



Executive Summary

Data centers are the backbone of our digital economy, powering everything from cloud services to critical business applications. This handbook provides a comprehensive overview of data centers for business professionals, covering their purpose, types, lifecycle, and key technical architectures. We delve into power and cooling systems, network connectivity, operational reliability standards, market hubs, and industry players, as well as emerging trends like AI-driven demand and sustainability initiatives. **Key takeaways:** Modern data centers come in various forms (from massive *hyperscale* campuses to *edge* micro-sites) and serve diverse users; they require robust power and cooling infrastructure with redundancy to ensure uptime; connectivity and interconnection within facilities enable low-latency data exchange; and the industry is evolving rapidly to meet growing, more compute-intensive demand while addressing power grid constraints and environmental impacts. For example, demand for digital services is rising both in volume and in compute intensity, pushing the limits of existing power and cooling systems (Uptime Institute, 2024)¹ ². At the same time, new AI workloads are driving unprecedented power densities, straining electrical grids and spurring innovation in cooling and on-site power solutions (DCD, 2025)³. In summary, data centers are critical infrastructure enabling today's digital services, and understanding their design, operation, and market context is essential for stakeholders in the technology and real estate value chain.

1. What is a Data Center?

A **data center** is a facility housing computing and network equipment (servers, storage devices, networking gear) that stores, processes, and distributes data for businesses and internet services. In essence, it is the physical home of the “cloud” – every email, streamed video, or online transaction ultimately runs on servers in one or more data centers. Data centers provide secure, uninterrupted environments for IT hardware, including power supply, cooling, physical security, and high-speed network connections. They are designed for *continuous operations*: reliable power (often with backup systems) and cooling are essential to prevent downtime⁴ ⁵. Large data centers typically contain rows of racks filled with servers in climate-controlled rooms, supported by dedicated electrical and mechanical systems (U.S. Chamber of Commerce, 2017)⁴ ⁶.

Purpose and Value Chain: Data centers enable organizations to deliver digital services and store critical data. Virtually every industry relies on data centers – from tech and finance (which use them for applications and transaction processing) to healthcare (for electronic records) and government (for public services and defense computing). The *value chain* includes data center **owners/operators** (who build and maintain the facility and its infrastructure), **tenants** or users (enterprises, cloud providers, content companies that put equipment in the data center), and **service providers** (network carriers, hardware vendors, etc.). For example, an e-commerce company might lease space in a data center to host its servers, while the data center operator provides the power, cooling, physical security, and network access those servers need. **Who uses them?** Broadly, two ownership models exist: **enterprise** data centers (built and operated by a single company for its own use, e.g. a bank's private data center) and **colocation** data centers (built and operated by a provider who leases space/power to multiple external customers)⁷ ⁸. Even large cloud and tech companies often use a mix – operating their own facilities and leasing from colocation providers to extend capacity or reach⁹. Small and mid-sized businesses typically use third-party data centers or cloud services

rather than maintaining their own, as over 90% of servers in the U.S. are housed in professional data center facilities (U.S. DOE estimate) ¹⁰.

Physical vs. Cloud Abstraction: The term “cloud” refers to on-demand computing resources often distributed across many data centers, but the cloud still runs on physical data centers. The distinction is that *cloud providers* abstract the physical infrastructure, so users consume computing as a service without worrying about the underlying hardware or location. From a user’s perspective, deploying an app “in the cloud” means it could be running in any number of data centers that the cloud provider operates. Cloud providers organize infrastructure by **regions** and **availability zones**. A **region** is a geographic area (e.g. US East, EU West) with multiple isolated **availability zones** (AZs) inside it. Each availability zone consists of one or more physical data center buildings with independent power and cooling – ensuring that if one AZ (data center) experiences an outage, others in the region can pick up the load ¹¹. For example, Amazon Web Services’ **us-east-1** region in Northern Virginia contains multiple AZs, each AZ comprising one or several nearby data center facilities ¹². Cloud providers often don’t publicize exact data center addresses, but they use coded naming (region and AZ identifiers) to abstract the physical sites. **Site naming** in enterprise contexts may use codes like DC1, DC2 for a company’s own facilities, whereas cloud uses region/zone codes. The key concept is that “the cloud” is not ethereal – it is a large network of very real data centers, but abstracted for ease of use.

2. Typology and Use-Cases

Not all data centers are alike. They vary widely in scale, ownership, and use-case. Here are common **types of data centers** and their characteristics:

- **Hyperscale Data Centers:** These are extremely large facilities (often 10 MW to 100+ MW of power capacity) designed for a single hyperscale operator’s needs – typically the major cloud providers or social media/tech giants (e.g. AWS, Google, Microsoft, Meta). Hyperscale data centers house *tens of thousands* of servers and are built for maximum efficiency and uniformity at scale. They often span hundreds of thousands of square feet, with massive compute and storage clusters. Hyperscale facilities may be self-owned by the operator or leased as a **wholesale** build-to-suit from a developer, but in either case are dedicated to one primary user (the hyperscaler). These support services like global cloud platforms and large-scale web services. A typical hyperscale site is highly automated and optimized for low cost per unit of computing; for instance, Google’s and Facebook’s self-built data centers use custom server hardware and cooling innovations to achieve high efficiency at scale ¹³ ¹⁴. Hyperscale data centers often forgo some of the ultra-rugged redundancy of enterprise data centers in favor of software resilience and distributed workloads – the applications are architected to tolerate a data center outage, so the facility itself can be streamlined for efficiency.
- **Wholesale Colocation Data Centers:** These are large multi-tenant data centers where the operator (e.g. Digital Realty, QTS, CyrusOne) leases substantial amounts of space and power to a few tenants, each of whom may take a dedicated suite or entire hall. Wholesale colo is “**wholesale**” in the sense of scale – customers typically lease **large footprints or power capacities** (often in the order of hundreds of kW to multi-MW deployments) under multi-year contracts ¹⁵ ¹⁶. The tenant often manages its own IT equipment in the space, and sometimes can customize the room build-out. Wholesale colo providers offer the core infrastructure (power, cooling, fiber connectivity, physical security) and often allow a high degree of customer control within their leased area. This model is popular with large enterprises and cloud providers that need significant capacity quickly without

building their own facility. A **hyperscale** data center can be considered a special case of wholesale – some colo providers offer “hyperscale suites” tailored to cloud giants, providing the scale and economics similar to self-build (Equinix, 2018) ¹⁷. In fact, hyperscalers often lease wholesale space: e.g., a cloud company may lease 10 MW in a multi-tenant campus instead of constructing from scratch, especially for quick expansion in a market (Uptime Institute, 2024) ¹⁸.

- **Retail Colocation Data Centers:** Retail colo facilities (often run by providers like Equinix, CoreSite, TierPoint, etc.) host *many smaller customers* in a shared environment. Clients can rent anything from a **single rack** or a few racks, up to a small cage or suite, with power capacities typically from <1 kW per rack up to a few hundred kW total. The provider supplies a turnkey environment: redundant power feeds to each rack, cooling, fire suppression, security, and on-site support. **Lease models** in retail colo are often shorter-term and charged per rack or per kW, plus fees for **cross-connects** (cables linking to other networks/services in the facility). This model is ideal for businesses that need a reliable data center without the scale to justify a whole building. Retail colos also serve as **interconnection hubs** – because they have many tenants, customers can directly connect to partners, carriers, or cloud on-ramps within the same facility. For example, a company might colocate its servers in an Equinix data center specifically to gain fast cross-connects to multiple Internet backbone carriers and cloud providers. Retail colocation accounts for the majority of deployments by count (over 70% of the colocation market by one estimate) and offers flexibility for customers to start small and scale up (Volico, 2024) ¹⁹.
- **Enterprise Data Centers:** This refers to privately owned and operated facilities dedicated to one organization’s internal IT needs. These can range from small server rooms in office buildings to large Tier III corporate data centers. Many enterprises historically built their own on-premises data centers for full control over security and customization. Enterprise data centers typically support mission-critical business applications, proprietary data storage, etc., and are tailored to the organization’s specific needs (for instance, a finance company’s data center might emphasize very high redundancy and security compliance for transactional systems). However, the trend in recent years is that enterprises are moving away from building new private data centers unless absolutely necessary, due to the cost and complexity – instead opting for colocation or cloud. Still, **legacy enterprise data centers** form a big portion of existing capacity, and some industries (like banking, government, defense) maintain their own facilities for data control or regulatory reasons. Enterprise DCs may also be smaller edge sites near offices, used for low-latency needs or data sovereignty.
- **Edge and Micro Data Centers:** These are small footprint facilities located closer to end-users or devices to reduce latency and network backhaul. **Edge data centers** might be a few racks up to a small container-sized facility, placed in second-tier cities, at base stations, or even on customer premises. They provide caching, content delivery, or compute power near where data is generated/consumed (for example, supporting IoT or 5G applications that need real-time processing). **Micro data centers** can be very small – sometimes a single enclosure with built-in cooling – and can be deployed in remote areas, factory floors, or retail locations. The drivers for edge/micro DCs are latency (each 1000 km distance can add ~10 ms of network delay) and bandwidth efficiency ²⁰. By processing data locally, edge sites avoid the delay and cost of sending everything to a central cloud. While each site is much smaller than a core cloud region, in aggregate they form a distributed layer. Use cases include content delivery networks (CDNs), smart city infrastructure, autonomous vehicle support, and any application where milliseconds matter. These edge sites often rely on standard designs (prefabricated modules) for quick deployment and can be managed remotely.

Cloud Regions, Zones, and On-Prem vs. Colo: It's worth noting how *cloud data centers* fit into this typology. A cloud provider's **region** typically consists of 2 or more large data center sites (each site corresponding to an AZ). Each site in turn may have multiple buildings but is operated as a single zone with independent infrastructure ¹¹. Cloud regions are essentially the hyperscale data centers of the cloud providers, sometimes self-built and sometimes in leased wholesale space. Meanwhile, *on-premises vs. colocation* is a decision many enterprises face: **on-prem** means the company runs its own data center (enterprise model), whereas **colo** means the company places its equipment in a third-party's facility. On-prem can offer more control and perhaps cost advantages at large scale, but colocation offers quicker time-to-market, fewer upfront costs, and the benefit of the provider's expertise and uptime guarantees. Many companies adopt a hybrid approach – keeping some critical systems on-prem or in a private data center, while moving others to colocation and cloud for flexibility.

Typical Densities and Lease Models: Data centers are often described by their **power density** (kW per rack or per square foot) and size (total power capacity or floor area). *Traditional enterprise* racks might be ~3–5 kW, while *modern averages* have crept up to 7–9 kW per rack (Uptime Institute, 2024) ²¹. Hyperscale cloud racks are often in the 5–15 kW range (balanced for efficiency), though specialized AI or HPC racks can draw 20–30 kW or more. Edge and micro data centers might have only a few racks but possibly higher density per rack if space is tight. **Footprints** vary: a wholesale colo tenant might lease a 10,000 sq ft hall for a 2 MW deployment, whereas a retail colo customer could rent a single 42U cabinet (about 2 sq ft of floor space but with cooling and power limits). **Lease models:** Wholesale deals often involve a *power commitment* (e.g. paying for 1 MW whether fully used or not, with a usage-based utility charge) over a long term (5–10 years), and the customer may handle their own IT equipment installation and even fit-out. Retail colo is usually billed per *rack* or per kW on shorter terms (1–3 years typical), with the provider including more hands-on support. Additionally, colo providers charge for **cross-connects** (physical cables linking one tenant to another or to a carrier in the meet-me room) – these are monthly fees per connection that allow access to telecom networks or cloud on-ramps in the facility (Digital Realty, 2025) ²² ²³. Cloud on-ramp connections (like AWS Direct Connect, Azure ExpressRoute) are offered in colocation sites where cloud providers have network nodes; customers can pay for a private virtual circuit into their cloud, which often reduces bandwidth *egress* costs and improves performance compared to using the public Internet.

3. Facility Lifecycle

Building and operating a data center is a complex process. The **lifecycle** of a data center facility spans from planning and construction through ongoing operations. Key phases include:

Site Selection and Entitlements

Choosing the right **site** is critical. Data centers need a location that can reliably support the facility's power, cooling, connectivity, and security requirements for decades. Important site selection factors include:

- **Power Availability:** Access to sufficient electrical power is the top priority. Sites near existing substations or high-voltage transmission lines are ideal to support large loads (tens of megawatts). A robust local power grid with redundancy and *reasonable cost per kWh* is desired. Many data center hubs are in regions with relatively low electricity rates and stable grids (for example, areas with abundant utility generation or cheap wholesale power). Sites might negotiate directly with utilities for dedicated feeder lines or even their own substation. Renewable energy availability is increasingly a factor too (sites in regions with wind, solar, or hydro are attractive for sustainability goals).

Entitlements with the utility – securing commitments for the needed megawatt capacity – must be arranged early, as utility lead times for new circuits or substations can be years.

- **Fiber Connectivity:** Proximity to **fiber optic networks** and Internet backbone nodes is essential for a data center. Ideally, multiple fiber routes run near the site so that diverse network connectivity can be established. Data centers thrive in **connectivity-rich locations** where many telecom carriers are present (e.g. near **Internet Exchange** points or **carrier hotels**). During site selection, developers will check for existing fiber conduits or PoPs (Points of Presence) in the area and may plan to bring in several carriers. Being near an **Internet exchange** or **cloud on-ramp** location can significantly enhance a site's value, as it lowers latency and cost for network traffic ²⁴ ²⁵. If a site is farther from fiber, the project must include extending fiber (or leasing **dark fiber**) to connect the facility to carrier networks, which adds cost and time.
- **Geographic and Environmental Factors:** Locations with low risk of natural disasters are favored. **Geographic stability** – minimal risk of earthquakes, hurricanes, floods, wildfires, etc. – reduces chances of catastrophic outages ²⁶ ²⁷. For example, Northern Virginia has very low seismic activity and modest weather extremes, one reason it became a huge data center hub. Site surveys examine floodplain maps (sites must be outside 100-year flood zones) and historical weather data. **Climate** matters as well: a temperate or cool climate can improve cooling efficiency (allowing use of outside-air economization) ²⁸. Very hot/humid climates mean higher cooling energy usage or more water consumption. That said, many data centers operate in hot regions by using advanced cooling or abundant water. **Altitude** (for example, high-altitude sites) can reduce air cooling effectiveness due to thinner air, which is also considered.
- **Land and Expansion:** Data centers need sufficient land – not only for the initial building but for future expansion phases. Campuses are often developed in phases, so having extra acreage for additional data halls, power equipment, or cooling plants is important ²⁹ ³⁰. The site should support heavy structures and possibly multiple buildings, so geotechnical conditions (soil stability, water table height) are checked. Industrial zoning or special use permits are usually required; the process of obtaining entitlements (zoning approval, environmental impact assessments, construction permits) is a major part of site development. Many jurisdictions have streamlined this for data centers due to economic benefits, but timelines vary.
- **Security and Compliance:** The site should allow for a **secure perimeter** and control of access. Often this means a relatively remote or fenced location where a buffer can be maintained (though urban data centers exist in high-rises too, they rely on building security). Some states or localities have specific building code requirements for mission-critical facilities (e.g. higher seismic design category, fire suppression standards). And local regulatory environment (tax incentives, as discussed later) can sway site choice. For example, states like Virginia, Arizona, and Oregon offer tax breaks for data center investments, making them attractive (VEDP, 2023) ³¹. *Entitlements* phase includes securing all necessary approvals: site plan approval, electrical interconnection agreement, water usage permits (if using city water for cooling), and so on.

