic value is high. The playbook is not to win the volume race but to secure a premium segment defined by high barriers to entry (regulatory trust, local relationships) and sticky, multi-year contracts with government and critical industry clients. This focus allows for premium pricing exceeding commoditised cloud offerings.

FIGURE 34: PROJECTED SURGE IN AI-DRIVEN DATA CENTRE CAPACITY



#### Bypass Strategy #4: The power-access focus

This strategy reframes energy from a simple opex item to the single most critical strategic asset. As the AI buildout collides with a structural power bottleneck, securing power access has become a primary competitive advantage, often superseding the cost of that power.

While electricity represents a major part of data centre opex, it is a secondary driver of TCO for premium AI services compared to hardware amortisation. The true bottleneck is access. Securing a large, stable power envelope is what enables economies of scale (e.g., >1GW campuses in the US) and unlocks strategic locations in constrained markets.

This new reality is forcing a strategic pivot. The hyperscalers' long-standing clean energy narrative is now under pressure, competing with the sheer, immediate need for power to fuel the AI race to dominance. While sustainability remains a key regulatory driver, especially in Europe, the bypass strategy itself is less about being "green" and more about being available.

This creates a playbook for specialised operators who bypass the congested power grids of primary hubs. This strategy is defined by two main approaches:

- Monetising stranded energy: This model targets non-latency-sensitive compute by deploying modular data centres directly at source of wasted or stranded energy. Crusoe Energy, for example, captures flare gas from oil fields, turning an environmental waste product into low-cost compute capacity, completely bypassing the grid. A second wave includes publicly traded miners (e.g., IREN, Hut 8, TeraWulf) leveraging their existing, gigawatt-scale data centre sites, which were already built on low-cost, stranded power, to pivot into AI cloud and colocation providers
- Relocating to power-rich regions: Instead of competing for scarce power in primary hubs (like Frankfurt or London), this model relocates non-latency-sensitive workloads to regions with abundant, stable, and often renewable energy. The Nordics are the prime example. Companies like EcoDataCenter in Sweden leverage the

region's significant hydropower resources, offering sustainable, large-scale capacity while bypassing the grid congestion and regulatory hurdles that plague the established FLAP-D markets.

#### Bypass Strategy #5: The "Picks & Shovels" approach

We note the picks & shovels (P&S) approach has been a dominant investor narrative, arguing that one should supply AI expansion rather than pick the uncertain beneficiary. We agree but find this view incomplete. This P&S market also benefits from being highly fragmented in some areas, creating distinct opportunities for specialised contenders. We see two major pockets of opportunity.

The first area concerns electrical and cooling equipment. The AI-driven shift to extreme power densities and liquid cooling is a major disruptive event. This market is not fully controlled by the large industrial incumbents like Schneider Electric, Vertiv, and Eaton. These giants, while dominant in traditional power and air cooling, are primarily geared for incremental optimisation, not disruptive R&D. This creates a classic innovator's dilemma, opening a significant gap for agile, specialised startups, particularly in high-growth niches like advanced liquid cooling. This dynamic fuels a robust M&A market: incumbents are forced to acquire this external innovation to keep pace, providing a clear exit path for successful, focused challengers.

The second area involves the data centre developers and operators. This is an equally critical P&S approach using the developer-operator model, which provides value to the hyperscalers themselves. Firms like STACK Infrastructure, EdgeConneX, and Compass Datacenters are not just builders; they are strategic enablers. Their value proposition is clear: they arbitrage the hyperscalers' structural limitations by offering (1) Speed-to-market, accelerating deployment in a time-critical AI race; (2) Capital allocation, allowing hyperscalers to deploy capex on high-ROI GPUs instead of concrete; and (3) Local expertise, navigating the complex, non-scalable challenges of regional real estate, permitting, and grid interconnection.

## A EUROPEAN DEEP DIVE

### SECTION 3

This section deep-dives into the European data centre market, arguing its path is structurally different from the US. Hampered by high energy prices, a lack of native tech giants, an AI compute gap, and strict regulatory mandates, Europe is ill-suited for the US-led, training scale race. This constrained environment creates unique opportunities, but only for companies that master a specific bypass playbook. We map the four specialised verticals that we believe are best positioned to benefit.



# 3.1 UNDERSTANDING THE EUROPEAN DATA CENTRE LANDSCAPE

Europe's data centre market, traditionally centred on the FLAP-D core, is being forcibly decentralised by new AI demands. Grid saturation is forcing significant-scale training workloads to power-rich peripheries (Spain, Nordics), leaving core hubs to specialise in latency. Stringent EU energy regulations and an AI compute gap are creating a unique, challenging environment, pivoting the core opportunity from training to localised, sovereign inference.

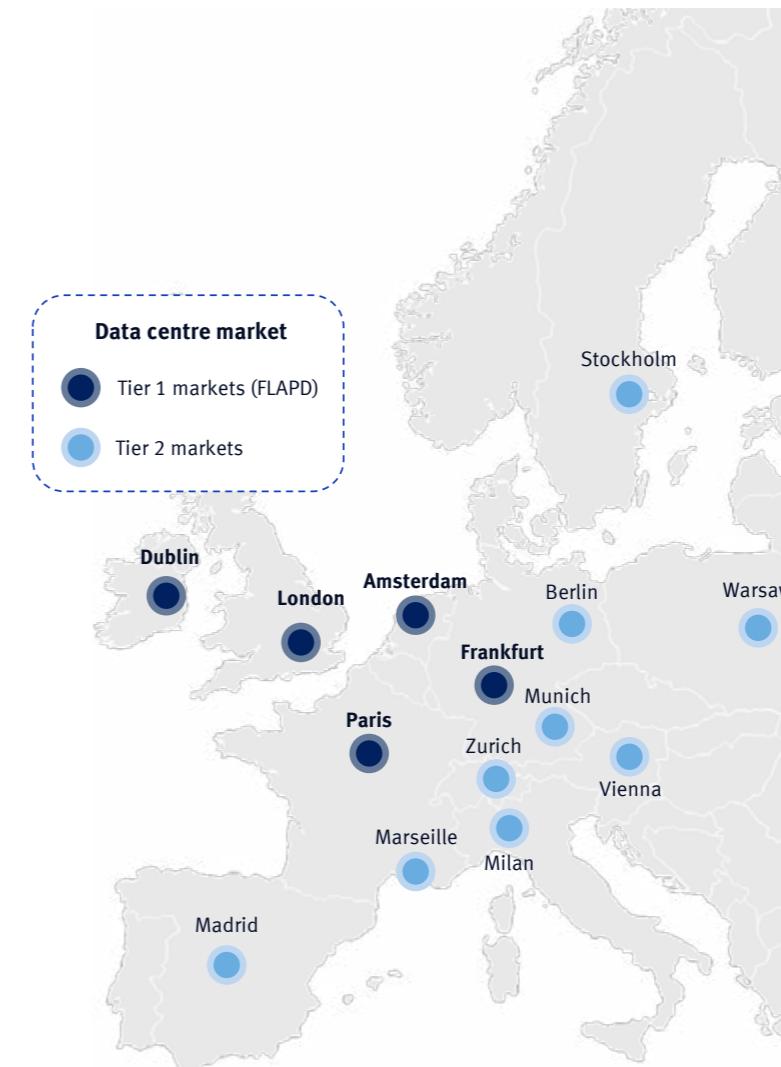
## FROM FLAP-D CORE TO POWER-RICH PERIPHERY?

Europe's data centre landscape remains architecturally centred on the "FLAP-D" core cities: Frankfurt, London, Amsterdam, Paris, and Dublin. The countries of these five Tier-1 markets collectively host approximately 72% of the continent's available IT power supply, according to Pb7 research. The dominance of these markets stems from decades of infrastructure entrenchment: dense fibre optic networks, superior interconnection ecosystems (AMS-IX in Amsterdam one of Europe's largest internet exchange), and critical submarine cable landings (MAREA, Dunant).

This concentration has evolved through key historical phases, each reflecting shifting technological and economic priorities. From the 1990s' enterprise-owned, on-premises facilities emphasising sovereignty and reliability, the market transitioned in the early 2000s to colocation models that pooled resources for efficiency and scalability, catalysing FLAPD's rise as connectivity nerve centres. The 2010s saw hyperscale cloud providers drive further expansion, introducing energy-optimised designs and the start of a modest diversification toward power-rich peripheries like the Nordics for its abundance of hydropower and natural cooling advantages.



FIGURE 35: MAPPING EUROPE'S DATA CENTRE CORE AND EMERGING MARKETS



Source: Stifel IRIS

The post-2022 AI boom, however, has introduced acute tensions, amplifying a debate over decentralisation toward Tier-2 metros and southern/northern peripheries. The primary constraint impacting European data centre markets is no longer local land or high costs, but structural power access. In legacy hubs, multi-year grid connection delays (in some markets extending to 10 years or more) effectively sterilise immediate expansion capacity. This constraint is most acute in Amsterdam and Dublin.