In summary, site selection balances **infrastructure** (power and fiber), **environment** (climate and hazard risk), **expansion capability**, and **regulatory climate**. A well-chosen site will have abundant power and connectivity, minimal risk profile, room to grow, and support from local authorities.

Power and Fiber Procurement

After picking a site, the data center developer must procure reliable power and network connectivity for the facility:

- **Utility Power Interconnect:** Data centers require high-capacity electrical feeds. Typically, the facility will connect to the grid at medium-to-high voltage (often 13.2 kV, 13.8 kV, or 34.5 kV lines in the US)³². Large campuses may justify building an on-site **substation** to step down utility transmission voltage (e.g. 115 kV or 230 kV) to medium voltage for distribution within the site. The developer works with the utility on an *interconnection agreement* specifying how many megawatts will be provided and in what timeframe. Redundant utility feeds (from separate substations or transmission lines) are ideal for reliability – many Tier III/IV data centers have two utility feeds (one primary, one secondary) entering the facility. However, redundant feeds depend on utility infrastructure being available. The utility will also have rules about **generator synchronization** or feeding back power; typically data centers are not feeding power to the grid except possibly in demand response situations.

The timeline for utility power can be a critical path. In some booming markets, utilities face backlogs – for instance, Dominion Energy in Virginia reported that connecting large new data centers could take 1–3 years longer than before due to the surge in requests (Dominion, 2023)³³. Developers often have to **front-load** investment: ordering long lead equipment like transformers early, and sometimes constructing portions of the grid connection themselves (then handing over to the utility). Reliable backup power systems (generators) mean data centers can theoretically run without grid power for a limited time, but long-term, inexpensive utility power is essential.

- **Fiber Network Connectivity:** In parallel with power, establishing network connects is crucial. Data center operators typically invite multiple **carrier** companies to build fiber into the facility's **Meet-Me-Room** (see Section 6). During construction, underground conduit paths are laid out to the property line for carriers to pull fiber. If the site is near a major fiber route, procurement might be as simple as arranging a lateral connection with providers. If not, the developer may need to invest in extending fiber cables from the nearest metro fiber intersection (or lease dark fiber from a provider). They often secure agreements with at least two diverse fiber providers for redundancy. High-count fiber cables (hundreds of strands) might be installed to future-proof capacity. The goal is to make the data center a network hub itself, offering tenants a choice of carriers and *low-latency* connections to exchange traffic.

In some cases, large customers (like a cloud provider leasing the site) might bring their own fiber networks or require specific connectivity. **Dark fiber** leases and **IRU (Indefeasible Right of Use)** agreements are common for connecting data centers – these give the operator dedicated fiber strands on long-haul routes to link between facilities or to other network hubs³⁴. For instance, an operator might lease dark fiber between their data center and a downtown carrier hotel to ensure high-bandwidth connectivity for all tenants.

- **Permitting and Construction of Utilities:** Building the power and fiber pathways may involve road boring, erecting new transmission poles, or expanding substations. These activities need permits and coordination with local authorities. Data center projects often include significant *infrastructure upgrades* in the area (sometimes funded partly by the developer), which can benefit the community (e.g. improved grid reliability, more fiber in the ground).

In summary, *procurement* of power and fiber means ensuring the site is wired with enough electricity and network bandwidth to meet current and future needs, with redundancy. These are foundational to delivering the promised uptime and performance to customers.

Design & Construction

Designing a data center is a specialized architecture and engineering task. The design must meet the desired **Tier/redundancy level**, capacity, and operational efficiency targets. Key points in design and build:

- **Phased Construction:** Most large data centers are built in phases or modules. For example, a campus might be master-planned for 4 x 10 MW buildings, but the developer will build the first 10 MW now and leave room for the others as demand grows. Even within a single building, fit-out might be phased: initial build-out of a portion of the hall with power/cooling, and additional **pods** or rooms are completed later. This *modularity* allows aligning capital spend with customer demand. It's common to see data halls commissioned one at a time, each maybe 1-3 MW, in a larger shell building that can hold several halls.
- **Layout:** Typical data center design includes a large one-story building (or multi-story in space-constrained markets) with raised floor or hard floor whitespace for IT racks, adjacent to dedicated **electrical rooms** (for switchgear, UPS systems, batteries) and **mechanical plants** (for chillers, cooling towers, CRAH units, etc.). The site will also have an exterior **generator yard** and possibly a separate **substation yard**. Designers optimize for short pathways (electrical and airflow) and segregating *critical* infrastructure from IT space. Modern designs also consider airflow containment in IT rooms (hot aisle/cold aisle containment) to improve cooling efficiency.
- **Standards:** Many data centers are designed to meet Uptime Institute Tier III or IV criteria, or the TIA-942 standard (which defines Rated-1 through Rated-4 facility levels). The design process includes selecting redundancy configurations (N+1, 2N, etc. – see Section 4) for power and cooling. The building itself is often fortified – e.g., high wind-rated walls if in hurricane-prone areas, seismic bracing if in quake zones, and enhanced fire protection (fire-rated walls separating rooms, very early smoke detection, etc.). Compliance with codes like **NFPA 75** (fire protection for IT equipment rooms) and **NFPA 70 (NEC)** electrical code for safe power distribution is mandatory (Dgtl Infra, 2024) ³⁵. Additionally, sustainability is increasingly baked into design: considering efficient cooling plant (like using water-side economization), and certification goals (some pursue LEED certification for green building).
- **Construction:** Data center construction is fast-paced – many use **prefabricated components** to save time, such as skids for electrical gear or modular block units for IT racks. A specialized *mission-critical* general contractor typically manages the build, coordinating electricians, mechanical contractors, and controls technicians who install the complex systems. Construction must be of high quality to ensure reliability (for instance, precise installation of busbars, proper redundancy wiring, etc.). The timeline for a large data center from groundbreaking to ready-for-service can be around 12-24 months, depending on size and complexity. As soon as one phase is up, often construction will start on the next phase if demand exists.

During late-stage construction, the focus shifts to testing and commissioning of systems (described next). Overall, the design/build phase aims to deliver a facility that meets capacity and **reliability targets** while

optimizing cost and efficiency. For example, a design might choose a particular cooling architecture (air-cooled chillers vs. water cooling) based on climate and water availability to hit a PUE goal, or adopt a flexible **hall design** that can accommodate either low-density or high-density racks as client needs evolve.

Commissioning, Operations, and Service Levels

Before a data center goes live with customer load, it undergoes rigorous **commissioning** to verify all systems work as intended under real conditions:

- **Commissioning Stages:** This typically includes Factory Acceptance Tests (FAT) for major equipment (e.g., testing a generator at the manufacturer), Site Acceptance Tests (SAT) once equipment is installed, and finally **Integrated Systems Testing (IST)**. In IST, all infrastructure – power, cooling, backup systems, controls – is tested together, often through simulated failure scenarios. For example, the commissioning team will cut utility power to simulate a blackout and ensure that UPS batteries carry the load until generators start, that cooling units switch to backup power, etc. They also simulate component failures (like failing one CRAC unit or one UPS module) to ensure the redundancy works (Uptime Institute calls this *concurrently maintainable* if maintenance or failure of one component doesn't disrupt IT load). These tests prove the design meets the desired Tier level performance. Professional commissioning agents often run this process to catch any installation issues or design flaws before customers are installed.
- **Operational Readiness:** Data centers staff must be trained and standard operating procedures (SOPs) and emergency operating procedures (EOPs) put in place. By opening day, the operations team will have methods for monitoring the environment, performing preventive maintenance, responding to alarms, and managing changes safely (often under a change management program to avoid human error outages). They also establish **Service Level Objectives (SLOs)** and **Indicators (SLIs)** – metrics like site availability, cooling setpoints, network latency, etc., which underpin customer **SLAs (Service Level Agreements)**.
- **Service Levels:** Most colocation data centers offer contractual uptime commitments (commonly 99.999% or “five-nines” availability for power and cooling to the rack). This equates to at most a few minutes of downtime per year and is usually aligned with Tier III or IV design. Meeting these high service levels requires robust **maintenance regimes** (e.g., regular generator testing, UPS battery checks, cooling preventive maintenance) without impacting live load. Many providers schedule maintenance during off-peak hours and use redundant components so that maintenance can happen on one system while another carries the load (this concept is *concurrent maintainability*).
- **Operations Practices:** Once live, data centers are run 24/7 by on-site staff including facilities engineers, security, and often remote hands technicians who can assist customers with their equipment. Monitoring systems (DCIM – Data Center Infrastructure Management software) track power usage, temperature, humidity, and alert staff to any anomalies. **Failure domain management** is an important concept: operators ensure that any single failure (a broken utility feed, a leak in a cooling pipe) remains isolated in its domain and doesn't cascade to full outage – this is achieved by clear physical segregation and robust emergency response procedures.
- **SLI/SLO Thinking:** Operators define SLIs like *power delivery continuity* or *cooling effectiveness* and set SLOs (targets) such as “temperature in cold aisles will remain within ASHRAE recommended range

99.9% of the time" or "packet loss on cross-connects < 0.1%". These are often internal targets, but they guide operational focus and continuous improvement. If any incident occurs (e.g. a momentary power slip causing IT reboot), a full root-cause analysis is performed and shared with customers.

In essence, after construction, the data center enters a long operational phase where the goal is **steady, predictable service**. Well-run data centers often achieve *years* of continuous uptime. Industry best practices for operations are codified in standards like **ISO 27001** (information security including physical security) and **SOC 2** audits (which check controls for security and availability) – many providers undergo audits to demonstrate good operational governance (Dgtl Infra, 2024) ³⁶ ³⁷.

The commissioning and operations phase is where design meets reality. A robust commissioning ensures that from day one, the facility can handle real-world stress. Then operational discipline keeps it running within specifications. This lifecycle repeats if expansions occur – each new phase goes through commissioning and then merges into ongoing operations.

4. Power Architecture

Power is the lifeblood of a data center. The **power architecture** encompasses how electricity is brought from the utility down to the servers, including all the conversion, backup, and distribution components. Key concepts include power capacity (MW or kW), density, redundancy levels, and the types of power backup systems employed.

Capacity and Density: Data center capacity is usually described in **megawatts (MW)**, which refers to the IT load supported (not counting cooling and losses). A 10 MW data center can power roughly 10,000 kW worth of running servers (plus overhead for cooling). Large hyperscale campuses can exceed 50–100 MW of IT load. Smaller enterprise or edge facilities might be under 1 MW. **Kilowatts (kW)** are used at the rack or small deployment level. For instance, a rack drawing 5 kW is fairly typical today; some high-density racks draw 15–30 kW or more. Power *density* can be expressed per rack (kW/rack) or per area (e.g., watts per square foot of white space). Traditional designs were ~50–100 W/sq ft, but modern designs easily exceed 150–200 W/sq ft due to packing more kW per rack. Data centers must ensure the power distribution and cooling can handle the peak density in any given area. Average server rack densities remain below 8 kW in many facilities, but a growing minority of data centers are deploying racks above 15 or 20 kW to support dense computing needs (Uptime Institute, 2024) ³⁸ ³⁹.

Block Sizing and Distribution: Large data centers often break the electrical architecture into **blocks** or **modules** – for example, a 20 MW facility might be divided into four modules of 5 MW each, each module having its own UPS systems, generators, and distribution panels. This modular approach localizes any fault and aligns with phasing (you can build one block at a time). Incoming utility feeders at medium voltage feed these blocks via transformers. Within each block, power is stepped down to *low voltage* (e.g. 415/240V AC or 480/277V AC in US systems) for distribution to server racks (often through Power Distribution Units, or PDUs, on the floor).

Redundancy Patterns: To ensure continuous power, data centers employ redundancy at various levels. Common redundancy topologies include:

- **N (no redundancy):** The capacity exactly meets the load requirement with no spare. For example, if you need 1000 kW, you have one 1000 kW UPS module (N). This maximizes efficiency (all capacity is used) but if that component fails, there's no backup – *not acceptable* for critical loads generally.
- **N+1:** You have one extra unit beyond N. For instance, if the load is 1000 kW and N is handled by 2 x 500 kW UPS modules (total 1000 kW), an N+1 system might have a third 500 kW module as spare. So you could lose one UPS and still meet full load. Similarly, an engine generator farm might have, say, 5 generators each capable of 2 MW for a 8 MW load (which is N=4, plus 1 spare). N+1 protects against single component failure and allows maintenance on one unit at a time. It's a common design for Tier III data centers ⁴⁰ ⁴¹. **Pros:** more economical than 2N (less extra capacity), while covering a single failure. **Cons:** if two failures occur or a failure happens during maintenance, load is lost.
- **2N (or N+N):** Two independent sets of infrastructure each capable of handling the full load (100% redundancy). For example, you have two UPS systems each rated for 1000 kW for a 1000 kW need. Normally each may run at 50% load, but if one fails, the other can instantly take the full load. Similarly, 2N generator design means every critical load has two generator feeds from separate generator sets. This is typical of Tier IV design – completely fault-tolerant systems ⁴² ⁴³. **Pros:** can tolerate any single failure without outage and often allows concurrent maintenance (you can service one side while the other carries load). **Cons:** highest cost (duplicate infrastructure) and lower normal utilization of components (less efficient if both sides are rarely fully loaded).
- **2N+1 (or 2(N+1)):** This is a rarity but denotes a 2N system where each side also has spare capacity. For example, 2N with each side having N+1 units. It's an extreme design for maximum fault tolerance – even during maintenance or a failure on one side, the system can handle an additional failure. Tier IV data centers sometimes approximate this (e.g., two utility feeds plus backup generators that are N+1) ⁴³ ⁴⁴. **Pros:** withstands multiple failures. **Cons:** very expensive and usually only used for ultra-critical apps.

There are also variants like **Distributed Redundancy** (e.g., 2N/2 = each load has half capacity from two systems, or block redundant designs like 3N/2, etc.) to optimize for cost while providing some failover. For example, a **catcher system** might have multiple modules where one spare can “catch” the load of any failing module via a bus tie arrangement. The trade-offs revolve around cost, complexity, and *risk tolerance*. Many enterprise and colo sites opt for N+1 as a good balance (concurrent maintenance and single fault tolerance), whereas some financial or military data centers might insist on 2N for every component for highest uptime. Uptime Institute Tier standards roughly correlate: Tier III requires N+1 for continuous operations and **concurrent maintainability**, and Tier IV requires 2N for **fault tolerance** ⁴⁰ ⁴².

Power Chain Components: The main elements in the power chain and their typical arrangements:

- **Utility Feeders & Transformers:** High or medium voltage utility power is transformed down to facility distribution voltage (e.g., 13.8 kV to 480 V). Often an **A-side and B-side** feeder for redundancy feed separate transformer banks (in 2N systems). In N+1, one utility feed and transformer may suffice if backup gens are present.

- **Switchgear and ATS:** Switchgear panels distribute power to large loads. An **Automatic Transfer Switch (ATS)** or *Static Transfer Switch* is used to switch between utility and generator supply. In many designs, an ATS is at each generator or each UPS, so that when utility fails, the ATS switches over to generator feed seamlessly (within a few cycles). Some designs use static (solid-state) switches for faster transfers or to maintain dual feeds to critical loads.
- **UPS (Uninterruptible Power Supply):** The UPS systems provide battery-backed power to bridge the gap between utility outage and generator start, and also condition the power. In large data centers, **double-conversion (online) UPS** units are common – these take AC input, convert to DC, charge batteries and invert back to AC, isolating the output from input disturbances. Double-conversion UPS provide very clean power and near-instant response to outages (the batteries supply power immediately on input loss) ⁴⁵ ⁴⁶. Another type is **line-interactive UPS**, more common in smaller installations, which normally pass through utility and use power electronics to buck/boost minor voltage changes, switching to battery inverter if the input fails. Large facilities might also use **rotary UPS** (DRUPS – Diesel Rotary UPS) which integrate a flywheel or motor-generator that provides ride-through until the diesel engine engages – these can eliminate the need for large battery banks. The UPS topology chosen affects efficiency and footprint. Modern static UPS can reach >95% efficiency in double-conversion mode. Lithium-ion battery adoption (discussed below) is also changing UPS designs.
- **PDUs (Power Distribution Units):** After UPS, power is distributed to server racks via PDUs. In a traditional setup, a PDU is a large cabinet that takes, say, 480 V from UPS and transforms it to 208/120 V for distribution within a room, with breaker panels feeding groups of racks. Many modern data centers now distribute at 415/240 V three-phase all the way to the rack (no need for a PDU transformer) to reduce losses – servers worldwide can typically accept 240 V AC supply. Rack-level **rack PDUs** (basically power strips with circuit breakers) then distribute to individual servers. Some designs use **busway** overhead systems, where bus bars run above racks and tap-off boxes feed racks, providing flexibility in power provisioning.
- **Remote Power Panels (RPPs):** These are smaller distribution panels, often used in large rooms, that take feeder circuits from a PDU and allow branch circuits to server racks in that vicinity ⁴⁷. They are basically a way to extend distribution closer to loads for manageability.

UPS and Battery Technologies: Historically, data centers used **Valve-Regulated Lead-Acid (VRLA)** batteries for UPS energy storage. VRLAs are like car batteries, proven but heavy and needing replacement every 3-5 years. Today, **Lithium-ion (Li-ion)** batteries are increasingly popular for UPS systems. Li-ion batteries offer several advantages: higher energy density (so they take up 50-70% less space and weight), longer lifespan (8-15 years, roughly triple VRLA life), and tolerance of higher temperatures which can reduce cooling needs for battery rooms (CoreSite, 2022) ⁴⁸ ⁴⁹. They also recharge faster (a Li-ion can recharge in ~2 hours vs 8-12 hours for VRLA) ⁴⁸. The main downside was cost, but Li-ion prices have dropped ~97% since 1990s and continue to fall ⁵⁰. By 2020, about 10% of data center UPS deployments used Li-ion; by 2025, it's predicted to reach ~35% as more operators make the switch (BloombergNEF via CoreSite, 2021) ⁵¹. Li-ion's better cycle performance and slower capacity degradation mean fewer battery replacements and potentially lower total cost of ownership despite higher upfront cost ⁵² ⁵³.

VRLA batteries remain common in older facilities. They are typically kept in dedicated battery rooms with strings of battery cabinets providing the DC ride-through time (often sized for ~5-10 minutes of backup,

though generators usually start in <1 minute). *Maintenance* of lead-acid batteries is significant (monitoring, testing, periodic replacement). Li-ion systems include battery management electronics at the cell/module level for safety (to prevent thermal runaway) ⁵⁴ ⁵⁵. Data center operators have been cautious with Li-ion due to safety concerns, but experience shows that with proper systems, Li-ion can be very safe and reliable (Uptime Institute, 2019) ⁵⁶ ⁵⁷. Many new-build data centers now opt for Li-ion UPS by default for space savings and longer life (Vertiv, 2023) – for example, a Li-ion battery might last the entire 15-year life of the UPS without replacement, whereas VRLA would be swapped 3 times in that span ⁵⁸.

Other energy storage technologies sometimes used: **Flywheels** (provide a few seconds of ride-through – often integrated with rotary UPS), and emerging options like **Nickel-Zinc batteries** or even **supercapacitors** for short bridging (less common). Some hyperscalers experiment with **grid battery systems** at large scale (e.g., using onsite utility-scale Li-ion batteries instead of generators for backup, as Microsoft has piloted with megawatt-scale batteries doubling for grid services). But diesel generators remain the dominant backup for longer outages.

Generators and Backup Power: Data centers nearly always have on-site **generator** sets, usually diesel engine generators, to supply power during extended utility outages. Generators are sized to handle the full IT load plus supporting infrastructure (cooling, etc.). They typically start automatically within seconds of power loss. A standard configuration might be one generator per power block or per UPS module – e.g., a 2 MW diesel gen for each 1.5 MW of UPS-protected load (N+1 or N+2 on generators). Generators can take around 10 seconds to come online and synchronize; the UPS bridge covers this gap ⁴⁵ ⁵⁹. Fuel storage on-site is critical – large facilities often have 24–48 hours worth of diesel fuel stored, with contracts for refueling in case of prolonged outages. For example, a data center with 10 x 2MW generators might have 50,000+ gallons of diesel stored to sustain full load for a day or more.

Run-time and Fuel Logistics: *Runtime* refers to how long the data center can run on generators before refueling. Many jurisdictions require a minimum (e.g., hospitals require 72 hours fuel for emergency power). Data centers typically design for at least 12–24 hours at full load, recognizing that in widespread outages they will need refueling. They establish fuel delivery contracts with multiple suppliers and test generators under load regularly (often monthly). Fuel quality maintenance (polishing, filtering) is also done to ensure diesel doesn't degrade in storage. Diesel generators can often run continuously for days or weeks if refueled, though maintenance (oil, filters) would eventually be needed after perhaps 72+ hours continuous.

Redundancy and Generator Configuration: Generators can be in N+1 or 2N configurations too. A common approach: N+1 generator farm, where one extra genset is available. If a site has eight generators for N and one fails, the ninth (spare) can cover. In a 2N utility design, generators might be tied to each utility path – effectively two separate generator banks each capable of the full load (which is a very robust but costly setup). Some data centers also cross-connect generator outputs so they can support each other in emergencies (with careful control systems to avoid overloading). High-tier sites test *full load transfer* regularly to validate that the entire facility can be carried by generators for an extended time.

Emissions and Sustainability: Diesel generators emit CO₂, NOx, particulate matter, etc., so environmental regulations limit their use. In many places, generators are permitted as **emergency standby** only, with a cap on non-emergency run hours (for example, air quality permits might allow ~50 hours/year of testing). Data center operators strive to minimize generator runtime for environmental and community reasons. In states like California, *Tier 4 emissions* standards or retrofits (advanced exhaust after-treatment) may be

required on new generators to cut pollution. There's also interest in greener backup solutions: some operators use **biogas** or **natural gas generators** (which have cleaner emissions but start slower and are less common for full UPS backup). Recently, there have been trials of **hydrogen fuel cells** to replace diesel gensets – Microsoft notably tested large hydrogen fuel cells for backup, emitting only water. While promising, fuel cells are not yet widely adopted for this purpose.

Another sustainable practice is limited use of generators in **grid support** (with permission): e.g., in some regions, data centers participate in demand response – they might intentionally switch to generator during peak grid demand to reduce strain on the public grid. However, this raises emissions concerns and regulatory scrutiny (generators used this way may be treated as power plants).

Distribution Paths and Fault Tolerance: Beyond main components, the electrical topology ensures power paths to each critical load are redundant. Many data centers provide dual power feeds to each rack (A and B feed), each from independent UPS/generator systems. Servers with dual power supplies can connect to both feeds, achieving redundancy at the IT device level. Thus even if one entire power train (UPS, PDU, etc.) fails or is under maintenance, the server stays powered from the other feed. This dual-corded approach is standard for enterprise and cloud hardware. For single-cord devices, sometimes automatic transfer switches at rack level are used, but dual-cord is preferred.

In summary, the power architecture is designed such that *no single point of failure* in the power supply chain will take down the IT load (for Tier III/IV objectives). Achieving this means redundant units (N+1 or more) and often two completely independent power streams (2N) running to the IT equipment. This architecture, combined with careful maintenance and testing, yields the high availability that enterprise customers expect – e.g., a Tier IV data center targets only about 26 minutes of downtime **per year** for power/cooling (99.995% uptime) ⁴² ⁴¹.

5. Cooling Architecture

Cooling is the other half of the equation in keeping a data center running. Servers produce heat that must be continuously removed to prevent overheating. The **cooling architecture** encompasses the methods and equipment used to keep IT equipment at safe operating temperatures. Key elements include the type of cooling systems (air vs. liquid, chilled water vs. direct expansion), efficiency measures like economization, and metrics like PUE and WUE that gauge effectiveness.

Air Cooling Basics: Traditionally, data centers use air as the cooling medium. Rows of server racks draw in cool air (at the front) and expel hot air (at the back). The facility's cooling system ensures there is a constant supply of sufficiently cool air to the front of racks and that hot exhaust air is removed and cooled down before recirculation. Two primary pieces of air-cooling equipment are:

- **CRAC (Computer Room Air Conditioner) units:** These are essentially air conditioners that cool air using a refrigeration cycle (compressor, evaporator coil). They often blow cold air under a raised floor or directly into cold aisles. **DX (Direct Expansion)** CRACs have refrigerant in the coils that directly cool the air (like a typical AC unit).
- **CRAH (Computer Room Air Handler) units:** Instead of having a built-in compressor, CRAHs use chilled water piped in from a central chiller plant to cool the air via a cooling coil. CRAHs have fans

and a chilled water coil – the heat from the air transfers into the chilled water, which then goes back to the chiller system to dump heat externally.

The choice of CRAC vs CRAH depends on whether a **chilled water** system is used. *Chilled water* plants have large chillers (like big industrial AC systems) that cool water which is circulated to CRAHs in the server rooms. *DX systems* don't need central chillers; each CRAC unit handles cooling with refrigerant and typically rejects heat to the outside via a condenser. Chilled water is common in large data centers because it scales well and can be more efficient especially if using economization. DX CRACs are simpler for small facilities (like a single room can have a few CRAC units and condensers outside).

Free Cooling and Economization: A big portion of data center energy is used in cooling. **Economization** refers to using favorable outside conditions to cool the facility with reduced chiller/compressor use. There are two main types:

- **Air-side Economization:** Bringing cool outside air into the data center to directly cool the space (and exhausting hot air out). For example, on a cool day, instead of running chillers, you can simply use fans to pull in outside air through filters and push it across the servers, as long as temperature and humidity are in acceptable ranges. Air-side economizers often have *dampers* and louvers that modulate outside vs return air. Some designs can operate in economizer mode a significant portion of the year (e.g., in temperate climates, maybe 50–70% of the time no mechanical cooling is needed). The challenge is filtering and controlling humidity of outside air, and avoiding introducing contaminants.
- **Water-side Economization:** Using cooling towers or dry coolers to reject heat without running chillers when ambient conditions permit. For instance, **cooling towers** can cool water by evaporation when outside air is cool enough, producing chilled water (or at least cooler water) for CRAHs without the chiller running. This often kicks in during cooler seasons or at night. Some systems use **dry coolers** (like giant radiators with fans) to cool water when air temps are low, avoiding water use. The net effect: reduce chiller use and save energy when possible.

Adiabatic and Evaporative Cooling: The term *adiabatic* cooling in data centers usually refers to using evaporation of water to enhance cooling. This can be part of air-side or water-side systems. For example, some modern air handling units use **adiabatic misting** or evaporative media – they spray water into the incoming outside air, which evaporates and cools the air (taking advantage of dry climates' capacity to absorb moisture). This significantly boosts cooling capacity with minimal mechanical energy, but uses water. Evaporative cooling works best in dry climates (like Phoenix or Denver) where the wet-bulb temperature is much lower than dry-bulb. It's less effective in humid places. Many hyper-efficient data centers (like Google's in certain regions) use evaporative cooling as the primary method, achieving very low PUEs but at the cost of water usage. **WUE (Water Usage Effectiveness)** is a metric to track this trade-off (see below).

In sum, **modern cooling architecture** often combines multiple modes: mechanical chilling when needed, and free cooling or evaporative assist when conditions allow. The control systems switch between modes to maintain target server inlet temperatures efficiently (Equinix, 2024)⁶⁰ ⁶¹. For instance, a facility might run in air economizer mode at winter, add evaporative cooling in spring/fall, and use full chiller operation on hot summer days – balancing water vs energy use depending on what's optimal (Equinix, 2024)⁶² ⁶³.