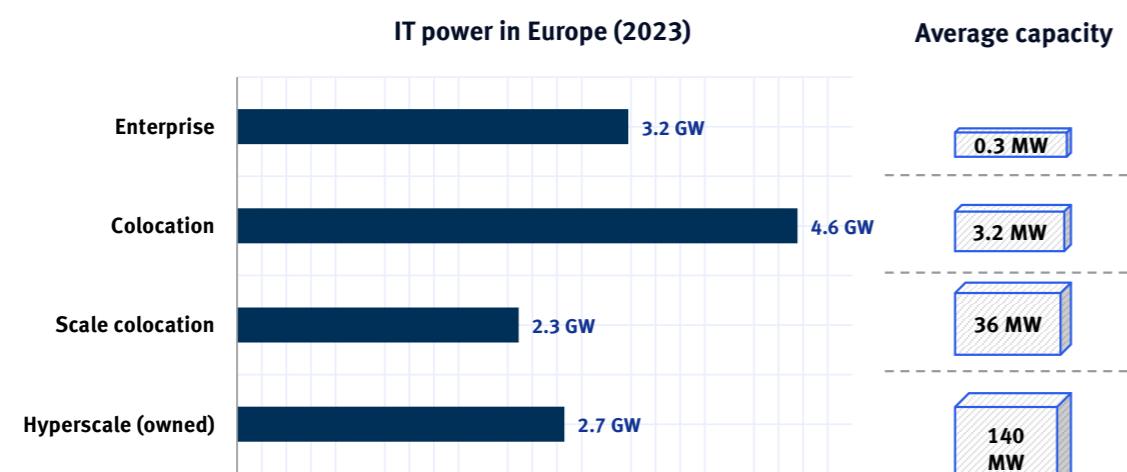
This energy crisis, reinforced by strict regulatory mandates, is the key driver of geographic shift. AI training workloads prioritise abundant renewables and regulatory

compliance over proximity. This has propelled Southern Europe as the continent's fastest-growing region, notably Spain and Italy. This acceleration is evident in gigawatt-scale projects, such as Spain's Extremadura attracting 1 GW campuses like Edged Energy's, leveraging solar/wind surpluses and subsea cable access. The Nordics benefit from similar green credentials, but have recently faced tempered expectations, with hydropower volatility and grid limitations as potential caps. We believe this current dynamic is compelling the industry to adopt a two-tier architecture, where specialisation in low-latency environments remains essential for resilient metro hubs, while significant scale is built in the power-rich periphery.

Europe's data centre capacity reached 12.7 GW of available IT power supply in 2023, housed across roughly 10,500 facilities. However, this total figure obscures a critical two-tier market structure. The landscape is dominated by a mass of over 9,000 smaller enterprise data centres, while a power elite of just ~1,500 colocation and hyperscale facilities control ~75% of IT power and capture nearly all market growth. Competitively, the European market remains notably fragmented. The European Data Centre Association (EUDCA) estimates that its member operators, including global leaders such as Equinix and Digital Realty, collectively hold no more than 45% market share within the European colocation sector, spanning retail, wholesale, and scale segments. In terms of market segmentation, colocation facilities dominate, accounting for over 50% of Europe's IT power. Hyperscalers have significantly expanded their footprint to represent nearly 40% of total capacity, executing a

dual strategy of leasing "Scale Colocation" and building their own "Hyperscale-Owned" facilities. This owned footprint, comprising just 19 facilities but averaging a significant 140 MW each, reveals a clear geographic strategy. Currently, this self-built capacity is concentrated in just two regions: Ireland (43%) and the Nordics (32%). However, the next wave of hyperscale-owned investment is already targeting new zones, with major projects announced for the United Kingdom, Spain, and Italy, signalling a strategic pivot to new power-rich areas. While Europe's share of global capacity (around 20%) should temporarily decline due to the intense, US-centric AI investment supercycle, the continent's own capacity is still forecast by several projections to at least double by 2030, an expansion projected to necessitate investments exceeding €100bn in European data centre infrastructure, exclusive of IT equipment.

FIGURE 36: DECONSTRUCTING EUROPE'S 12.7 GW: TOTAL CAPACITY AND AVERAGE SCALE PER SEGMENT



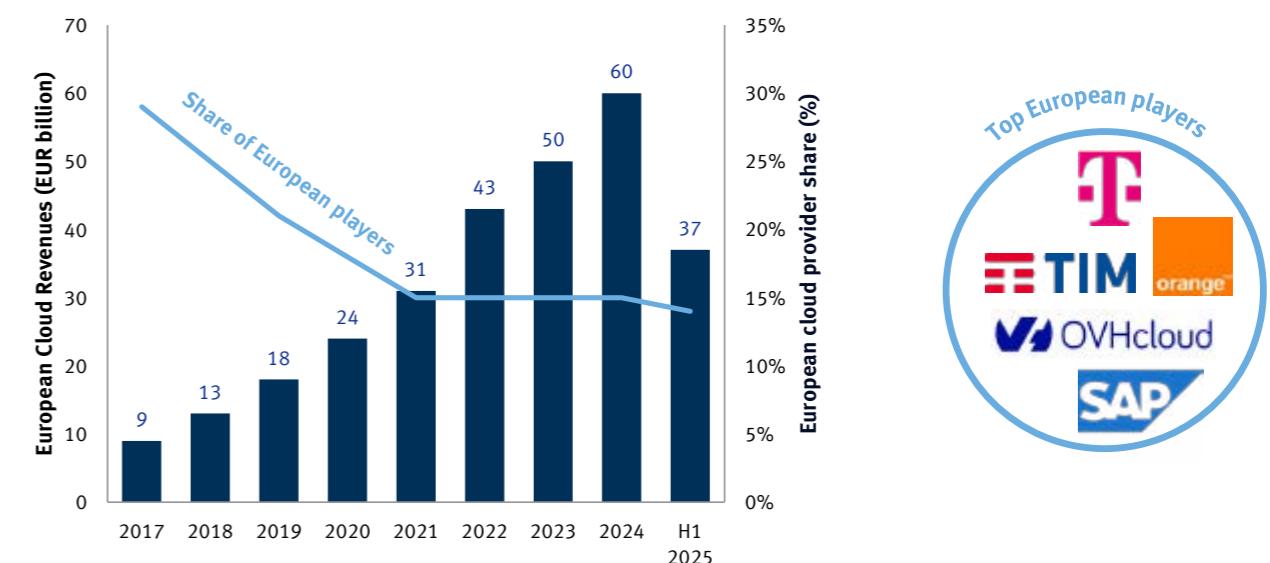
Source: Pb7 Research, Stifel IRIS

## EUROPE'S DISTINCTIVE CLOUD AND REGULATORY LANDSCAPE

The European cloud infrastructure services market (IaaS, PaaS, Hosted Private Cloud) reached €61bn in 2024, having expanded six-fold since 2017. This growth, however, has been overwhelmingly captured by non-European entities. According to Synergy Research Group, the collective market share of European providers collapsed from 29% in 2017 to just 15% by 2022. The main beneficiaries (Amazon, Microsoft, and Google) now command a combined 70% of the regional market. However, this narrative of structural erosion masks two critical nuances. First, this occurred even as absolute revenues of local providers tripled over the same period, highlighting a market expanding so rapidly that it could accommodate growth even for share-losers. Second, this market share collapse has since plateaued,

holding steady at ~15% since 2022. This stabilisation has resulted in a highly fragmented landscape for local companies: the remaining 15% is not a bloc, but a long tail of national champions operating in the shadow of the hyperscalers. This landscape is fragmented, defined by incumbent telecommunications giants (including Deutsche Telekom's T-Systems, Orange Business Services, and Telecom Italia) leveraging their established enterprise customer base, and cloud-native pure play providers (like OVHcloud, Hetzner, and IONOS) competing aggressively on price-performance or developer focus. Pan-European software leaders, notably SAP and Dassault Systèmes (Outscale), also hold a distinct position.

FIGURE 37: EUROPEAN CLOUD MARKET (IAAS, PAAS, HOSTED PRIVATE CLOUD)



Source: Synergy Research Group, Stifel IRIS

The European data centre market is characterised by a distinctive confluence of energy imperatives, shaped by the region's ambitious climate goals and digital sovereignty agenda. The primary physical constraint stems from the high penetration of intermittent renewables (accounting for over 44% of EU electricity mix in 2024 according to Ember) which exacerbates grid instability and price volatility. This poses a critical challenge for data centres to be uninterrupted 24/7. Compounding this is an increasingly stringent regulatory landscape, where the EU is codifying efficiency standards to align with net-zero targets. This strategy is being deployed in two distinct layers:

- The EU framework:** The cornerstone is the revised Energy Efficiency Directive (EED), which entered into force in 2023 and mandates enhanced reporting from 2025 onward. Specifically, data centres exceeding 500 kW in installed IT power must annually disclose standardised key performance indicators (KPIs) on energy consumption, water usage effectiveness (WUE), renewable energy factor (REF), and waste heat reutilisation. While not yet a punitive tool, it creates the transparent data foundation for national-level enforcement.
- The national sticks:** The second layer is the national implementation, which transforms this transparency into hard law. The most stringent example is

Germany's Energieeffizienzgesetz (EnEfG). This law moves from reporting to mandating, imposing new facilities (operational after July 1, 2026) must achieve a PUE of 1.3 within two years, while existing facilities must reach 1.5 by July 1, 2027 and 1.3 by July 1, 2030. Comparable frameworks are emerging elsewhere, such as France's RE2020 building codes mandating carbon lifecycle assessments and Ireland's grid connection moratoriums tied to sustainability criteria, underscoring a pan-EU push toward accountability.

EU data centre regulations are poised to disproportionately favour large-scale operators over smaller ones, due to the difficulty of achieving energy efficiency targets and more generally compliance costs. This framework should compel operators to innovate aggressively in energy efficiency (liquid cooling, heat reuse) and explore alternative power solutions, such as nuclear SMRs, to ensure stability. However, compliance costs could deter investments, exacerbating Europe's lag behind the US and pushing AI firms towards more flexible markets like the UK. This cloud creates a core contradiction with the ambition from the European Commission to triple EU data centre capacity by 2035 as part of its AI Continent Action Plan announced in April 2025, framing these facilities as essential infrastructure for AI competitiveness and digital sovereignty.



## EUROPE'S POSITION IN THE AI COMPUTE RACE

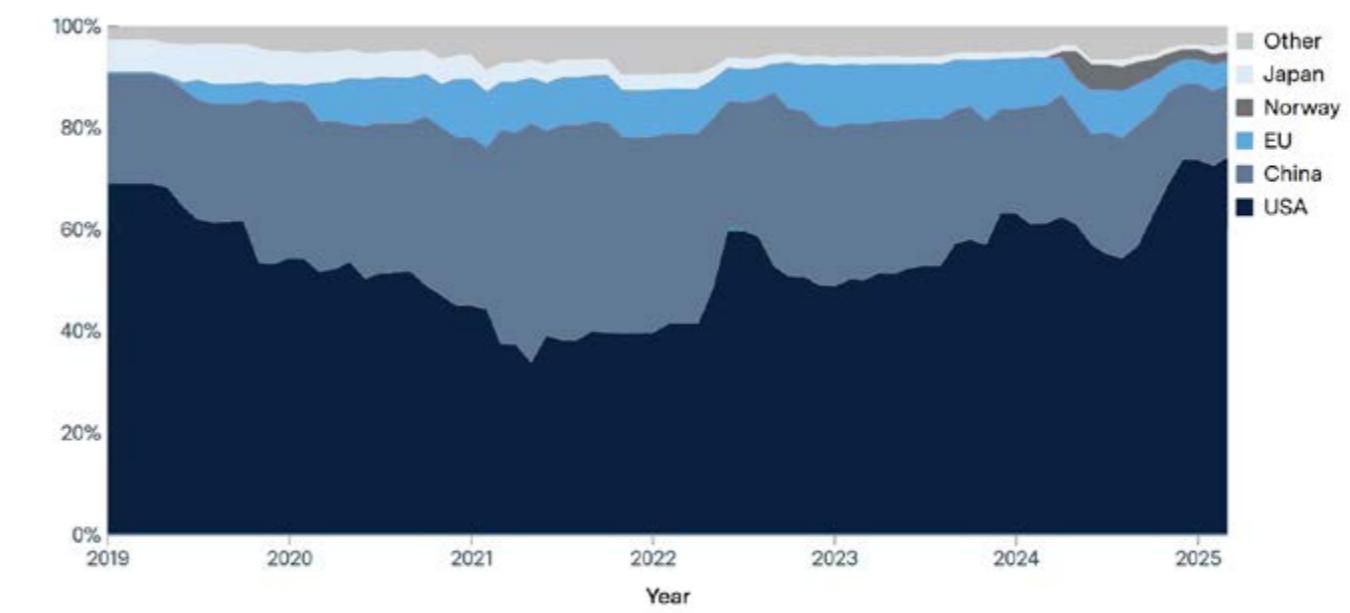
Europe's position in the AI infrastructure supercycle is defined by a structural compute gap. Research published by Epoch AI estimates that as of March 2025, the EU holds only 4.8% of global AI supercomputer performance, trailing far behind the United States (74.4%) and China (14.1%). While these figures must be interpreted carefully (they capture only an estimated 10-20% of total AI compute and over-represent public systems), they signal a clear and widening disparity in scale and, critically, in execution velocity. The EU's most ambitious plans, such as the "AI Gigafactories" initiative, aim to deploy 100,000+ advanced chips per site. In contrast, xAI's Colossus cluster, featuring 200,000 H100 GPUs, was reportedly built in just 122 days. Planned facilities like the UAE-US AI Campus in Abu Dhabi are targeting a scale of over 2 million chips by 2030, an order of magnitude beyond current European ambitions.

This lag is not accidental. Europe is fundamentally ill-suited for the significant, centralised AI training factories seen in the US due to compounding structural headwinds. Brute-force, gigawatt-scale build-outs are deterred by strained power grids, with hubs like Dublin and Amsterdam facing multi-year grid connection queues, reportedly extending 7-10 years. High energy

costs: persistently high wholesale energy costs make brute-force training economically challenging compared to other regions. While mandatory, the regulatory environment (EED, permitting complexity) imposes a financial and temporal drag on large-scale development that the US market, despite its own grid constraints, does not replicate. This reality has so far cemented Europe's role as a "strategic demand sink" for foundational models trained elsewhere. It exacerbates a dangerous compute sovereignty gap, where the continent's future economic activity risks becoming reliant on foreign-owned infrastructure.

This dynamic highlights an opportunity to address a compute sovereignty gap through targeted investments. The continent's economic future may increasingly depend on localised infrastructure, but this challenge could stimulate demand for AI compute, potentially backed by government incentives such as the EU's €1.2bn allocation for AI supercomputing under the EuroHPC Joint Undertaking. Ultimately, Europe's lag in centralised training clusters creates a strategic pivot toward capturing the burgeoning demand for localised inference.

FIGURE 38: SHARE OF AGGREGATE GPU CLUSTER PERFORMANCE (MARCH 2025)

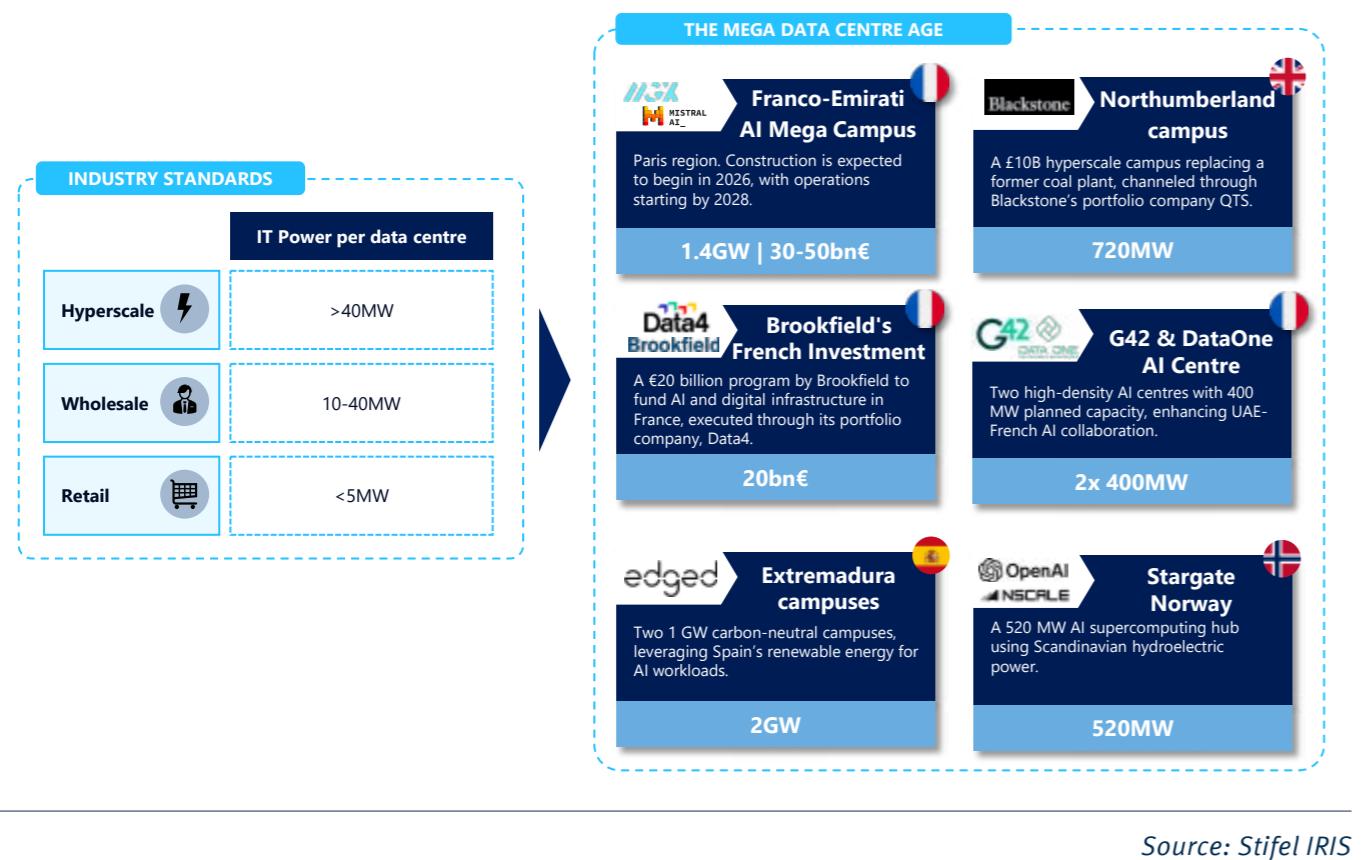


Source: Epoch AI, Stifel IRIS

While Europe's 2025 AI infrastructure announcements pale in scale compared to U.S. announcements such as the \$500bn Stargate (though we note that significant new projects are being announced almost every week, in what has been called a “bragawatts” phenomenon), European countries have so far been able to attract major investments, signalling viability in the global race through sovereign, efficiency-focused builds amid regulatory and energy hurdles. Key initiatives include Data4’s expansive Paris-Marcoussis campus, backed by Brookfield’s €20bn commitment (announced February 2025), which plans to triple capacity to over 1.5 GW

across France. In Spain, MERLIN Properties and Edged Energy’s two gigawatt-scale AI data centre campuses in Extremadura, announced in January 2025 with regional government backing, leverage low-cost renewables and waterless cooling for ultra-efficient operations, positioning the region as a green AI hub. France’s joint venture with the UAE, unveiled at the May 2025 Choose France Summit, establishes Europe’s largest AI campus in Île-de-France via MGX, Bpifrance, Mistral AI, and NVIDIA, scaling to 1.4 GW to drive sovereign AI amid bilateral cooperation.

FIGURE 39: MEGASCALE DATA CENTRE ANNOUNCEMENTS ACROSS EUROPE



## 3.2 IDENTIFYING EUROPEAN CHAMPIONS

The AI revolution’s disruption of legacy cloud stacks creates a structural opening for European companies. Success hinges on mastering a new playbook: developing AI-native stacks, weaponising European DNA like sovereignty and sustainability as a commercial moat, and executing bypass strategies on high-value niches. We map the resulting investment theses across four key verticals.