Liquid Cooling (Direct-to-Chip, Immersion, etc.): As rack densities rise with high-performance computing (HPC and AI workloads), traditional air cooling becomes challenging beyond ~20–30 kW per rack. Liquid cooling, which involves bringing a liquid (usually water or a dielectric fluid) in contact with the heat sources, offers much higher heat removal capacity. There are a few approaches:

- **Rear-Door Heat Exchangers:** These are passive cooler doors on the back of server racks. They have water-cooled coils that absorb the hot air as it leaves the rack, removing most heat before it enters the room. This can allow racks up to perhaps 30–50 kW without overheating the room, as most heat is transferred to water in the rear door. It's a retrofit-friendly solution – servers still use air internally, but the rack door cools that air effectively.
- **Direct-to-Chip (Cold Plate) Cooling:** Here, cold plates (small water blocks) are attached directly to the CPUs/GPUs and other hot components inside servers. A piping system circulates coolant (often water with additives) to each server node, picking up heat straight from the source and carrying it to a CDU (cooling distribution unit) and then to building heat exchangers. This is very efficient since water's thermal conductivity far exceeds air. Systems like this can cool extremely high power chips (like 400W+ CPUs, GPUs) that air could not. Some designs still have air cooling for secondary components, others are fully liquid for all high-heat parts. Data centers using direct liquid typically run a warm-water loop (the water can be, say, 40°C supply) because even warm water absorbs lots of heat from chips that run at ~70–80°C. Warm liquid cooling enables potential **heat reuse** (since the coolant comes out hot, it can be used for heating nearby buildings, etc., more easily than warm air could).
- **Immersion Cooling:** This cutting-edge approach places servers (with minimal modification) into baths of dielectric fluid (non-conductive coolant). The fluid directly touches all components and carries away heat. There is single-phase immersion (fluid stays liquid and pumped to external heat exchanger) and two-phase (fluid boils on hot components, the vapor condenses on a cooled lid and drips back). Immersion can handle extremely high densities (100 kW+ per tank easily) and can be very quiet (no fans needed). It's being explored for crypto mining and HPC clusters. However, it requires specialized tank enclosures and maintenance can be messy (need to pull servers out of liquid for service).

Liquid cooling adoption is still limited but growing as AI/ML drive up densities. In 2023, fewer than half of operators had any liquid cooling deployed, but many plan to incorporate it for future high-density requirements (Uptime Institute, 2024) – the reliability of liquid cooling at large scale is still being proven ⁶⁴. By using liquid cooling, data centers can keep PUE in check even as per-rack power skyrockets, and importantly, avoid the need to blast extremely cold air which is inefficient. Some hyperscalers have announced plans for direct liquid cooling pods specifically for AI workloads.

Key Metrics: PUE and WUE: Data center infrastructure efficiency is commonly measured by **Power Usage Effectiveness (PUE)** – the ratio of total facility power to IT equipment power. PUE = 1.0 would mean all power goes to IT equipment (no overhead), which is ideal but not attainable fully. The closer to 1, the better (meaning minimal power lost to cooling, lighting, UPS losses, etc.). In practice, the *industry average PUE* has flattened around ~1.55–1.6 in recent years (Uptime Institute, 2024) ⁶⁵. This means roughly 0.55 kW overhead per 1 kW IT load on average. Best-in-class new data centers, especially large hyperscales in cool climates, can achieve PUE ~1.1–1.3 (very little overhead) ⁶⁶. Older or smaller sites might have PUE of 2.0 or higher (inefficient cooling or low utilization). Uptime data shows that while new designs are very efficient,

the global fleet average improves slowly because of legacy facilities ⁶⁷ ⁶⁸. Operators monitor PUE as a gauge of cooling and power chain efficiency; improvements come from things like using economization, raising server inlet temperatures, using more efficient components (UPS, fans, etc.), and higher utilization of capacity. One caution: PUE doesn't account for *IT efficiency* (efficient servers) – it's purely an infrastructure metric ⁶⁹.

Water Usage Effectiveness (WUE) has emerged to measure water efficiency. WUE = liters of water used per kWh of IT compute. Using evaporative cooling can significantly cut electrical PUE but at the cost of water. For example, a hyperscale data center might have a WUE around 0.2 L/kWh (meaning for every 1 kWh of IT, 0.2 liters of water evaporated for cooling) – Amazon reported ~0.19 L/kWh global average for its data centers (2022), a 24% improvement over the previous year by optimization (AWS, 2023) ⁷⁰. A facility using only air cooling (no water) would have WUE = 0 but likely higher PUE. Many operators choose based on local conditions: in water-scarce regions, they avoid high water use and accept a bit higher energy use; in cooler or water-rich regions, using water can be okay to slash energy use (Equinix, 2024) ⁷⁰ ⁷¹. This is a critical design trade-off: **water vs energy**. Regulations are starting to address it – e.g., some areas considering limits on data center water use due to drought risks, while others might impose efficiency standards that indirectly push water use. Large data centers can consume **millions of gallons per day** if using evaporative cooling heavily. One report noted some U.S. hyperscale facilities may each use up to 5 million gallons of water per day for cooling during peak, equivalent to a town's water use (EESI, 2025) ⁷² ⁷³. Nationwide, U.S. data centers were estimated to consume about 449 million gallons per day (1.7 billion liters) of water as of 2021 ⁷⁴. This has raised concerns, especially as AI loads increase both energy and water demands ⁷⁵ ⁷⁶. To mitigate this, some data centers use non-potable water sources (reclaimed wastewater) for cooling, and implement recycling (running water through multiple cycles in cooling towers). Others shift to **air-cooled chillers** or **liquid cooling** that use less or no water. The Equinix example shows WUE can range widely: an evaporative-cooled site could have WUE of 1.5–2.0 (if using a lot of water in hot, dry climate) while an air-cooled site has WUE ~0 but higher PUE ⁷¹. Best practice now is to consider *both* PUE and WUE together to truly gauge sustainability ⁷⁷.

Design for AI/HPC: *AI and high-performance computing* trends are affecting cooling design significantly. These deployments have a few characteristics: very high rack power (sometimes 30–60 kW per rack of GPUs), high heat flux in small areas, and often more tolerance for warmer water cooling (some GPU clusters can be water-cooled). For such loads, many new facilities or retrofits are enabling **liquid cooling loops**. For example, some colocation providers are adding liquid cooling options like coolant distribution units to support customers with dense AI training pods. Cooling for AI also has to consider *hot spots* (local concentrated heat). Traditional cooling might not evenly cool a 50 kW rack – liquid is better here. Additionally, the energy needed to cool such dense loads can push up PUE if done with air alone (because you'd need very cold air or very high airflow). Using liquid cooling can keep PUE lower by reducing air cooling needs and even allow higher coolant temperatures (since CPUs/GPUs can run warmer with liquid). Interestingly, higher coolant temperatures mean you can reject heat without compressors (wider economizer window) or even reuse heat – for instance, 60°C water coming off an AI cluster could be used for district heating of nearby buildings.

Thermal guidelines: The industry follows guidelines like **ASHRAE TC 9.9** recommendations for data center temperature and humidity. ASHRAE's *recommended* envelope for most data center IT is about 18°C to 27°C (64°F to 80°F) inlet temperature ⁷⁸ ⁷⁹, and relative humidity 40–60% (with allowable ranges that are wider) ⁸⁰ ⁸¹. Modern servers can run fine at 27°C inlet or higher, which has allowed many facilities to raise thermostat setpoints and rely more on economization, cutting cooling energy. Every 1°C increase in server

inlet can save significant energy on cooling. Many data centers now operate in the mid-70s°F inlet range rather than blasting cold air at 60°F as was common 15 years ago. ASHRAE also defined classes A1 through A4 for environments, where A4 allows up to 45°C (113°F) inlet in emergencies⁸¹ – some telecom gear or specially designed kit can survive that, but typical enterprise IT stays within A1/A2 recommended range.

Containment: One technique widely adopted is **hot aisle or cold aisle containment**. By using physical barriers (panels, doors), mixing of hot and cold air is minimized. For example, hot aisle containment encloses the hot aisles and ducts hot air directly back to AC units or ceiling return, preventing it from spilling into cold aisles. This way, the cold air supplied can all go to server intakes without being diluted by hot exhaust, allowing higher supply temperatures and thus more efficient cooling. Containment can significantly improve cooling effectiveness and reduce energy waste.

In conclusion, data center cooling architecture is a multi-faceted design balancing act. It must ensure reliable removal of heat under worst-case conditions (full load on a hot day), yet be efficient by taking advantage of outside conditions and advanced cooling methods. **Power Usage Effectiveness (PUE)** remains a primary measure – with state-of-the-art facilities achieving PUE ~1.2 or better via techniques like free cooling – but **Water Usage Effectiveness (WUE)** is now also front-and-center to ensure water resources are managed responsibly⁶⁰ ⁷¹. With emerging high-density workloads, **liquid cooling** is transitioning from niche to mainstream in certain environments, heralding a shift in data center design akin to the shift in performance the workloads demand. The next generation of data centers, especially for AI, will likely blend air and liquid cooling and push the boundaries of heat reuse and sustainable operations.

6. Network & Interconnection

Beyond power and cooling, a data center's value is deeply tied to its **network connectivity**. Data centers are nodes in the vast fabric of the internet and private networks. The design of network and interconnection facilities within a data center determines how easily and quickly data can flow between the hosted systems and the outside world (or between customers within the facility). Key concepts include meet-me rooms, cross-connects, carrier-neutral facilities, Internet Exchanges, and latency considerations.

Meet-Me-Rooms (MMRs): A *meet-me-room* is a dedicated space in a data center where telecommunications carriers and customers interconnect their network cables and equipment. Essentially, it's the physical nexus of network connectivity in the facility. Multiple carriers will have fiber terminating in the MMR, and customers can run cross-connects to those carriers to obtain internet or WAN services. The MMR is secure and managed by the facility operator – it's the place where one network *meets* another. According to one definition, an MMR is "a physical location within a data center where multiple telecommunications and network service providers interconnect their equipment to facilitate exchange of data and traffic" (Sunbird, 2021)⁸². By centralizing connections in an MMR, the data center provides a *neutral* playing field: any customer can connect to any carrier present, usually via short fiber patch cords. MMRs often have overhead cable trays or under-floor routes where dozens or hundreds of cross-connects run. In large data centers, there may be two MMRs for redundancy (diverse paths so if one room is compromised, networks can switch to the other).

Carrier Hotels and Connectivity Hubs: Historically, some buildings in major cities became **carrier hotels** – colocation sites packed with network service providers interconnecting their networks. Examples: 60 Hudson Street in New York, One Wilshire in Los Angeles. These sites have enormous network density. Modern data centers often collaborate with or emulate carrier hotels by attracting many carriers. A *carrier-*

neutral data center means it hosts multiple carriers and doesn't favor one – customers have choice. **Interconnection** is a major draw for retail colocation: companies colocate specifically to gain access to a marketplace of providers and partners. Facilities like Equinix IBX centers are known for rich ecosystems of networks, clouds, and enterprises all cross-connecting (Equinix, 2024) ⁶⁰.

Cross-Connects: A cross-connect is simply a physical cable linking two parties' equipment within the data center. For instance, if a bank's router in Rack A needs to connect to an ISP's router in the MMR, the data center will run a fiber patch or copper cable from the bank's patch panel to the ISP's port in the MMR – this is a cross-connect. Cross-connects can be fiber (most common for network links), copper Ethernet, or even coax for certain legacy telecom. They are typically low-cost to install but providers charge a monthly fee per cross-connect (often ~\$100-300/month each, depending on market). Cross-connects are delivered within days and are much faster and cheaper than connecting at an external site via local loop. Having multiple carriers in one building reduces latency and cost because you're not going over long external circuits – you directly patch to the provider. Digital Realty notes that by interconnecting within the meet-me-room, customers can streamline WAN performance and lower costs (Digital Realty, 2025) ⁸³. Cross-connects inside a facility are usually *redundant* if needed (customers might run two cross-connects via diverse paths for redundancy).

Internet Exchange (IX) Fabrics: Many data centers also host **Internet Exchange** points. An IX is essentially an ethernet switch fabric where dozens or hundreds of networks meet to freely exchange (peer) traffic. Instead of every network buying transit from telecom carriers, they can directly swap traffic with each other at an IX, improving performance and reducing cost. Data centers in major hubs often have one or more IX providers (e.g., Equinix Internet Exchange, DE-CIX, LINX) present. To join, a customer orders a cross-connect to the IX switch and establishes BGP peering with other participants. This is important for content providers, ISPs, cloud providers, etc. Some IXs span multiple data centers within a metro via trunked fiber links, so location is less critical, but often the largest participant density is in just a few facilities.

Cloud On-Ramps: Cloud on-ramps are private connections into cloud providers' networks offered at colocation sites. For example, AWS Direct Connect, Microsoft Azure ExpressRoute, Google Cloud Interconnect – these services allow a customer to connect directly from their colocation environment to their cloud virtual network, bypassing the public internet. The result is lower latency, higher reliability, and often significantly lower bandwidth charges (cloud providers reduce **egress fees** for traffic sent via a direct connect vs. over the Internet). These on-ramps are set up via either a carrier (who has connectivity into the cloud) or via the cloud provider's router in the facility. Colocation data centers advertise cloud connectivity ecosystems because enterprises running hybrid architectures want dedicated links to their cloud resources. This has been a big driver of interconnection growth in recent years – companies will lease small colo footprints just to be near cloud on-ramp nodes and avoid paying hefty cloud egress costs and suffering internet latency. Cross-connects to cloud on-ramps or using an exchange like Megaport or Equinix Fabric provide these links in a flexible way.

Latency Considerations: *Latency* is the time delay for data to travel between points, very important for real-time applications. Physical distance is a major factor (speed of light in fiber is ~5 microseconds per km, ~0.5 ms per 100 km). So, placing servers in a region closer to end-users reduces latency. This is why we have East Coast vs West Coast data centers, etc. Financial trading firms may colocate in the same facility as an exchange matching engine to get sub-millisecond latency. Similarly, web services often distribute across regions to serve users faster. Within a data center, latency is negligible (a cross-connect might be a few microseconds). The key is connecting to networks that provide fast routes to end-users.

By having a presence in a well-connected data center, a company can reach *major backbone nodes* with minimal hops. For example, a server in Ashburn, Virginia (the world's densest internet hub) can reach most US and European Tier 1 networks in 1–2 ms. If that server were in a remote area, it might incur tens of milliseconds just to get to a major hub. Therefore, data center markets like Northern Virginia (NoVA), New York/New Jersey, London, etc., became popular partly due to network centrality (each additional network that comes improves connectivity options for all, a network effect).

Metro vs. Long-Haul Connectivity: Data centers often distinguish between **metro connectivity** (local fiber rings connecting facilities within a city/region) and **long-haul connectivity** (connections between cities via terrestrial or submarine cables). Big data center providers offer *metro interconnect* services (dark fiber or wavelengths) between their facilities so customers can distribute infrastructure across multiple sites but link them as if local. Long-haul connections typically come via carriers who land in the meet-me-rooms. If a data center is near a **submarine cable landing station**, it might become an international hub (e.g., LA, New York, Marseille for Europe, Singapore). Those with long-haul fiber routes on-site attract hyperscalers for use as network hubs.

Network Redundancy: Just as power has redundancy, network links do too. Data center operators often ensure diverse fiber entry paths into the building (two or more entry points, so a fiber cut in one spot won't isolate the facility). Carriers will also run their fibers in diverse routes. Tenants can contract multiple carriers or routes to ensure resilient connectivity. Some advanced data centers have an "A-side" and "B-side" meet-me-room with completely separate fiber routes, akin to dual power feeds.