### THE CRITERIA FOR SUCCESS IN THE EUROPEAN DATA CENTRE MARKET

The European data centre landscape is defined by the entrenched dominance of US companies and tech giants, whose platforms penetrate nearly every layer of the value chain. Yet, this has not precluded the success of local companies. As detailed in Section 2, European champions have historically thrived by applying proven bypass strategies, capturing valuable niches, from low-cost simplicity to regional sovereignty, within a rapidly expanding market.

The generative AI revolution, far from cementing this status quo, now represents a structural opening.

The legacy cloud stack, built for general-purpose workloads, is ill-suited for the demands of high-performance AI. This technological shift creates a rare opportunity for new architectures to “leapfrog” established companies, evidenced by the rise of specialist European neoclouds and the growing potential for public-private support for sovereign AI initiatives. The game is no longer just about taking slices of the old pie; it’s about identifying and capturing the new, high-value niches created in the giants’ shadows.

We believe that for local European companies, future

success and outperformance will be defined by mastering one or more of these three capabilities:

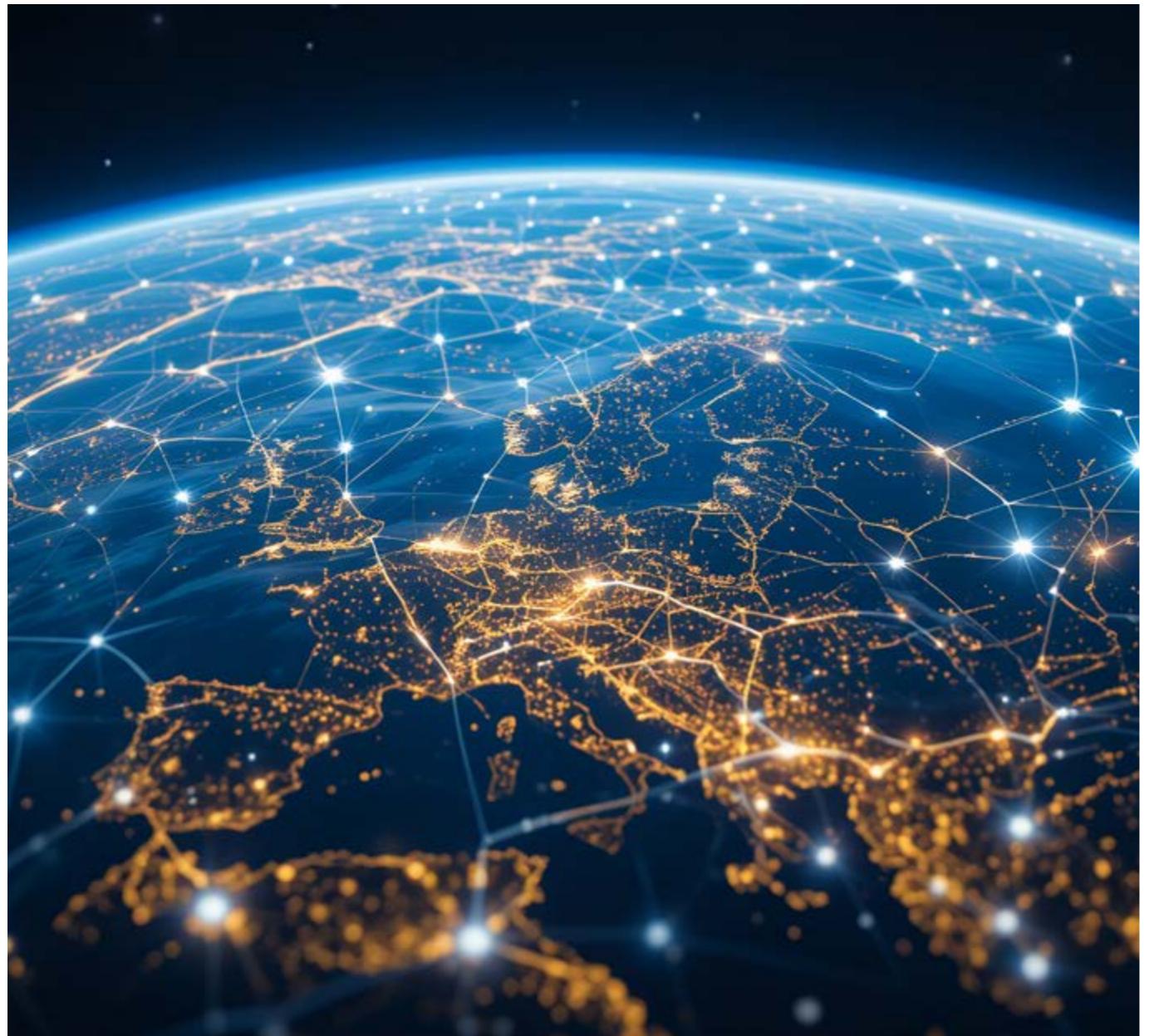
- **Mastery of the AI-native stack:** The technological baseline has fundamentally shifted. Success now demands deep expertise not only in high-density physical infrastructure (liquid cooling, advanced electrics) but also in the complex software stacks required for GenAI services. This is the new barrier to entry. European specialists like Nebius and Fluidstack have already demonstrated this capability, winning major contracts by delivering the operational excellence and AI-native architecture that even US giants require.
- **Weaponising European market DNA:** What were once seen as regulatory burdens must be wielded as commercial advantages. Mastering the complexities of sovereignty (Data Act, GDPR, national certifications like SecNumCloud) and sustainability (EED, heat reuse mandates, carbon lifecycle tracking) creates a powerful, defensible moat. This positions local companies as essential, trusted partners for European enterprises and public bodies navigating these non-negotiable requirements.

- **Aggressively executing the new bypass playbook:**

The AI wave is causing the market to fragment, and creates new gaps. Value will be captured not by competing head-on, but by identifying and dominating defensible niches. The “bypass strategies” (Section 2.3) remain the core playbook, but the targets are new: “no-frills” inference platforms, latency-critical edge AI, sovereign-by-design clouds, and specialised solutions for Europe’s unique industrial and R&D landscape. The core mission remains: find the specialised niche the hyperscalers cannot, or will not, serve.

This optimism, however, must be qualified. We must not be naive about the formidable headwinds Europe faces. The region contends with structurally high energy costs, a fragmented tech ecosystem lacking the integration of Silicon Valley or Shenzhen, and a clear absence of native hyperscale giants.

Yet, this very constraint creates the market’s most powerful competitive moat. Agile companies that have secured access to reliable, sustainable power, or those that develop superior energy-efficiency technologies, can weaponise this scarcity. This power-access advantage will define the beneficiaries.



## MAPPING EUROPEAN COMPANIES

We believe the most compelling, high-potential theses in the European data centre landscape lies in four specialised verticals. We have identified: 1) the AI-Native Neoclouds, outperforming on sovereignty and operational excellence; 2) the specialised data centre operators, leading on power-access and AI-optimised infrastructure; 3) foundational software, excelling by abstracting commoditised compute; and 4) industrial “Picks and Shovels”, solving high-density power and cooling bottlenecks.

### Vertical #1: The AI-native neoclouds

This is the primary vector for challenging US hyperscaler dominance in high-value AI workloads. The advantage is no longer mere hardware access, but the deep, AI-native software layers that delivers superior TCO and operational excellence. In Europe, this model has a unique right-to-win by embedding sovereignty-by-design, targeting the critical compute gap for national AI models and regulated industries that cannot, or will not, use US-based platforms. They are transforming from low-cost alternatives into high-value, strategic enablers. For investors, these operators represent a high-margin, service-oriented opportunity in the most demanding segment of the AI infrastructure market.

### Vertical #2: Specialised data centre operators

This is a critical infrastructure opportunity defined by strategic bypass. While hyperscalers pursue brute-force scale, the European opportunity for local operators is specialisation. The beneficiaries will be those who master the power-access bypass, such as by securing hyperscale sites in power-rich peripheries (Nordics, Spain), and those who monetise sovereignty and interconnection in core hubs. They provide the essential, trusted physical ground for the hybrid and private AI deployments that European enterprises and governments structurally demand.

### Vertical #3: Foundational software

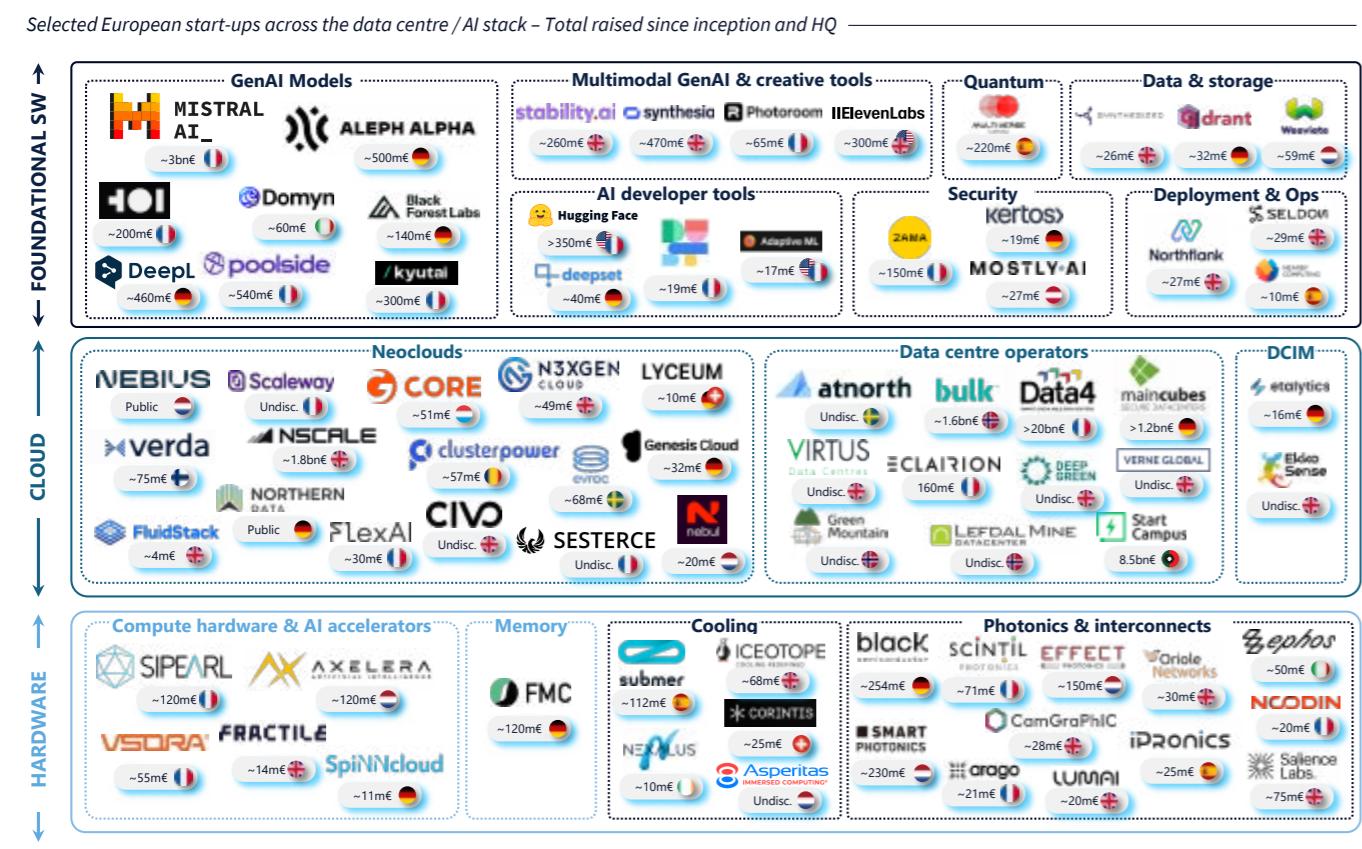
We posit that the AI investment cycle will follow a pattern analogous to previous IT cycles: infrastructure first,

followed by platforms and then applications. Currently, the prolonged infrastructure build-out, coupled with underwhelming enterprise adoption of AI applications, means significant value remains concentrated in the infrastructure layer; the era of widespread applicative AI is yet to dawn. Nevertheless, we identify a compelling opportunity in the foundational software layer, a crucial intermediary between physical infrastructure and end-user applications. The most compelling prospects lie in hardware-agnostic, pure-play software platforms governing the entire AI lifecycle, encompassing data management, developer tooling, and security. Being decoupled from hardware dependencies, these platforms are emerging as a powerful driver of vendor lock-in and sustainable, high-margin revenue across the European AI market.

### Vertical #4: Industrial “Picks and Shovels”

The industrial thesis of AI lies within the disruptive technologies it enables and the bottlenecks it creates. The AI-driven shift to extreme rack densities (100kW+) creates a classic innovator’s dilemma, rendering decades of traditional air-cooling expertise obsolete. While large incumbents (Schneider Electric, Vertiv, Legrand) are geared to optimise a legacy model, the most valuable opportunities lie with agile, specialised challengers focused on disruptive liquid cooling (Direct-to-Chip, immersion) and advanced power management. These are the technologies that enable the AI roadmap. A similar, albeit far more challenging, ‘picks and shovels’ opportunity lies in the core IT equipment stack itself (processors, networking). This is an immense, high-growth market, but it remains dominated by formidable US incumbents. For new European entrants, the capital and R&D barriers are staggering. Yet, the strategic rewards for success are equally high, making the emergence of a future European champion in AI silicon a high-risk, high-reward aspiration.

FIGURE 40: A SNAPSHOT OF EUROPE'S DATA CENTRE & AI STARTUP ECOSYSTEM



Source: Stifel IRIS. Funding totals are not strictly comparable and may include, debt, committed capital or public grants.

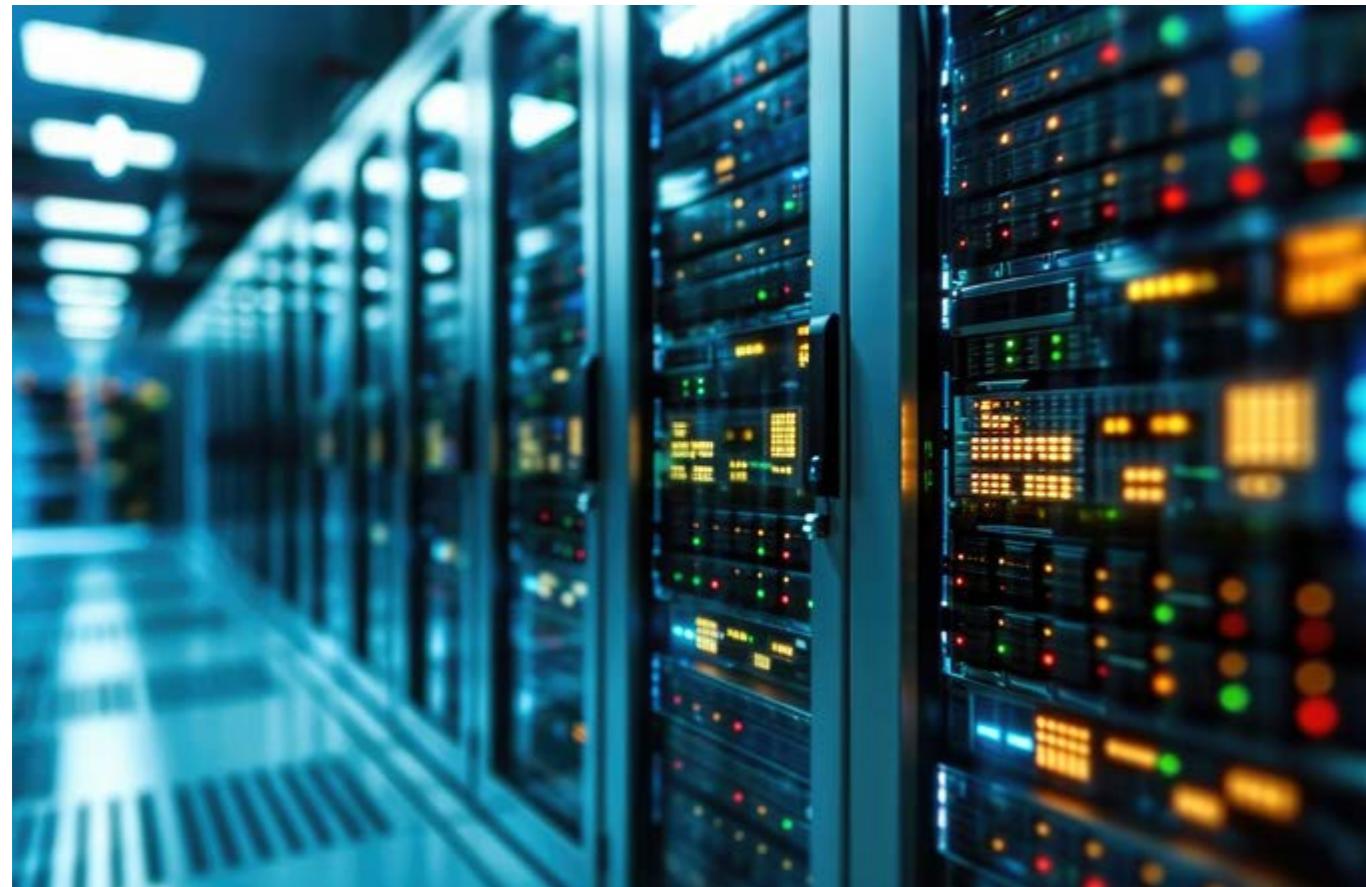
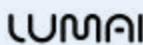
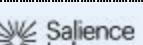
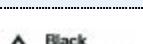
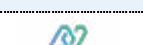