Physical Security of Networks: MMRs are high-security areas since all customers' circuits converge there. Only authorized personnel and escorted vendor technicians can enter. Cross-connects are labeled and tracked to avoid mix-ups. Modern data centers use automated infrastructure management to track every patch and connection.

In summary, the network and interconnection capabilities of a data center define how valuable it is for communication-intensive deployments. Retail colo providers differentiate themselves by the richness of their ecosystems – e.g., Equinix's success is largely due to having *over 460,000 interconnections* among its customers globally (Equinix, 2024) ⁸⁴ ⁸⁵. For a business reader: locating your infrastructure in a well-connected data center can drastically improve user experience (through low latency) and reduce network costs (through local peering and direct cloud links). It provides the flexibility to quickly connect with partners and scale connectivity as needed, in contrast to on-prem data centers where you might have to order telco circuits that take months. As digital businesses rely on multi-cloud and real-time data exchange, interconnection has become a key selling point of modern colocation.

7. Reliability, Safety & Compliance

Data centers are engineered for **high reliability** and must adhere to rigorous safety and compliance standards. This section covers how reliability is classified (e.g. Uptime Institute Tiers), the approaches to ensure safety of personnel and equipment, and key compliance frameworks (both infrastructure standards and regulatory compliance like security certifications).

Uptime Tiers and Reliability Objectives

The **Uptime Institute's Tier Classification** is a widely-referenced system to describe data center resiliency. There are **Tier I** through **Tier IV**:

- **Tier I:** Basic site infrastructure. Single path for power and cooling distribution, no redundant components ⁸⁶. This offers around **99.67% uptime** annually (about 28.8 hours of downtime) ^{86 41}. Tier I is essentially a server room with a UPS and a backup generator, but if anything fails or needs maintenance (power or cooling), the IT load must shut down. Suitable for non-critical needs due to its limited redundancy.
- **Tier II:** Redundant components with a single distribution path. This means some backup capacity (like N+1 on major components: UPS, generator, cooling units) but still a single path for power/cooling delivery ⁸⁶. Expected uptime ~**99.74%** (about 22 hours downtime per year) ^{86 41}. Tier II protects against some component failures (e.g. one generator can fail and another takes over), but **not against distribution path failures** or maintenance that affects the sole path. There may still be downtime for certain maintenance or if a path needs to be worked on.
- **Tier III:** Concurrently maintainable infrastructure. It has **multiple power and cooling distribution paths**, typically only one active at a time, but the other is available (so you can switch over for maintenance) ⁴⁰. Also includes redundant components (N+1) so that any single component failure does not impact IT operations. Tier III sites can undergo planned maintenance without shutdowns – for example, you can service a UPS or chillers by transferring load to the redundant module or secondary path. Expected uptime ~**99.982%** (about 1.6 hours downtime/year) ^{40 41}. Most enterprise-grade and colocation data centers aim for Tier III. This level assumes outages would mostly only come from very unlikely multiple failures or major events.
- **Tier IV:** Fault-tolerant infrastructure. This is the highest level, requiring **2N redundancy** for all systems (or sometimes 2N+1) and physically isolated dual distribution paths ^{42 43}. Tier IV can sustain a failure of any single system or distribution path without impact, and even tolerate a fault during maintenance of the other path (hence the 2N+1 concept) ^{43 44}. Uptime is ~**99.995%**, which is only ~26 minutes of downtime per year ^{42 41}. Achieving this means compartmentalization – for instance, separate UPS rooms for each path so a fire or leak in one doesn't affect the other, etc. ^{87 88}. Tier IV sites also need on-site 96-hour fuel for generators, among other stringent criteria ^{89 90}. This Tier is chosen by organizations that cannot afford any downtime (some banking transaction systems, certain military or emergency systems, etc.), though it comes at high cost. Many cloud and colo providers find Tier III sufficient because application-level redundancy can cover the rest.

It's important to note the tier classification focuses on infrastructure resilience, not on IT or operational processes. Also, a higher Tier isn't always better for every business – each tier fits different needs and budgets (PhoenixNAP, 2021) ^{91 92}. For instance, a Tier IV might be overkill for a typical enterprise that could tolerate a short outage or that can failover to another site. Most colocation providers advertise Tier III equivalent designs, with some offering Tier IV in select facilities (often charging premium).

Beyond Uptime tiers, reliability is also quantified by **SLA** percentages (e.g., "five nines" 99.999% which is ~5 min downtime/year). It's interesting that even Tier IV is not 100% - stuff like simultaneous maintenance

errors or multiple failures can still happen ⁹³. Ultimately, data center operators try to minimize outages of any kind. Uptime Institute surveys show that while major outages (> downtime or big impact) still happen industry-wide, best practices and tiered designs have reduced their frequency somewhat. Common causes of outages include human error, power system failures, and cooling failures. Robust procedures and testing are as important as design for maintaining reliability.

Failure Domains & Testing: Tiering aside, operators partition systems to limit failures. A *failure domain* might be one data hall on one UPS system – if something goes wrong there, it doesn't affect other halls on independent systems. Large cloud facilities might treat an entire data center as a failure domain and rely on software redundancy across multiple data centers (availability zones). Within a facility, routine **Integrated Systems Tests (IST)** are done (usually annually or semi-annually) to verify that backup systems kick in properly. This might involve simulate power failures, etc., similar to initial commissioning but periodically to ensure nothing has drifted.

Safety Measures

Operating a data center involves significant electrical power and critical equipment – safety is paramount for both personnel and equipment. Key safety considerations:

- **Electrical Safety:** Data centers have high capacity power systems (buses carrying thousands of amps). There is risk of arc flash – a dangerous blast that can occur during a fault or if equipment is mishandled. To protect workers, facilities do arc flash hazard analysis and mark equipment with incident energy levels, requiring appropriate Personal Protective Equipment (PPE) for certain operations. Often, live work is minimized; many modern designs include *make-before-break* maintenance features or use **flywheel UPS** that allow safer maintenance. The Uptime Institute Tier Standard also expects that you can isolate and de-energize systems for maintenance (concurrently maintainable for Tier III). Additionally, all sites follow the **National Electrical Code (NEC)** for proper grounding, wiring, and breakers to prevent electrical fires and shocks.
- **Fire Detection & Suppression:** Data centers use very sensitive fire detection – typically **VESDA** (Very Early Smoke Detection Apparatus) which continuously samples air for minute smoke particles. This allows catching a developing issue before it becomes a fire. For suppression, most server rooms use **clean agent** fire suppression systems (like FM-200, Novec 1230, or inert gas blends). These agents can extinguish fires without damaging electronics (unlike water). They are usually deployed via pipes in the ceiling, discharging gas if smoke or heat triggers the system. Some use **dual-interlock pre-action sprinkler** systems as backup – meaning pipes are filled with air and only fill with water if both heat and smoke triggers occur (to avoid accidental water leaks). NFPA standards (like NFPA 75, NFPA 76 for telecom facilities) guide these protections, as noted in compliance standards (Dgtl Infra, 2024 ³⁷ ⁹⁴). Fire suppression is tricky: you need to ensure any agent release doesn't harm people or cause more damage than the fire. Thus, typical approach is alert, attempt human intervention (operators are trained to use handheld extinguishers), and only automatically dump suppression if fire is confirmed.
- **Physical Security:** Data centers implement layered security to prevent unauthorized access (which could cause intentional or accidental harm). **24/7 security staff**, multi-factor access control (badge + biometric like fingerprint or iris scan), mantraps (one door must close before another opens), CCTV surveillance of all areas – these are standard. Server cabinets can be locked individually, and cage

areas for customers have their own badge readers. Security protects not just against malicious intruders but also ensures only qualified personnel enter critical areas (e.g., an untrained person shouldn't wander into a UPS room). Most providers adhere to standards such as **SOC 2 Type II** for security controls, which require strict access management and monitoring (SOC 2 checks that physical access to systems is controlled and logged) ⁹⁵ ⁹⁶.

- **Environmental Safety:** Cooling systems involve liquids (water or refrigerant). Data centers use things like leak detection sensors under raised floors or near CRAC units to catch any water leaks early. Also, backup batteries (especially large lead-acid banks) give off hydrogen gas when charging – battery rooms have hydrogen detectors and ventilation to prevent gas buildup and explosion risk. Any chemical (like diesel fuel, coolant chemicals) is stored and handled per safety regs. Diesel storage tanks have to meet fire code (often a maximum quantity indoor vs outdoor, with containment for spills).
- **Emergency Procedures:** Staff are trained on emergency power-off (EPO) procedures in case of an electrical fire (hitting EPO will kill power to protect people, though it brings the site down), evacuation plans, and how to safely shut systems if needed. Drills are often conducted. Some data centers have on-site *Operations Control Centers* that coordinate during incidents and have communication protocols to inform customers and authorities if needed.
- **Personnel Safety and 24/7 Staffing:** Data centers are often staffed at all hours. Ensuring staff safety means providing things like backup lighting (for when power fails, though critical areas often lit by emergency generator-backed circuits), redundancy in cooling to avoid dangerously high temps if one system fails (server aisles could get very hot quickly in a cooling outage). Many sites require at least two people on site at all times for safety (buddy system) when doing potentially hazardous tasks. Also, some tasks require certified professionals (e.g., only qualified high-voltage electricians operate certain switches).

Data center construction and operations follow OSHA regulations and any local health and safety laws strictly – with heavy equipment and high voltage, the environment necessitates robust safety culture.

Compliance and Certifications

Data centers often pursue various **compliance certifications and attestations** to demonstrate that they meet industry standards and regulatory requirements. Some key ones:

- **SOC 2 (Service Organization Control 2) Type II:** This is an audit framework by the AICPA for service organizations (like data centers or cloud providers) to report on controls related to security, availability, processing integrity, confidentiality, and privacy. A *SOC 2 Type II report* basically confirms that the data center has appropriate controls (e.g., physical security, network monitoring, incident response, etc.) and that these controls have been operating effectively over a period (usually 6-12 months). Many customers, especially in finance or SaaS, require a SOC 2 report from their colocation provider for their own compliance. Data center operators typically undergo annual SOC 2 audits and provide results to customers under NDA.
- **ISO/IEC 27001:** This is an international standard for information security management systems (ISMS). It ensures the organization has a systematic approach to managing sensitive information

(which includes physical and logical security). Data center companies get ISO 27001 certified to show they adhere to best practices in securing customer data and systems. It covers risk assessment, security policies, access control, incident management, etc. For example, Equinix and Digital Realty have many sites ISO 27001 certified (Equinix, 2023) ⁹⁷. Achieving this involves an external audit and maintenance of the ISMS with periodic reviews and continuous improvement.

- **PCI DSS (Payment Card Industry Data Security Standard):** This is relevant if the data center hosts systems storing or processing credit card data (cardholder data environment). While PCI primarily applies to the *systems*, some colocation providers also get a PCI attestation for their facility controls (physical security, etc.). Typically, the customers must ensure PCI compliance of their own systems, but having a PCI-compliant data center (with strong physical security, CCTV retention, access logs) helps. Some providers offer PCI audited spaces or secure cages for this purpose.
- **HIPAA:** For healthcare data (PHI – Protected Health Information), data centers can assert they are HIPAA compliant. HIPAA isn't a formal certification but providers sign Business Associate Agreements (BAA) and ensure they meet required safeguards (like controlled access, breach notification processes).
- **FISMA / FedRAMP:** If hosting US government systems, data centers might align with FISMA (Federal Information Security Management Act) requirements or go through FedRAMP (which includes physical infrastructure checks). Some colocation providers have FedRAMP Moderate or High ready facilities, meaning they've implemented the needed controls (continuous monitoring, personnel clearance, etc.) to simplify compliance for government customers.
- **Uptime Institute Certifications:** Separate from tier *design* certification, Uptime also offers **Tier Certification of Design Documents** and **Tier Certification of Constructed Facility** – some operators get their facilities formally certified as Tier III or IV by Uptime. They also have an **Operational Sustainability** certification (Gold, Silver, Bronze ratings) which assesses ongoing practices beyond just design. These are optional but can be a marketing advantage (showing a third-party validated the redundancy and operations). Additionally, organizations like TIA (Telecom Industry Association) have **ANSI/TIA-942** certification for Rated-3/Rated-4 which parallels Tier system and includes other aspects like telecom, safety, etc.
- **Environmental and Energy Management:** Data centers increasingly pursue **ISO 50001** (Energy Management) or **ISO 14001** (Environmental Management) certifications to demonstrate commitment to efficiency and environmental responsibility ⁹⁸ ⁹⁹. They track PUE, carbon footprint, implement energy optimization, and have processes to minimize environmental impact.
- **Green Building Codes/Standards:** Some data centers get **LEED** certification (Leadership in Energy and Environmental Design) to showcase building sustainability (e.g., use of efficient power and cooling, recycling water, sustainable construction materials). There are also standards like **BREEAM** in Europe similarly.
- **NERC standards:** If a data center is classified as critical infrastructure (for example, if it provides services to the power grid operations or is part of a utility's control system network), it might need to comply with **NERC CIP (Critical Infrastructure Protection)** standards which mandate rigorous cyber and physical security controls. Generally, commercial data centers aren't under NERC CIP

unless tied to the electric sector, but with grid interconnection issues, we see some interplay (e.g., if data centers agree to provide demand response, they coordinate with NERC guidelines to ensure reliability).

- **Local Codes & Standards:** Data centers abide by building codes (with perhaps enhancements). For example, California's building code might require certain seismic resilience for mission-critical facilities, Florida's might demand high wind ratings. Fire codes (NFPA) require things like emergency power off buttons and suppression systems. Often, authorities having jurisdiction (AHJs) treat data centers as essential facilities, which can require higher standards in construction (like a hospital would have).
- **Audit and Reporting Requirements:** Many customers audit their colo providers. A data center might need to support customer audits of physical security or provide reporting on environmental controls. The better providers preempt this by maintaining certifications so customers accept that instead of doing their own audit.

In essence, compliance is about **trust**: enterprises want assurance that a data center will keep their equipment secure and meet regulatory obligations. Achieving frameworks like ISO 27001 or SOC 2 shows an operator follows industry best practices for security and availability (Dgtl Infra, 2024) ⁹⁵ ³⁵. For example, ISO 27001 certification indicates a robust security management program (Equinix, 2023) ⁹⁷, and SOC 2 Type II report covers detailed controls on who can access the data center, how climate and power are monitored, how incidents are handled, etc.

Safety & Compliance Summary: Data centers pair **engineering for reliability** (redundant power/cooling, robust building) with **process and people reliability** (maintenance, training, security, compliance checks) to achieve high uptime. They must also ensure the **safety** of everyone and everything inside – through design (safe electrical systems, fire suppression) and procedures (lockout-tagout, emergency drills). And by aligning with compliance standards, they provide transparency and assurance to customers that their valuable data and operations are in a controlled and secure environment. An outage or security breach can be catastrophic, so the industry's ethos is very much *prevention and preparedness*. That's why you hear of multi-year streaks without downtime in top facilities and why data centers often look fortress-like – they're built and run to avoid failure even under duress.

8. Markets & Location Drivers

Data center development is concentrated in key geographic **markets** that offer the right mix of power, connectivity, economic incentives, and other factors. In the U.S., several metro areas have become major data center hubs, each with its own drivers. We will profile some top markets and what makes them attractive:

- **Northern Virginia (Ashburn/"Data Center Alley" – IAD):** Northern Virginia, especially Loudoun County (Ashburn), is the largest data center market in the world by capacity ¹⁰⁰. Over 2.9 GW (2900 MW) of data center load is located here as of 2024, more than the next several markets combined (DCD, 2025) ¹⁰⁰. **Why NoVA?** Historically, MAE-East (one of the earliest Internet exchange points) was in this region, seeding excellent fiber connectivity. Today, Ashburn has dozens of carriers and the highest density of **fiber routes** – it's essentially the East Coast internet hub. **Power** is relatively affordable and scalable: Dominion Energy provides abundant electricity (often ~\$0.05–\$0.07/kWh

range for large users) and has built many substations for new campuses. Virginia also offers a **state sales tax exemption** on data center equipment (servers, generators, etc.) for qualifying investments ¹⁰¹ ¹⁰², which massively reduces costs for operators – Virginia's program cost the state ~\$750M in 2023 in foregone taxes due to the booming construction ¹⁰³ ¹⁰⁴. Additionally, Loudoun and neighboring counties have welcomed data centers with expedited permitting ("Fast-Track" programs) and set lower property tax rates on data center equipment to attract business ¹⁰⁵ ¹⁰⁶. **Land** was plentiful in this outer D.C. suburb (though now premium), and critically, the region is low-risk (no common natural disasters apart from occasional storms). Water is moderately available for cooling (though there are concerns as usage grows). The typical deployment in Ashburn is a multi-building campus – companies like Equinix, Digital Realty, Amazon, Microsoft, and many colocation providers have huge footprints. It's divided into "clusters" like Ashburn, Sterling, Manassas, and now expansion further west. Latency from Ashburn to major East Coast cities is excellent (~6–8 ms to New York, <2 ms within the region to D.C.), making it ideal for serving Eastern US traffic and government clients. **Nickname:** "Data Center Alley" – with some of the world's largest internet exchanges and reportedly up to ~70% of global internet traffic passing through at some point (often cited anecdotally due to major interconnects there).

- **Dallas-Fort Worth (DFW, Dallas Plano Garland – DAL):** The Dallas metro area is one of the top data center markets in the U.S. (historically ranked #2 by inventory). It boasts a **central location** in the U.S. which is geographically advantageous – roughly equidistant latency to East and West coasts (~20–30 ms), making it a good hub for serving North America. **Power** in Texas is deregulated and ample: utilities like Oncor or CoServ supply DFW with relatively low-cost power (Texas has lots of wind energy and natural gas). The state has no income tax and often municipalities provide property tax incentives for data center projects. For example, the town of Garland or Plano have offered tax abatements to attract facilities. Texas also has a sales tax exemption similar to Virginia's for data center equipment if investment thresholds are met. **Fiber connectivity:** Dallas has long been a telecommunications hub – it's a crossroad for many fiber backbones (north-south from Chicago to Houston, east-west from Atlanta to Los Angeles). The downtown carrier hotel (Infomart building) anchored connectivity, and now many carriers are present in suburban campuses too. **Land** in suburban Dallas (e.g., in Garland, Plano, Allen) is relatively affordable and flat, and zoning is generally business-friendly. **Risk:** Dallas has some tornado risk but most data centers there are built to withstand it (concrete structures). The climate is hot, but not extremely humid, and many facilities use air or evaporative cooling (with careful eye on water use). Notably, Texas's independent power grid (ERCOT) has had issues (like the 2021 freeze outages), but data centers in DFW generally rode through those on generators where needed. **Latency:** DFW is ~20ms to Ashburn, ~25ms to Silicon Valley – a good midpoint to aggregate or redistribute traffic. Also, Dallas is a gateway to Latin America with some fiber routes heading south.
- **Phoenix, Arizona (PHX):** Phoenix has risen as a major market due to a combination of low cost and low natural disaster risk. **Power** is reasonably priced and often comes from a mix including nuclear (Palo Verde) and solar; Arizona Public Service (APS) and Salt River Project (SRP) are utilities known to support data center growth. Notably, Arizona offers **tax incentives**: a transaction privilege tax (sales tax) exemption on data center equipment for large investments, similar to VA/TX. **Climate:** It's very hot in summers, but extremely dry – which ironically can be good for data centers if they use evaporative cooling (very effective in dry air, albeit with water usage). Many Phoenix data centers leverage **adiabatic cooling** for much of the year to get PUEs ~1.3 despite high ambient temperatures. The dryness and design allows lots of **free cooling hours** in cooler months. **Water:**

ironically, water scarcity is a concern. Some data centers in Phoenix use reclaimed water for cooling to alleviate potable water use. **Land:** Phoenix (especially suburbs like Chandler, Mesa, Goodyear) has large tracts of cheap land, often in industrial zones. Big players like CyrusOne, Microsoft, Google have large campuses there. **Connectivity:** Phoenix is a key interconnection point for the southwest. It's on routes connecting California to Texas, and houses cloud region infrastructure (e.g., one of the major US-West regions for some clouds). It's also within ~10 ms latency of Los Angeles, making it a viable DR site or alternate hub for West Coast traffic. There's an Internet exchange (Phoenix-IX) and growing carrier presence. **Risk:** Very low risk of earthquakes (solid ground, unlike California). No hurricanes. Some dust storms (haboobs) but those mainly affect filters on cooling systems. Generally stable geology and weather beyond the heat. These factors, plus the incentives, have made Phoenix one of the fastest growing markets; indeed by some 2025 figures, Phoenix became the #3 largest market in North America by capacity, overtaking some older markets (DCD, 2025) ¹⁰⁷ ¹⁰⁸. Companies are drawn to its combination of *cost-efficiency and resiliency*.

- **Atlanta, Georgia (ATL):** Atlanta has historically been a network hub of the Southeast U.S. – it's where many long-haul fiber routes intersect (Florida up the coast, east-west across the south). **Connectivity:** Major carriers and cable routes go through Atlanta; it's also a major peering point for the SE region. In 2024, Atlanta saw a *huge* jump in data center investment – it reportedly led the U.S. in net absorption, even surpassing NoVA for the year ¹⁰⁹ ¹⁰⁷. This growth is partly driven by new hyperscale builds (e.g., Microsoft building big campuses nearby) and availability of land/power. **Power:** Georgia Power offers relatively low electricity rates, and Georgia has a robust grid with diverse generation (nuclear at Plant Vogtle, lots of natural gas, growing solar). The state implemented tax exemptions for data centers (sales tax on equipment) for qualifying large projects, making it financially attractive. **Land:** Atlanta's suburbs (like Douglas County, Lithia Springs area, or northeast in Suwanee) have ample land. The state and local governments actively court data centers, touting low natural disaster risk and low cost. **Risk:** Minimal – inland location (no hurricanes direct, though can get remnants), mild seismicity, and while tornadoes occur in GA, Atlanta metro is not in Tornado Alley's core. **Water:** abundant water resources (no desert here) which makes cooling easier with evaporative methods if desired. **Use cases:** Atlanta is popular for both enterprise colocation and hyperscale cloud regions (Google and Microsoft have regions here). Latency: ~12 ms to Ashburn, ~20 ms to Chicago, ~8 ms to Miami – so it's well positioned to serve the southeast and as a junction to Florida and Latin America (some Latin American connectivity goes via Miami though). One thing boosting Atlanta: it's seen as an alternative to Northern Virginia as that area's power becomes constrained. Indeed, industry reports highlight Atlanta adding more capacity in 2024 than it had in total prior, pushing it among top markets (CBRE, 2024) ¹⁰⁹ ¹⁰⁷.
- **Silicon Valley (San Jose/Santa Clara, CA - SJC):** The Silicon Valley area (Santa Clara, San Jose) is a historic data center hub given its proximity to countless tech companies. **Connectivity:** It's a major West Coast network hub, tying in trans-Pacific submarine cables (many land in nearby Hillsboro, OR or L.A. and connect up) and domestic backbones. Equinix's early flagship sites were in Silicon Valley (like SV1 in San Jose) and many networks peer there. **Challenges:** Silicon Valley has some of the most expensive power (~\$0.12+ per kWh for commercial) and land in the country. The local utility (Silicon Valley Power in Santa Clara) actually offers relatively cheaper power than PG&E and is a key reason many data centers cluster in Santa Clara city. Santa Clara's municipal utility historically provided reliable, coal-free power at rates half of PG&E's, making it attractive. However, **power availability** has become a problem – the grid in that area is nearing capacity. By 2021–22, Silicon Valley Power had to temporarily halt new large data center hookups until new substations could be built. This

effectively slowed growth. **Environment:** Earthquake risk is significant (the area is near major fault lines). Thus, data center buildings are engineered with seismic bracing/base isolation sometimes. Water is moderately available (not like Phoenix, but CA does have drought concerns; some use recycled water for cooling towers in Santa Clara). **Climate:** mild Mediterranean, which is great for air economization – lots of free cooling days. **Tax:** California doesn't have broad data center incentives like VA or AZ, and it has high taxes, but some local incentives exist (e.g., Santa Clara sometimes negotiates deals on utility user taxes). Many data centers persist here because of adjacency to corporate HQs and talent, plus latency – *ultra-low latency* to the tech companies' offices or between financial exchanges in SF and trading platforms needed being local. But many companies have also expanded elsewhere (like in Oregon or Arizona) to avoid the high cost. As of 2025, Silicon Valley's growth has been relatively slower and it was possibly overtaken by Phoenix and Atlanta in total capacity ranking ¹⁰⁷ ¹¹⁰. Still, it remains a crucial interconnect location and houses key infrastructure (like cloud region sites – AWS us-west, etc.).

- **Seattle, Washington (SEA) & Portland, Oregon (PDX):** The Pacific Northwest has two notable markets: **Seattle** and **Hillsboro (Portland metro)**.

Seattle (including cities like Quincy, WA where Microsoft and others have large farms, and the eastern Seattle suburbs) became significant largely due to cloud (Microsoft and Amazon's home bases). **Power:** Washington has abundant cheap hydroelectric power (like from Columbia River dams) – electricity can be very low cost (in Quincy, < \$0.04/kWh historically). Quincy's Grant County PUD had special rates for large loads. The climate is cool much of the year, aiding PUE via economization. **Connectivity:** Seattle is a landing for some trans-pacific cables (though more land in Oregon now). It's also a regional hub connecting Alaska, Asia, and providing an alternate path to California. **Latency:** ~25 ms to Los Angeles, ~30 ms to Ashburn. **Risk:** some seismic risk, though less than SF; volcanic (Mt. Rainier) is extremely low probability but present; generally stable. Seattle area has significant cloud region presence (AWS, Azure). However, not as many multi-tenant colo as SV or VA – it's more dominated by self-builds or a few key colos (e.g., Sabey's Intergate campuses).

Hillsboro, Oregon (outside Portland) has surged recently (sometimes counted as "Portland" market, PDX). **Drivers:** It's a major landing point for *submarine cables to Asia* (new cables from Japan, etc., terminate in Hillsboro). That turned it into an intercontinental hub with several IXs and cloud on-ramps. Oregon also offers **zero sales tax** (no state sales tax at all), which acts as an incentive on all that expensive equipment. On top, certain Enterprise Zones in Oregon provide property tax abatement for data centers for up to 5 years. **Power:** Portland General Electric (PGE) and other local utilities have a lot of hydro and wind in mix; power rates moderate (~\$0.05–0.06). **Climate:** mild and cool, ideal for free cooling (and lots of rain means water for cooling towers). **Risk:** low natural disaster incidence (some seismic risk from Cascadia subduction zone exists, but infrequent; volcano risk minimal short-term). Hillsboro hosts multiple large colocation campuses now (e.g., Digital Realty, Flexential, NTT) and hyperscale builds. It effectively became the **connectivity gateway** for US-West to Asia, rivaling Los Angeles in that niche, because of new cables like Jupiter, New Cross Pacific, etc. So companies needing Asia connectivity might favor Hillsboro (latency to Tokyo ~120 ms via direct cable, which is excellent). This location's growth shows how **submarine cable landings + tax benefits** can seed a market.

- **Chicago, Illinois (CHI):** Chicago has long been a primary U.S. data center market due to its central location and status as a financial hub. **Connectivity:** Chicago is the crossroads of many U.S. fiber routes (coast-to-coast links often pass through, and north-south links too). It's home to big carrier

hotels like 350 E Cermak (the iconic downtown carrier hotel in a former printing plant), which is extremely network-dense. Chicago's financial sector (options/futures exchanges) also drives demand for ultra-low-latency connectivity; many trading firms put infrastructure in data centers near the exchanges (some are in Aurora, IL for proximity to CME's matching engine). **Power:** Illinois historically had decent power rates (ComEd utility) and relatively robust grid. The state in 2019 also enacted sales tax incentives for data centers (exemption on equipment and materials for 20 MW+ projects that create jobs), which boosted growth. **Climate:** Cold winters, warm summers – but plenty of cool periods for economization. Water from Lake Michigan is ample for cooling. **Risk:** Low natural disaster risk overall – minimal quake risk, some chance of tornadoes but less frequent than Plains states, and no hurricanes. **Tax:** besides sales tax breaks, some suburbs offer property tax incentives. Chicago is an interesting case: while still big, its growth has been a bit slower relative to sun-belt markets. However, it remains top 5 in size. There is also a trend of some data centers in nearby states (like in Iowa where Google, Facebook have large sites due to cheap power and tax breaks) that serve similar central U.S. roles. Latency: from Chicago ~15 ms to Ashburn, ~40 ms to Silicon Valley, making it a good hub for nationwide service.

- **Other Notables:** The prompt specifically listed IAD, DFW, PHX, ATL, SJE (San Jose/Silicon Valley), SEA, PDX, CHI. These indeed cover the major U.S. markets. For completeness: **New York/New Jersey** is also large (with financial sector driving it), **Los Angeles** (with content/media and as Pacific gateway) is significant, and internationally, places like London, Frankfurt, Singapore, etc., are major – but the focus here is U.S.