FIGURE 41: SELECTED DATA-CENTRE-RELATED FUNDING ROUNDS IN EUROPE SINCE JANUARY 2023

Deal date	Company	HQ location	Sector	Deal amount	Financing round	Key investors
Nov-2025	<b>FMC</b>	DEU	Memory	100m€	Series C	HV Capital, DeepTech & Climate Fonds, VSquared Ventures
Nov-2025	<b>NCODIN</b>	FRA	Photonics	16m€	Seed Round	MIG AG, Elaia Partners, Earlybird Venture Capital, OVNI Capital, Verve Ventures, Maverick Silicon, PhotonVentures
Oct-2025	<b>NSCALE</b>	UK	Neocloud	433m\$	Pre-Series C	Nvidia, Nokia, Dell Technologies, Blue Owl Capital
Oct-2025	<b>synthesia</b>	UK	Multimodal GenAI & creative tools	200m\$	Series E	GV
Oct-2025	<b>etalytics</b>	DEU	DCIM	16m€	Series A	ebm-papst St. Georgen, M12, ALSTIN Capital
Sep-2025	<b>Asperitas IMMERSED COMPUTING</b>	NLD	Liquid cooling equipment	Undisclosed	Scale up investment round	Stecon Group Global, Invest-NL, Shell Ventures, PDENH (Participatiefonds Duurzame Economie Noord-Holland)
Sep-2025	<b>CORINTIS</b>	CHE	Liquid cooling equipment	24m\$	Series A	BlueYard Capital, Acequia Capital, Founderful, AICONIC VENTURES, A&E Investments, Celsius Industries and XTX Ventures
Sep-2025	<b>MISTRAL AI</b>	FRA	Foundational AI models	1,700m€	Series C	ASML (Lead), DSL Global, Andressen Horowitz, Bpifrance, General Catalyst, Index Ventures, Lightspeed, NVIDIA
Sep-2025	<b>NSCALE</b>	UK	Neocloud	1100m£	Series B	Aker, Nokia, Dell Technologies, T.Capital, Nvidia, Sandton Capital Partners, G Squared, Point72 Ventures, Fidelity Management & Research Company, Blue Owl Capital
Sep-2025	<b>Salience Labs</b>	UK	Photonics	31€m	Later Stage VC	Borski Fund, Cambridge Innovation Capital, Oman Investment Authority, Oxford Sciences Enterprises
Sep-2025	<b>SCINTIL PHOTONICS</b>	FRA	Photonics	54€m	Series B	NVIDIA, NGP Capital, Bosch Ventures, Bpifrance, Applied Ventures, Banque Publique d'Investissement, BNP Paribas Développement, Innovacom, Supernova Invest, Yotta Capital Partners
Sep-2025	<b>verda</b>	FIN	Neocloud	64m\$	Series A	Varma Mutual Pension Insurance, Tesi, byFounders, Skaala, J12 Ventures
Sep-2025	<b>kertos</b>	DEU	Security	14m€	Series A	Portage, 10x Founders, Pi Labs, seed + speed Ventures, Redstone VC
Jul-2025	<b>arago</b>	FRA	Photonics	21m€	Seed Funding	C4 Ventures, Earlybird Venture Capital, GenerativelQ, Protagonist, Visionaries Tomorrow, Co-Founder of Hugging Face
Jul-2025	<b>EFFECT PHOTONICS</b>	NLD	Photonics	57€m	Series D	Innovation Industries, Invest-NL, PhotonVentures, b2venture, Brabantse Ontwikkelings Maatschappij, Matterwave Ventures, PhotonDelta, Smile Invest, ViaSat, Beek Capital, Innovatiefonds Brabant, Optiverder
Jul-2025	<b>Zephos</b>	ITA	Photonics (Glass based PIC)	43m€	Grant	European Chips Act
Jul-2025	<b>SIPPEARL</b>	FRA	AI accelerators	130m€	Series A	Atos, Cathay Venture, Arm, Caisse d'Epargne Rhone Alpes, European Innovation Council, Bpifrance, France 2030
Jun-2025	<b>LYCEUM</b>	DEU, CHE	Neocloud	10.3m€	Seed Round	redalpine, 10x Founders
Jun-2025	<b>MULTIVERSE</b>	ESP	Quantum AI software	189m€	Series B	HP Tech Ventures, Toshiba, CDP Venture Capital, Quantonation, SET Ventures, Forgepoint Capital, Santander Climate Fund, Spri Group
Jun-2025	<b>ZAMA</b>	FRA	Security	54m\$	Series B	Blockchange Ventures, Pantera Capital, Robot Ventures
May-2025	<b>SYNTHEZIZED</b>	UK	AI Data infrastructure software	16.7m£	Series A	Redalpine, Deutsche Bank Ventures, Mercia Ventures, Seedcamp, IQ Capital Partners, Northern Venture Trust, UBS Group

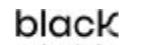
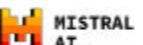
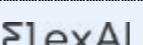
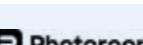
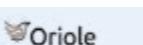
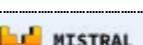
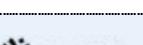
Source: PitchBook, Stifel IRIS

FIGURE 41: SELECTED DATA-CENTRE-RELATED FUNDING ROUNDS IN EUROPE SINCE JANUARY 2023

Deal date	Company	HQ location	Sector	Deal amount	Financing round	Key investors
Apr-2025	 LUMAI	UK	Optical Processor technology	10m\$	Early stage VC	Constructor Capital, State Farm Ventures, Journey Ventures, Qubits Ventures, TIS INTEC Group, IP Group, PhotonVentures
Apr-2025	 N3XGEN CLOUD	UK	Neocloud	45m\$	Series A	Moore & Moore Investments, EIM Capital
Apr-2025	 VSORA	FRA	AI accelerators	40m€	Later Stage VC	Otium Capital, Omnes Capital, Adelie Capital, European Innovation Council
Mar-2025	 CamGroPhIC	UK	Photonics	25€m	Series A	Bosh Ventures, CDP Venture Capital, Frontier IP group, Indaco Venture Partners, Join Capital, NATO Innovation fund, Sony Innovation fund
Mar-2025	 EVROC	SWE	Neocloud	55m€	Series A	blisce/ & Giant Ventures, EQT Ventures, Norrsken VC
Feb-2025	 ECLAIRION	FRA	Data centre operator	50m€	PE Growth	Tikehau Capital
Feb-2025	 Salience Labs	UK	Photonics	24m€	Series A	Applied Ventures, ICM Group, Oxford Science Enterprises, Cambridge Innovation Capital, Braavos Investment Advisers, Silicon Catalyst, Jalal Bagherli, Strategic Investment Fund
Feb-2025	 IIElevenLabs	USA, UK	Multimodal GenAI & creative tools	250m\$	Series C	Andreessen Horowitz, ICONIQ Growth, Salesforce Ventures, Deutsche Telekom, LG Technology Ventures, HubSpot Ventures, NTT Docomo Ventures, RingCentral, FT Ventures, KPN Ventures, T.Capital, Smash Capital
Jan-2025	 iPronics	ESP	Software-defined programmable photonics	20m€	Series A	Triatomic Capital, Bosch Ventures, Amadeus
Jan-2025	 synthesia	UK	Multimodal GenAI & creative tools	147m€	Series D	Atlassian Ventures, New Enterprise Associates, PSP Partners, GV, Adobe Ventures, Nvidia, MMC Ventures, FirstMark, Kleiner Perkins, Accel, WIL, AI Futures Fund, Octant Ventures, UCL Business, Schroders, First Serve Ventures
Dec-2024	 NSCALE	UK	Neocloud	155m\$	Series A	Sandton Capital Partners, VentureSouq, Florence Capital Advisors, BlueSky Asset Management, Kestrel Investment Management
Nov-2024	 Black Forest Labs	DEU	Text-to-image models	126m\$	Series A	Andreessen Horowitz, BroadLight Capital, Creandum, Northzone Ventures, SV Angel, Cherry Ventures, LEA Partners
Nov-2024	 LUMAI	UK	Optical Processor technology	11.5m\$	Early stage VC	Constructor Capital, TIS INTEC Group, State Farm Ventures, Journey Ventures, Qubits Ventures, PhotonVentures, IP Group, LIFT
Nov-2024	 Northflank	UK	MLOps & cloud deployment	17.3m€	Series A	Bain Capital Ventures, Vertex Ventures US, Explorer34, MongoDB Ventures, Tapestry VC, Uncorrelated Ventures, Pebblebed, Kindred Ventures
Oct-2024	 poolside	FRA	Generative AI coding platform	500m\$	Series B	Bain Capital Ventures, Rashmi Gopinath, Nvidia, Citi Ventures, Capital One Ventures, HSBC, LG Technology Ventures, eBay Ventures, Ericsson Ventures, SoftBank Investment Advisers, Premji Invest (US), BAM Elevate
Oct-2024	 submer	ESP	Liquid cooling equipment	55.5m€	Series A	M&G, Planet First Partners, Norrsken VC, Mundi Ventures
Oct-2024	 verda	FIN	Neocloud	20m€	Seed Round	byFounders, J12 Ventures, Oskari Saarenpää, Tuomo Riekki, Lasse Espeholt, Nal Kalchbrenner, Henrik Rosendahl, Anders Bo Pedersen, Ari Tulla, Maaike Bryon, LocalTapiola, Nordea Funds
Sep-2024	 Zephos	ITA	Photonics (Glass based PIC)	8.5m€	Seed Round	Exor, Collaborative Fund, Unruly Capital, 2100 Ventures, Green Sands Equity, Inflexor, Silicon Roundabout Ventures, Joe Zadeh, Diego Piacentini, Simone Severini, Club degli investitori
Aug-2024	 Black Forest Labs	DEU	Text-to-image models	31m\$	Seed Round	Andreessen Horowitz, General Catalyst, Mäatch VC, Vladlen Koltun, Timo Aila, Brendan Iribe, Michael Ovitz, Garry Tan
Aug-2024	 Oriole Networks	UK	Photonics	17.5m€	Series A	Plural, XTX Ventures, Dorilton Ventures, UCL Technology Fund, and Clean Growth Fund
Jul-2024	 CORE	LUX	Neocloud	60m\$	Series A	Wargaming, Constructor Capital, Han River Partners