Each market's attractiveness can be summarized by **power, fiber, taxes, land, and water** as the prompt suggests:

To illustrate succinctly:

- **NoVA (Ashburn): Power:** High availability (Dominion) but facing future constraints; cost low-moderate; **Fiber:** unparalleled density, MAE-East heritage; **Taxes:** strong incentives (sales tax exempt); **Land:** previously plentiful, now tighter but still expanding (some concerns about suburban sprawl and pushback from residents on aesthetics/noise); **Water:** available but huge usage is raising local issues (Loudoun is planning reclaimed water systems for data centers). **Latency hub:** main hub for U.S. East, transatlantic connectivity (to Europe) excellent, hub for Federal government networks too.
- **DFW:** **Power:** Plentiful (Texas generation capacity large), competitive providers, some risk from ERCOT issues but data centers mitigate with generators; **Fiber:** national crossroads, many carriers (Lumen, AT&T etc. have big presence); **Taxes:** beneficial (no income tax, and Freeport exemptions in Texas counties for business personal property can reduce tax on movable servers, plus state incentives for big projects); **Land:** ample in suburbs, cheap; **Water:** not scarce (but summers hot mean large water use if evaporative); **Latency:** central hub, also house to some hyperscalers (though Google notably doesn't have region there, AWS & Azure do), and a growing edge for Mexico/LatAm connectivity.
- **PHX:** **Power:** Utility investment by APS/SRP to support large campuses, heavy use of solar potential; **Fiber:** decent – not as many carriers as LA but enough long-haul links (many go through Phoenix toward Texas or up to Utah); **Taxes:** strong state incentives (TPT exemption); **Land:** vast and cheap, especially in West Phoenix (Goodyear/Buckeye) which is exploding with hyperscale builds; **Water:**

major concern – data centers in Phoenix must innovate (e.g., using reclaimed water or running higher allowable temps to use less water; some use liquid cooling to avoid evaporative towers).

- **ATL:** **Power:** Expanding generation, nuclear (Plant Vogtle new reactors coming online gives stable supply); **Fiber:** Southeastern hub, at intersection of many long-hauls and cable landing going to Virginia or Florida; **Taxes:** yes, GA tax credit for job creation and equipment; **Land:** available around metro (new big builds in Douglas County for example); **Water:** ample (but must manage with city infrastructure); **Latency:** main in Southeast, sub-20ms across much of Eastern half of US, connectivity to Latin America via Florida.
- **Silicon Valley (San Jose/Santa Clara):** **Power:** Provided by Silicon Valley Power (for Santa Clara) with historically 100% uptime to data centers but now limited new load capacity; expensive regionally; **Fiber:** huge – original MAE-West and PAIX were here, lots of carriers; **Taxes:** no special incentive, high baseline cost (though Santa Clara has slightly lower property tax on data center equipment relative to SF due to enterprise zone historically); **Land:** extremely expensive and now basically fully utilized in Santa Clara (some providers converting old office lots to multi-story data centers to maximize space); **Water:** moderate (SC has some recycled water lines for cooling towers), risk of drought restrictions in CA always something to monitor; **Latency:** key West Coast hub, needed for minimal latency to SF tech and APAC cables.
- **Seattle:** **Power:** cheap hydro, but some areas now queue for substation upgrades (like in Quincy, large loads forced upgrades), overall positive; **Fiber:** good, Seattle Internet Exchange (SIX) is big, carriers like CenturyLink (headquartered in Monroe, LA but huge presence in Seattle from old Qwest/Worldcom), regional networks; **Taxes:** WA also has data center sales tax exemptions for rural areas (Quincy benefitted, but in King County not as much); **Land:** Seattle proper limited, but outlying (Eastern WA) abundant; **Water:** plentiful (rain, rivers); **Latency:** stepping stone to Asia, and to Vancouver/Canada.
- **Hillsboro/Portland:** **Power:** moderate cost, some direct access to BPA (Bonneville Power Admin) hydro possibly; **Fiber:** superb for transpacific connectivity (new cables mean an operator in Hillsboro can directly peer with Asian networks with one less hop), carriers like Wave, Zayo heavily present; **Taxes:** no sales tax big plus, enterprise zone abatements; **Land:** available in Silicon Forest area and welcomed by local authorities; **Water:** plentiful; **Latency:** key for US-West to Asia, slightly higher latency to LA (~20ms) and to Seattle (<10ms) but negligible; a good alternative to SV.
- **Chicago:** **Power:** ComEd improving grid, nuclear-heavy mix yields stable supply; **Fiber:** top-tier, 350 Cermak and others are big peering points (Chicago houses one of largest internet exchange points in US – Equinix's and community ones); **Taxes:** IL introduced sales tax breaks in 2019, making it more competitive; **Land:** in suburbs like Elk Grove Village, lots of data center campuses (EGV actively fostered a cluster with tax incentives); **Water:** Lake Michigan – basically infinite cooling water; **Latency:** ideal central location, e.g., Chicago to London ~40ms via transatlantic from east coast plus backhaul, to Tokyo ~100ms via west coast – often used as a mid-point aggregation location and disaster recovery site for coasts.

The markets above represent strategic location decisions: companies may choose a market for proximity to users (latency), network access (peering), cost optimization, and risk diversification. Increasingly, **emerging**

markets or secondary cities (e.g., Denver, Minneapolis, Miami, etc.) also see growth for edge deployments or regional coverage. But the listed hubs concentrate the lion's share of capacity and connectivity.

Finally, mention that these U.S. hubs align with where the big **hyperscale self-builds** and **colocation footprints** go. Northern Virginia, for instance, hosts many of the east coast availability zones of AWS, Azure, etc., and is the default location for federal cloud regions. Dallas and Chicago often function as interior hubs, Phoenix and Hillsboro as alternative west coast hubs given California's constraints, and Atlanta as the south's answer to Ashburn's dominance. This geographic distribution also ensures *lower latency to end-users nationwide*: e.g. someone in Florida gets better latency if served from Atlanta than from Ashburn, etc.

To conclude the section, one could highlight that these markets also have *massive economic impacts* locally – e.g., Loudoun County's tax revenue from data centers was \$663 million in 2022 ¹¹¹ ¹¹², and in general, data centers are now an important part of real estate investment and infrastructure strategy in these regions.

9. Operator Landscape (neutral overview)

The data center industry includes a range of **operators** with different business models. Broadly, we can categorize them into a few groups:

- **REIT (Real Estate Investment Trust) Data Center Companies:** These are publicly traded companies focused on data center real estate and services, many structured as REITs for tax purposes. Examples: **Equinix**, **Digital Realty**, **CoreSite** (was public REIT, now part of American Tower), **CyrusOne** (was a REIT, went private in 2022), **QTS Realty Trust** (went private under Blackstone in 2021), **Iron Mountain Data Centers** (a division of Iron Mountain REIT). REIT operators often have global or national footprints, offering colocation space to enterprises, content companies, and cloud providers. Many REITs specialize: Equinix is known for retail colocation and rich interconnection ecosystems (with 251 data centers globally as of 2024) ¹¹³, focusing on many customers in each site (Equinix's average site has hundreds of networks and tenants, making it a "marketplace"). Digital Realty (312 data centers globally) ¹¹⁴ historically leans toward larger-footprint customers and wholesale deals, though it also runs an interconnection business (especially after acquiring Interxion in Europe). REITs have access to capital and tend to acquire as well as build – the industry has seen consolidation (e.g., Digital buying DuPont Fabros, Equinix buying Verizon's data centers, etc.). **Positioning:** These companies emphasize **neutrality and scale**. They provide infrastructure to any customer, rather than being tied to one cloud or one carrier. For instance, Equinix hosts all major clouds within its facilities as well as thousands of enterprises – essentially becoming hubs for hybrid IT. Digital Realty similarly serves a broad range (and runs campuses where a cloud provider might lease a full building alongside enterprise colocation in another).
- **Private Data Center Operators:** There are many significant operators not publicly traded. Some are backed by private equity, infrastructure funds, or large parent companies. Examples: **Aligned** (US-based, focused on efficient cooling and ability to scale customers' power density up – backed by Macquarie's infrastructure fund), **Vantage Data Centers** (builds globally, backed by DigitalBridge and others), **Flexential** (mid-size colocation provider in US, PE-owned), **TierPoint** (regional retail colo player), **Cyxtera** (was public via SPAC but filed Chapter 11 in 2023, now restructuring privately), **Switch, Inc.** (large high-density campuses in Nevada and Michigan, went private acquired by DigitalBridge in 2022). Also **NTT Global Data Centers** (part of NTT, Japan – one of the largest

globally with 95 data centers, ~1.1 GW capacity) ¹¹⁵, **Compass Datacenters** (build-to-suit and edge focused developer, private). These players often specialize in either **hyperscale leasing** (building whole campuses for one or few cloud customers) or **enterprise colocation** or both. For example, Vantage and Aligned primarily do wholesale/hyperscale deals providing custom large suites to single tenants, whereas TierPoint and Flexential focus on retail colo and managed services in second-tier markets. **Positioning:** Private operators often highlight flexibility, customer service, or specific innovations. Aligned, for instance, touts its patented cooling system and ability to **expand power density on demand** (cooling and power infrastructure that can support variable loads without overbuilding) – appealing to customers worried about future AI loads, for example. Switch (before acquisition) marketed its Tier IV **ultra-secure** designs and 100% renewable power use.

- **Cloud/Hyperscale Self-Build:** The major cloud providers – **Amazon (AWS)**, **Microsoft (Azure)**, **Google (GCP)**, **Meta (Facebook)**, **Apple** – all build and operate their own massive data center campuses for their internal use. These are not multi-tenant (except cloud tenants in a logical sense). Often these hyperscalers also lease significant capacity from the colo operators above, essentially as single-tenant leases. For instance, in 2022 about half of the wholesale data center leasing in North America was driven by cloud companies expanding (they lease from REITs or private operators when it's faster or more flexible than building). **Hyperscaler build vs. lease:** Amazon tends to lease a good portion of capacity (especially in colocation for edge locations), Microsoft and Google also lease but also construct many of their own sites. Meta and Apple predominantly self-build as they want control and custom specs (Meta has 24 huge data center campuses globally) ¹¹⁶. These self-builds are very large (often hundreds of MW) and highly optimized for specific needs (e.g., Meta's are tailored for massive storage and caching, with their own Open Compute server designs). They may not have the same redundancy level as a commercial colo – often built to somewhere between Tier II and III (since the applications provide resilience across sites). **Positioning:** Hyperscalers generally don't resell capacity (except where they open it for government like AWS GovCloud etc., but that's still their service). However, their presence influences the whole ecosystem – they set standards in innovation (like Google pioneered high-voltage DC distribution in some data centers, Microsoft is testing hydrogen fuel cells, etc.), and their demand drives much of the construction in the industry.
- **Telecom and Edge Operators:** Historically, telcos like AT&T, Verizon, CenturyLink (now Lumen) operated many data centers (some for internal, some retail colo) but most have divested their colo businesses to the above companies. Verizon's went to Equinix, CenturyLink's to Cyxtera, etc. Now, **tower companies** are edging in: e.g., **American Tower** bought CoreSite in 2021 – interestingly bridging cellular tower infrastructure with data centers, potentially to develop edge computing at tower sites tied into core colocation. Another example: **EdgeConneX** (founded 2013) built smaller data centers in second-tier markets and at network edge locations (close to ISP hubs) – it's now an important private operator (majority owned by EQT Infrastructure) with both edge sites and some large hyperscale builds globally. There are also startups focusing on **micro data centers** often at telecom sites or 5G edge (but that market still evolving; many are experimental or in partnership with carriers).
- **International vs. Regional:** Some operators have global footprints (Equinix and Digital Realty are worldwide; NTT covers Asia, US, EU; Iron Mountain has multi-country now; CyrusOne started US but now in Europe too, etc.). Others are regional specialists – e.g., **Cologix** focuses on Canada and some US cities (carrier hotel sites), **Interxion** (now Digital Realty) was Europe specialist, **ST Telemedia** in

Asia, **Chindata** in China/Asia, etc. In highly regulated markets like China or Russia, local operators dominate due to data sovereignty rules.

Industry Structure: Many large colocation providers (especially REITs) are *neutral*, meaning they provide space and power and let customers manage their own equipment (and carriers interconnect freely). Some also offer value-add services like managed hosting, cloud on-ramps, or consulting, but generally avoid competing with their customers. That neutrality is key to their positioning as platforms where ecosystems form (Equinix is explicit about this ecosystem play).

Consolidation trends: The last decade saw lots of M&A – e.g., Equinix and Digital Realty together now account for a significant share of multi-tenant data center (MTDC) market globally (by some estimates, Equinix has ~8% of global DC space market and Digital ~6% by 2024) ¹¹⁷ ¹¹⁸, though cloud self-build is a bigger portion in terms of sheer capacity. Private equity has taken many players private (CyrusOne, QTS, Switch, etc.), indicating investors see stable, utility-like returns. There are also specialized players like **DCI Data Centers** in APAC or **NEXTDC** in Australia focusing on specific regions.

High-Level Positioning Summary:

- *Equinix*: Largest colo operator globally (by revenue and footprint), retail-focused, interconnection leader. (Operates 250+ data centers, heavy in network-dense sites) ¹¹³.
- *Digital Realty*: Very large global provider, historically more wholesale but also retail, runs huge campuses and has presence in all major metros; also invested in interconnection with its "PlatformDIGITAL" concept and ServiceExchange fabric.
- *Other REITs (pre-privatization)*: CoreSite (strong in network-centric US markets e.g. LA, Denver, DC – now under American Tower, integrating with tower edge?), CyrusOne (focused on large enterprise/hyperscale wholesale in US/Europe – now under KKR and GIP, so private), QTS (did a lot of hyperscale deals, now under Blackstone).
- *Private & Infra-backed*: Aligned (innovation in cooling, selling sustainability – e.g., they use a pay-for-use model where customers aren't forced to pay for unused capacity, and boast very low PUE via their cooling tech), Vantage (rapidly expanded via investment funds, building mainly large campuses for cloud in NA and Europe), NTT GDC (leveraging NTT comms customer base and aggressive expansion, especially in APAC and now US after buying RagingWire), Iron Mountain (leverages brand trust from its records management biz to attract compliance-focused clients; also invests heavily in green power purchase agreements for 100% renewable goal).
- *Hyperscalers*: AWS, Azure, Google – not providers of colo but *major consumers* of space from others plus operators of their own large facilities. They sometimes partner (e.g., AWS Outposts or Azure Stack allow placing cloud-managed racks in colocation sites for hybrid, and these providers partner with Equinix, Digital, etc., to host them). Also, cloud on-ramps cause cloud providers to have presence in colo (often a single cage where they put their networking to connect to customers, not an entire DC).
- *Categories summary*: **Retail Colo Providers** (e.g., Equinix, CoreSite, smaller ones like 365 Data Centers regionally) target many customers in each site, offering connectivity and often managed

services. **Wholesale Colo Providers** (e.g., CyrusOne, Digital Realty, Aligned) target fewer, larger deals (whole rooms or buildings for single tenant). Many companies do both or have hybrid models now. **Specialized Edge Providers** (EdgeConneX, Vapor IO, DataBank – DataBank runs many edge sites and also cloud adjacents) focusing on emerging needs for edge computing near metro populations or 5G.