Source: PitchBook, Stifel IRIS

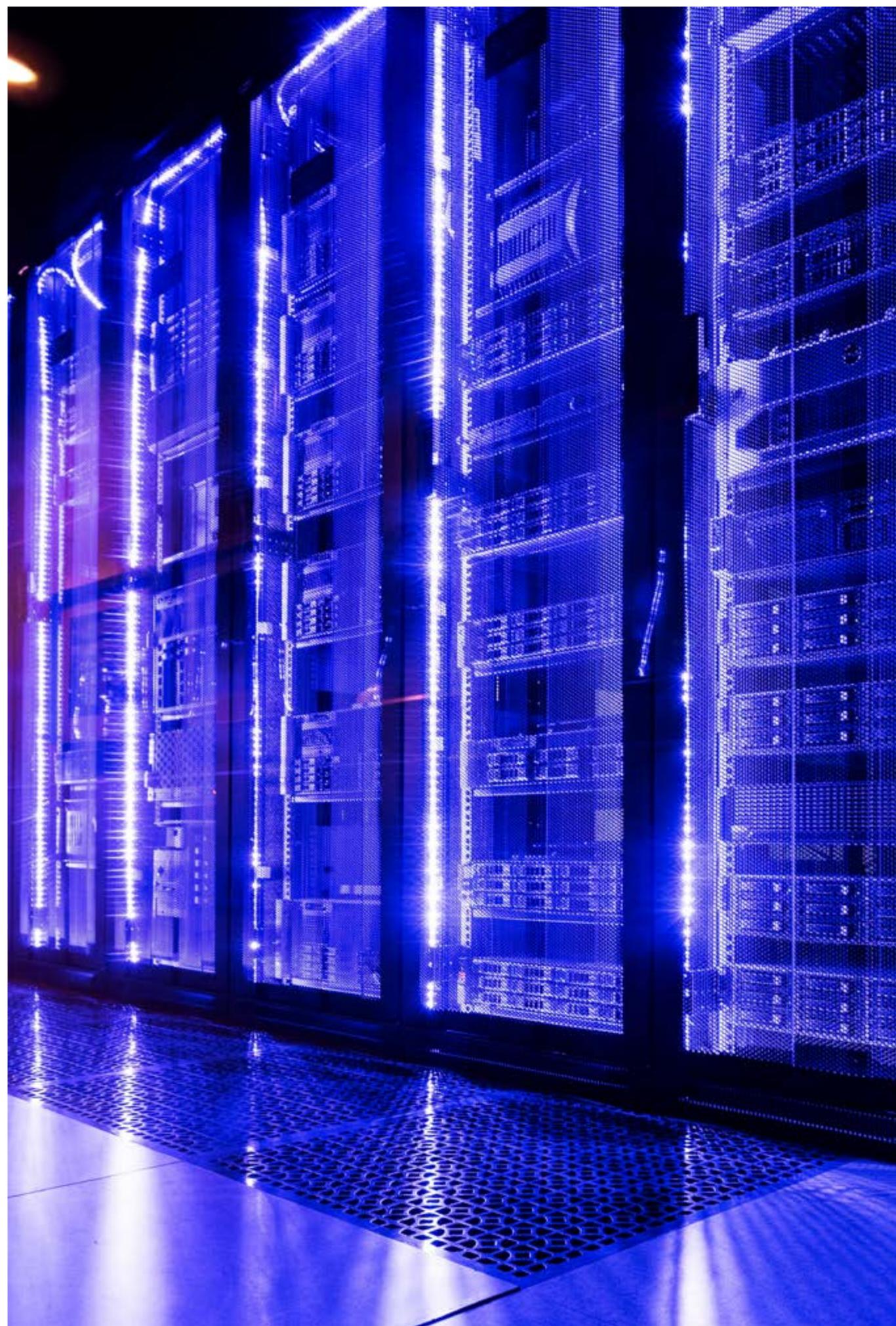
FIGURE 41: SELECTED DATA-CENTRE-RELATED FUNDING ROUNDS IN EUROPE SINCE JANUARY 2023

Deal date	Company	HQ location	Sector	Deal amount	Financing round	Key investors
Jun-2024	 AXELERA	NLD	AI accelerators	75m\$	Series B	Samsung Catalyst Fund, Innovation Industries, CDP Venture Capital, Verve Ventures, European Innovation Council, Invest-NL (Investment Company), Kinled Holding
Jun-2024	 black semiconductor	DEU	Graphene-based photonic interconnects	254m€	Series A	Porsche Ventures, Project A Ventures, Vsquared, Scania Growth Capital, Cambium Capital Management, Capnamic Ventures, East Hill Equity, NRW Bank, Onsight Ventures, TechVision Fund
Jun-2024	 Dynamilis	FRA	AI Infrastructure & developer tools	16m\$	Series A	Sequoia Capital, Seedcamp, Connect Ventures, XYZ Venture Capital, Prototype Capital, Motier Ventures
Jun-2024	 MISTRAL AI	FRA	Foundational AI models	600m€	Series B	General Catalyst, Lightspeed Venture Partners, Andreessen Horowitz, Nvidia, Samsung Venture Investment Corporation, and Salesforce Ventures
Jun-2024	 nebul	NLD	Neocloud	20m€	PE Growth	BeStacking
May-2024	 DeepL	DEU		320m\$	Series B	IVP, Index Ventures, Teachers' Venture Growth, Atomico, WiL (World Innovation Lab), ICONIQ Growth, Ontario Teachers' Pension Plan
May-2024	 FRACTILE	UK	AI accelerators	12m€	Seed round	Oxford Science Enterprises, Kindred Capital and NATO Innovation Fund
May-2024	 4AI	FRA	Foundational AI models / Agentic AI	220m\$	Seed	Eric Schmidt, Xavier Niel, Yuri Milner, Aglaé Ventures, Motier Ventures, Accel, Bpifrance, Creadum, Elaia Partners, Eurazeo, FirstMark Capital, Visionaries Club, Amazon, Samsung
Apr-2024	 FlexAI	FRA	Neocloud	30m\$	Seed Round	Alpha Intelligence Capital, Elaia Partners, Heartcore Capital, Frst Capital, Partech, Better Angle, Karim Beguir, Bpifrance, Motier Ventures
Mar-2024	 Adaptive ML	USA, FRA	AI Infrastructure & developer tools	20m\$	Seed Round	Index Ventures, ICONIQ Capital, Databricks, IRIS, Kima Ventures, Factorial Capital, Xavier Niel, Olivier Pomel, Dylan Patel, Tri Dao, Financière Saint James, SemiAnalysis, Motier Ventures
Mar-2024	 MULTIVERSE COMPUTING	ESP	Quantum AI software	27m€	Series A	Quantation, Columbus Venture Partners, Audacia, INDI Partners, CDP Venture Capital (Rome), QAI Ventures, Redstone VC, European Innovation Council
Mar-2024	 ZAMA	FRA	Security	73m\$	Series A	Protocol Labs, Multicoin Capital, Stake Capital, Metaplanet Holdings, VSquared Ventures, Blockchange Ventures, Portal Ventures, Juan Benet, Anatoly Yakovenko, Gavin Wood, Tarun Chitra, Julien Bouteloup
Mar-2024	 Photoroam	FRA	Multimodal GenAI & creative tools	55m\$	Series B	Balderton Capital, Aglaé Ventures, Meta Platforms, FJ Labs, Kima Ventures, Adjacent, Y Combinator, Yann LeCun, Zehan Wang
Jan-2024	 deepset	DEU	AI Infrastructure & developer tools	30m€	Later stage VC	Balderton Capital, GV, Harpoon Ventures, System.One, Celero Ventures, Lunar Ventures and Aequia Capital, SignalRank
Jan-2024	 Oriole Networks	UK	Photonics	10m€	Seed Round	XTX Ventures, Dorilton Ventures, Clean Growth Fund, UCL Technology Fund, Innovate UK
Jan-2024	 drant	DEU	AI Data infrastructure software	28m\$	Series A	Spark Capital, 42CAP, Unusual Ventures
Jan-2024	 IIElevenLabs	USA, UK	Multimodal GenAI & creative tools	80m\$	Series B	Andreessen Horowitz, Nat Friedman, Daniel Gross, Salesforce Ventures, Smash Capital, BroadLight Capital, Credo Ventures, SV Angel, Sequoia Capital
Dec-2023	 Domyn	ITA	Foundational AI models	11m\$	Accelerator / Incubator	Endeavor Italy
Dec-2023	 MISTRAL AI	FRA	Foundational AI models	415m\$	Series A	Andreessen Horowitz, Lightspeed, Salesforce, BNP Paribas, CMA-CGM, General Catalyst, Elad Gil, Conviction
Dec-2023	 NSCALE	UK	Neocloud	30m\$	Seed Round	Undisclosed
Nov-2023	 ALEPH ALPHA	DEU	Foundational AI models	>500m\$	Series B	Innovation Park Artificial Intelligence (IpaI), Bosch Ventures, Schwarz Group

Source: PitchBook, Stifel IRIS

FIGURE 41: SELECTED DATA-CENTRE-RELATED FUNDING ROUNDS IN EUROPE SINCE JANUARY 2023

Deal date	Company	HQ location	Sector	Deal amount	Financing round	Key investors
Nov-2023	<b>ECLAIRION</b>	FRA	Data centre operator	110m€	PE Growth	Tikehau Capital
Oct-2023	<b>Genesis Cloud</b>	DEU	Neocloud	20.5m\$	Series A	Undisclosed
Aug-2023	<b>Hugging Face</b>	USA, FRA	AI Infrastructure & developer tools	235m\$	Series D	Salesforce Ventures, Amazon.com, Intel, AMD, IBM, Qualcomm Ventures, GV, NVIDIA GPU Ventures, Mozilla Ventures, Nvidia, Premji Invest, Firebreak Ventures, Sound Ventures, Alspace Ventures, MarketX, Alphabet
Aug-2023	<b>N3XGEN CLOUD</b>	UK	Neocloud	11.2m\$	Seed Round	Undisclosed
Aug-2023	<b>poolside</b>	FRA	Generative AI coding platform	97m\$	Series A	Felicis, Redpoint Ventures, Xavier Niel, Air Street Capital, Scribble Ventures, First Capital, New Wave VC, Abstraction Capital, Point Nine Capital, Bain Capital Ventures, Rodolphe Saadé, Motier Ventures, Bpifrance
Jul-2023	<b>SMART PHOTONICS</b>	NLD	Photonics	111€m	Later Stage	ASML, NXP, Invest-NL, PhotonDelta
Jun-2023	<b>EVTEC</b>	SWE	Neocloud	13m€	Seed Round	EQT Ventures, Norrsken VC
Jun-2023	<b>MISTRAL AI</b>	FRA	Foundational AI models	113m\$	Seed	Lightspeed, Redpoint, Index Ventures, Xavier Niel, JCDecaux, Rodolphe Saadé, Motier Ventures, La Famiglia, Headline, Exor Ventures, Sofina, First Minute Capital, LocalGlobe
Jun-2023	<b>VSORA</b>	FRA	AI accelerators	13m€	Later Stage VC	Undisclosed
Jun-2023	<b>synthesia</b>	UK	Multimodal GenAI & creative tools	71m€	Series C	Accel, GV, NVentures, Setcoin Group, MMC Ventures, Kleiner Perkins, FirstMark Capital, Alex Wang, Olivier Pomel, Amjad Masad
May-2023	<b>AXELERA</b>	NDL	AI accelerators	50m\$	Series A	Innovation Industries, imec.xpand, CDP Venture Capital, Verve Ventures, Netherlands Enterprise Agency, Federal Holding and Investment Company
May-2023	<b>poolside</b>	FRA	Generative AI coding platform	26m\$	Seed Round	Undisclosed
May-2023	<b>ElevenLabs</b>	USA, UK	Multimodal GenAI & creative tools	18.5m\$	Series A	Andreessen Horowitz, Nat Friedman, Daniel Gross, Embark Studios, Storytel, TheSoul Publishing, Concept Ventures, Credo Ventures, Creator Ventures, SV Angel, Mike Krieger, Brendan Iribe, Mustafa Suleyman, Tim O'Reilly
Apr-2023	<b>Weaviate</b>	NLD	AI Data infrastructure software	50m\$	Series B	Index Ventures, ING Ventures, Battery Ventures, Zetta Venture Partners, Cortical Ventures, New Enterprise Associates
Mar-2023	<b>EFFECT PHOTONICS</b>	NLD	Photonics	40m\$	Series C	Innovation Industries, Invest-NL
Feb-2023	<b>SELDON</b>	UK	MLOps & cloud deployment	16.5m€	Series B	Bright Pixel Capital, Amadeus Capital Partners, Cambridge Innovation Capital, AlbionVC, Kings Arms Yard VCT, Uncharted Ventures
Jan-2023	<b>DeepL</b>	DEU		100m\$	Series B	IVP, Bessemer Venture Partners, Atomico, b2venture, Benchmark Capital Holdings, Wil, Flat Capital, Eniac Ventures, Allison Pickens Ventures, Illusian Founder Office



Source: PitchBook, Stifel IRIS

## WHITE PAPER AUTHOR



**Antoine Lebourgeois**

Associate

Paris

[antoine.lebourgeois@stifel.com](mailto:antoine.lebourgeois@stifel.com)

## IRIS TECHNOLOGY RESEARCH



**Cédric Rossi**

Vice President

Paris

[cedric.rossi@stifel.com](mailto:cedric.rossi@stifel.com)



**Valentin-Paul Jahan**

Vice President

Paris

[valentin-paul.jahan@stifel.com](mailto:valentin-paul.jahan@stifel.com)



**Gabriel Santier**

Associate

Paris

[gabriel.santier@stifel.com](mailto:gabriel.santier@stifel.com)



**Aurélien Deside**

Associate

Paris

[aurelien.deside@stifel.com](mailto:aurelien.deside@stifel.com)



**Clément Genelot**

Vice President

Paris

[clement.genelot@stifel.com](mailto:clement.genelot@stifel.com)



**Thomas Mordelle**

Associate

Paris

[thomas.mordelle@stifel.com](mailto:thomas.mordelle@stifel.com)



**Mahaut Arnaud**

Associate

Paris

[mahaut.arnaud@stifel.com](mailto:mahaut.arnaud@stifel.com)

## IRIS LEADERSHIP



**Paul de Froment**

Managing Director  
Co-Head of IRIS

Paris

[paul.defroment@stifel.com](mailto:paul.defroment@stifel.com)



**Damien Choplain**

Managing Director  
Co-Head of IRIS

Paris

[damien.choplain@stifel.com](mailto:damien.choplain@stifel.com)

# LEGAL DISCLAIMER

This white paper is provided on a confidential basis for informational purposes only and is not intended to, and does not, constitute a recommendation with respect to any potential transaction or investment. Any opinions expressed are solely those of Stifel and applicable only as at the date of this white paper. This white paper is necessarily based upon economic, market, financial and other conditions as they exist on, and on the information made available to Stifel as of, the date of this white paper, and subsequent developments may affect the analyses or information set forth in this white paper. This white paper does not purport to give legal, tax or financial advice. Recipients should not rely on the information contained in this white paper and must make their own independent assessment and such investigations as they deem necessary. Stifel is not soliciting any action based upon this white paper. This white paper does not constitute or form part of any offer or invitation to sell, or issue, or any solicitation to any offer to purchase or subscribe for, any shares, financial instruments, or other securities, nor shall it (or any part of it), or the fact of its distribution, form the basis of, or be relied on in connection with or act as any inducement to enter into, any contract whatsoever relating to any securities, financial instruments or financial services of Stifel or of any other entity or constitute an invitation or inducement to any person to underwrite, subscribe for or otherwise acquire securities. The information in this white paper is not complete and is based upon information that Stifel considers reliable, but it has not been independently verified. Stifel does not represent, guarantee, or warrant, expressly or implicitly, that this white paper or any part of it is valid, accurate or complete (or that any assumptions, data or projections underlying any estimates or projections contained in the white paper are valid, accurate or complete), or suitable for any particular purpose, and it should not be relied upon as such. Stifel accepts no liability or responsibility to any person with respect to or arising directly or indirectly out of the contents of or any omissions from this white paper. As a multi-disciplined financial services firm, Stifel regularly seeks investment banking assignments and compensation from companies, which could include companies mentioned within this paper. Our European Policies for Managing Conflicts of Interest are available at [www.stifel.com/institutional/ImportantDisclosures](http://www.stifel.com/institutional/ImportantDisclosures).

The distribution of this white paper may be restricted by law. Accordingly, this white paper may not be distributed in any jurisdiction except in accordance with the legal requirements applicable to such jurisdiction. Persons into whose possession this document comes are required to inform themselves about and to observe any such restrictions. This white paper is only be addressed to and directed at specific addressees who: (A) if in member states of the European Economic Area (the "EEA"), are persons who are "qualified investors" within the meaning of Article 2(e) of Regulation (EU) 2017/1129 (as amended) (the "Prospectus Regulation") ("Qualified Investors"); (B) if in the United Kingdom, are Qualified Investors within the meaning of Article 2(e) of the Prospectus Regulation as it forms part of domestic law by virtue of the EU (Withdrawal) Act 2018 (as amended from time to time) and who are: (i) persons having professional experience in matters relating to investments who fall within the definition of "investment professionals" in Article 19(5)

of the Financial Services and Markets Act 2000 (Financial Promotion) Order 2005 (the "Order"); or (ii) high net worth entities falling within Article 49(2) (a) to (d) of the Order; or (C) are other persons to whom it may otherwise lawfully be communicated (all such persons referred to in (B) and (C) together being "Relevant Persons"). This white paper must not be acted or relied on in (i) the United Kingdom, by persons who are not Relevant Persons; (ii) in any member state of the EEA by persons who are not Qualified Investors; or (iii) in the United States ("U.S.") by persons who are not Qualified Institutional Buyers ("QIBs") as defined in and pursuant to Rule 144A under the U.S. Securities Act of 1933, as amended. Any investment activity to which this white paper relates (i) in the United Kingdom is available only to, and may be engaged in only with, Relevant Persons; (ii) in any member state of the EEA is available only to, and may be engaged in only with, Qualified Investors; and (iii) in the U.S. is available only to, and may be engaged in only with, QIBs. If you have received this white paper and you are (A) in the United Kingdom and are not a Relevant Person; (B) are in any member state of the EEA and are not a Qualified Investor; or (C) are in the U.S. and are not a QIB, you must not act or rely upon or review the white paper and must return it immediately to your Stifel representative (without copying, reproducing or otherwise disclosing it (in whole or in part).

No person shall be treated as a client of Stifel or be entitled to the protections afforded to clients of Stifel, solely by virtue of having received this document..

### Independence of Research

Stifel prohibits its employees from directly or indirectly offering a favourable research rating or specific price target, or offering to change a rating or price target, as consideration or inducement for the receipt of business or for compensation.

### Basis of Presentation

References herein to "Stifel" collectively refer to Stifel, Nicolaus & Company, Incorporated, Stifel Nicolaus Europe Limited ("SNEL"), Stifel Europe AG ("STEA"), Stifel Europe Advisory GmbH, Stifel Nicolaus Canada Incorporated, Stifel Europe Limited, Stifel Europe Securities SAS, Stifel Europe GmbH, Stifel Europe AS and other affiliated broker-dealer subsidiaries of Stifel Financial Corp. SNEL also trades as Keefe, Bruyette & Woods ("KBW"). For a list of Stifel affiliates and associated local regulatory authorisations please see here: [www.stifel.com/disclosures/emaildisclaimers](http://www.stifel.com/disclosures/emaildisclaimers). References herein to "Stifel Financial" refer to Stifel Financial Corp. (NYSE: SF), the parent holding company of Stifel and such other affiliated broker-dealer subsidiaries. Unless otherwise indicated, information presented herein with respect to the experience of Stifel also includes transactions effected and matters conducted by companies acquired by Stifel (including pending acquisitions publicly announced by Stifel), or by Stifel personnel while at prior employers.

If you no longer wish to receive these marketing communications, please e-mail [StifelEurope.GDPR@stifel.com](mailto:StifelEurope.GDPR@stifel.com) and we will arrange to have you taken off the relevant mailing list(s).

Copyright 2026 Stifel. All rights reserved.

# STIFEL

**London, UK**

150 Cheapside  
London, EC2V 6ET  
Tel: +44 20 7710 76

**Paris, France**

26, Avenue des Champs-Elysées  
75008 Paris  
Tel: +33 1 56 68 75 00

**Munich, Germany**

Königinstraße 9  
80539 Munich  
Tel: +49 89 242 262 11

**Frankfurt, Germany**

Skyper Taunusanlage 1  
60329 Frankfurt am Main  
Tel: +49 69 247 4140

**Stockholm, Sweden**

Mäster Samuelsgatan 1,  
111 44 Stockholm