All operators differentiate on things like *uptime track record, scalability, network richness, security features*, and increasingly *sustainability* (e.g., who can offer carbon-neutral options, renewable energy sourcing – many top operators have 100% renewable energy goals or already match use with renewables (Equinix, Digital, Iron Mountain all claim >90% renewable power for their portfolio in recent years)).

Because the user asked for “neutral overview,” we avoid promotional tone. So for example, we wouldn’t say “Equinix is the best at interconnection” as an opinion, we’d say “Equinix has a large global footprint and is known for dense interconnection, operating over 460k cross-connects (Equinix, 2024)⁸⁴”. Similarly, mention Digital’s revenue size to neutrally illustrate scale: (Digital Realty had about \$5.5B revenue in 2024, vs Equinix \$7.3B)¹¹⁹ ¹²⁰.

Important note: not all REIT are independent now (some taken private but still operate similarly, just under private ownership but we’ll still mention them as operators).

Neutral positioning: In summary, the landscape ranges from *carrier-neutral multitenant providers* (the colos) to *single-tenant hyperscale builds*, with some hybrid in between (like wholesale providers who effectively have one big tenant per site). The existence of neutral colos fosters the ecosystems where networks and enterprises meet (e.g., at Equinix). Meanwhile, hyperscalers rely on both their own sites and leasing from those neutral providers to meet their explosive growth – there’s a symbiotic relationship.

For a business reader: this means customers can either go directly to these providers for data center space or consume services from cloud providers who behind the scenes work with these operators.

10. Procurement & Commercial Models

When companies procure data center capacity, there are different **commercial models** depending on the scale and needs. The two broad categories of colocation we discussed – wholesale vs retail – correspond to different contract and pricing structures. Additionally, data centers have various charges like power commit and cross-connect fees, and cloud usage has its own cost considerations (like egress fees).

Wholesale Colocation Contracts: In a wholesale deal, a customer (often a large enterprise or cloud operator) leases a dedicated space (could be a suite, a full data hall or even a whole building) from the provider. These leases resemble real estate leases in many ways. Key features:

- **Longer Term:** Wholesale leases are typically **5 to 10+ years**. The customer makes a longer commitment because the provider is often customizing build-out for them and giving them a better \$/kW rate in return for predictability.
- **Power Commitment:** Pricing is heavily based on the amount of power reserved (kW or MW). For example, a deal might be “3 MW at \$x per kW per month, with an annual escalator of 2%”. The

customer often pays for the committed capacity regardless of actual usage (like a minimum bill). This *reserved capacity* is critical because the provider invests in infrastructure to support that load.

- **Rent and Utility Separation:** Many wholesale contracts are structured like triple-net leases: the customer pays a base rent for the space/infrastructure and then also covers variable costs such as electricity usage (passed through at cost). So a bill might have a fixed monthly charge for space & infrastructure, plus a metered power charge. Some deals just incorporate a power allowance into a flat rate but have surcharges if usage exceeds certain threshold.
- **Customization and Fit-out:** The customer often can customize layout, maybe choose their own racks or even their own chillers in some cases (some wholesale is almost like *building to spec* for the tenant). The demarcation of responsibility is defined – e.g., provider up to the PDU, customer beyond that, etc.
- **SLA and Reliability Guarantees:** Wholesale providers still offer SLAs (usually 99.999% power uptime or similar). If an outage happens, typically a remedy is a rent credit for affected time. Wholesale clients often negotiate stricter terms because of their scale.
- **Expansion Options:** Contracts may include options or rights of first refusal on adjacent space for growth. This is important for scalability – e.g., a lease might say the tenant can expand into an adjacent 1 MW within 2 years at predefined rates if needed.

Retail Colocation Contracts: These are more like service agreements than real estate leases, often shorter term and more all-inclusive.

- **Term and Size:** Retail colo deals can be as short as 12 months, often 24-36 months standard for a rack or small cage. They cover anywhere from a single rack (the unit might be 42U rack with a certain power circuit) to a few cages of racks. The scale is smaller, but some retail-focused providers do allow multi-MW in a retail model for certain clients – the line blurs if the client wants the provider to manage a lot for them.
- **Pricing Units:** For a single cabinet, pricing might be a flat monthly fee per cabinet including a certain power circuit (e.g., \$X per month for a locking 4kW cabinet and a 20A power feed). If the client draws more power than allocated (some allow burst up to a limit, but generally circuits have a capacity), they might need to upgrade to a higher kW provision which costs more. For cages (multiple racks), pricing can be per kW of power allocated plus some fee per square foot. It's common to see quotes like "\$1200 per kW per month" in major markets for small deployments (this would include space, cooling, electricity, remote hands availability, etc. – note per kW pricing lumps things).
- **Power Circuit & Usage:** In retail, sometimes it's a flat fee for a circuit of certain amperage (assuming a certain usage profile). Some colos meter the actual power consumption and charge for energy used (especially for larger deployments), often plus a mark-up or just pass-through. Many will include let's say the first 200 kWh per month per kW in the base price then charge excess. The complexity varies by provider.

- **Cross-Connect Fees:** Very important in retail colocation – each connection (fiber or copper) to another party or to a carrier in the meet-me room typically has a setup fee (perhaps \$200 one-time) and a recurring fee (commonly \$100-\$300 monthly per cross-connect). These can add up if a company connects to multiple carriers and to other systems. For example, connecting to two ISPs and a cloud on-ramp and an IX could mean 4 cross-connect fees. Equinix has a large portion of revenue from interconnection fees. These fees might seem nominal individually but encourage an ecosystem (and also lock customers in, because moving would mean re-establishing all those connections).
- **Bandwidth/Internet Services:** Some retail colos also sell IP transit or blended bandwidth as an add-on if customers prefer a one-stop-shop (especially smaller businesses might just buy bandwidth from the colo provider). Otherwise, the customer buys from carriers and uses cross-connects.
- **Remote Hands & Other Services:** Retail colos usually offer “remote hands” support – technicians to do basic tasks for customer equipment (power cycle a server, plug a cable, swap a tape, etc.) when the customer’s staff can’t be on-site. These are often charged hourly (e.g., \$150/hour) or some contract might include a few hours per month. Wholesale providers offer it too, but retail customers use it more heavily since they often fully rely on the colo provider for physical presence. Other services include managed firewall, rack and stack services, migration assistance, etc., depending on the provider’s portfolio (some stick purely to space&power, others upsell managed services).

Power Usage and Billing Models: A critical aspect is **efficiency and over-subscription**. In wholesale, if a customer reserves 1 MW, they pay for 1 MW whether or not they use it (ensuring the provider builds that capacity ready). Some deals allow a bit of burst above commit (with a fee) or tiered pricing if they consistently underuse. In retail, providers often over-subscribe power assuming not all customers will draw peak at same time – but they monitor to ensure the total draw stays within facility limits. If a facility fills up, sometimes providers will ask high-usage customers to re-evaluate their circuit (maybe sell them a higher kW plan).

Comparison in simple terms: Retail is like renting *seats on a bus* – flexible, you pay per seat per month, and you can get more or fewer easily, with the “bus operator” providing a lot of the maintenance. Wholesale is like leasing an entire bus or fleet of buses – you commit to it and manage more of the specifics, but you get economy of scale in pricing. Retail unit costs (per kW) are higher than wholesale because of the included services and fractional nature. For example, a large wholesale might be, say, \$100-150 per kW-month plus energy costs, whereas a small retail could effectively be paying \$300+ per kW-month all-in (just as an illustrative number – actual prices vary widely by market).

Hybrid and Cloud Connectivity – Commercials: Many enterprises now adopt **hybrid cloud**, meaning they collocate some gear and use cloud for other workloads. Colos encourage this by offering direct cloud connectivity services. Commercially, an Equinix might sell a “Cloud Exchange Port” at a certain capacity for a fee, letting the customer connect to AWS/Azure, etc., through that exchange fabric. Cloud providers in turn reduce their egress fees if using these direct paths – e.g., AWS Direct Connect data transfer might be \$0.02/GB instead of \$0.09/GB from region to internet (just as an example). That can save a heavy-data user enormous sums. So, a pattern emerges: customers put certain systems in colo (perhaps for performance or data control reasons) and link to their cloud, saving on egress and achieving more consistent latency. The cost trade-off includes paying cross-connect/port fees but often those are far outweighed by egress savings for large volumes.

Interconnection and Egress: Cloud **egress** fees (charges for data leaving the cloud to the internet) are notoriously high and a pain point for many businesses ("data gravity" issues). Colocation can help by moving some systems out of the cloud to a neutral site or at least establishing private links that lower those fees. For example, if a company has to move 100 TB a month from cloud to on-prem, at \$0.09/GB that's ~\$9,000. Over a year, \$108k just in egress. Through a Direct Connect 1 Gbps port (which might cost a few hundred a month) and then \$0.02-\$0.03/GB, it could be say \$20k total – a big difference. This saving partially subsidizes the colocation cost.

Procurement Process: Enterprises typically go through RFPs (requests for proposals) when selecting a colocation provider – comparing costs (space, power, cross-connect, installation fees), locations, compliance, etc. Key negotiables are usually: power pricing, any included remote hands hours, caps on power rate increases (since electricity tariffs can vary), and flexibility for growth. Wholesale often has custom agreements with clauses for expansion, early termination (with penalties), etc. Retail is often a more standardized service contract with SLAs.

Edge and New Models: With edge deployments, some providers offer **pay-as-you-go** colo models or shorter terms (month-to-month for a rack, etc.), but those are not mainstream for larger deals.

Resilience Options and Cost: Redundancy tier might influence cost – e.g., a customer could choose to take an N+1 resiliency space vs. 2N space (some providers offer "different tier level halls"). If a customer is fine with slightly lower SLA (maybe because they have redundant sites), they might get a better price. Some wholesale providers now offer "**Zone**" **pricing** like cloud AZs – if you take space in two separate halls (two zones) with independent systems, they might guarantee one can fail but you still operate, similar to multi-AZ in cloud, but you'd run active-active. Not common but conceptually.

Network Services Procurement: To complete their solution, colo customers also procure network (IP transit, point-to-point circuits, etc.). Those are separate contracts with carriers but facilitated by the colo environment. The presence of **marketplaces** (e.g., Digital Realty's ServiceFabric or Equinix Fabric portals) allows ordering virtual connections quickly. These network costs become part of the total cost of hosting solution (and often cheaper than trying to bring the same networks to an on-prem site).

Cloud Procurement Contrast: To contrast, if a customer uses pure cloud instead of colo, they are on an OPEX fully variable model – paying per hour or per GB etc. This is flexible but can be more expensive at scale compared to a fixed colo cost. Colocation is more fixed-cost (lease-like) plus the responsibility of managing hardware. So often, critical stable workloads might be cheaper in colo (especially if company already owns hardware and can amortize it), whereas spiky or short-term workloads benefit from cloud's elasticity. So many organizations do hybrid for cost optimization – e.g., steady state baseline in colo, cloud bursting for peaks.

Commercial Model Trend: We've seen some shift to **power-based pricing** even for retail (i.e., selling by kW rather than per square foot), reflecting that power is the prime resource. Also, renewable energy options are part of procurement – some customers demand that their portion is powered by renewables (some providers offer renewable energy credits or direct green power options as part of contract, sometimes at slight premium or as a differentiator included).

Cross-Connect and Interconnect Charging Debate: Some have criticized cross-connect fees as exorbitant (the actual cost of a fiber patch is low). But providers defend it as part of facility revenue model that enables

building those meet-me ecosystems. In any case, customers should account for those costs, especially if they plan to connect to many partners (like if you put your infra in Equinix to connect with 10 others, you'll pay 10 cross-connect fees monthly).

In summary, **procurement** of data center services is about balancing the technical needs (space, power, connectivity, SLA) with the commercial terms (cost per kW, length of term, flexibility, included services). Wholesale vs retail is a core choice: *Wholesale* suits large requirements where the client can manage more and wants economy of scale, *Retail* suits smaller or distributed requirements or those that need rich connectivity out-of-the-box and perhaps more provider support. Many large companies use both: e.g., deploy big cloud nodes in wholesale colo but also maintain some smaller deployments in retail sites for network peering or specific use-cases.

11. Emerging Trends

The data center industry in mid-2020s is facing new **trends and challenges** that shape future designs and strategies. Some of the most impactful emerging trends include the rapid growth of AI/HPC workloads, the adoption of liquid cooling, constraints in power availability, and a push for greater sustainability and efficiency.

AI and High-Performance Computing (HPC) Loads: The boom in artificial intelligence – particularly training large machine learning models – has created *unprecedented demand* for computing power. AI training clusters deploy thousands of high-wattage GPUs (graphics processing units) or specialized AI accelerators. This leads to **much higher rack power densities** than traditional enterprise IT. Where an average rack was ~5 kW a few years ago, AI racks fully populated with GPUs can draw 30 kW or even 50+ kW. Such dense loads strain conventional data center power and cooling setups. The trend is forcing operators to rethink architecture: **power distribution** must support higher amperage per rack (e.g., use of 415V power to reduce current, larger gauge wiring, more powerful busways), and **cooling** often needs to transition to liquid methods (rear-door heat exchangers or direct liquid to chips) because air cooling for 50 kW would require impractically high airflow. Uptime Institute's 2024 survey indicates that while most data centers still have few racks above 30 kW, many operators expect that to change in coming years as they prepare for AI hardware (Uptime Institute, 2024) ²¹ ³⁹.

This AI wave also drives huge **scale** needs – new large language model training facilities may be tens of megawatts each just for one training cluster. Cloud providers and colocation operators are racing to build capacity for these. For example, Microsoft and NVIDIA announced plans for AI supercomputer infrastructure, and colo firms like Digital Realty and Equinix have formed partnerships to host "AI hubs" with ready power and cooling for GPU deployments. Another aspect is **low latency interconnects** – HPC clusters require very high-bandwidth, low-latency networking (like NVIDIA InfiniBand) between servers. Data center designs might need to accommodate high-density cabling and networking gear in new ways to support that east-west traffic.

Impact on sites: HPC loads can alter site power usage effectiveness (PUE). On one hand, high utilization of servers can increase cooling demand (thus raising PUE slightly if cooling isn't upgraded). On the other hand, if liquid cooling is used, it can actually improve efficiency of heat removal. There's a dynamic: if an HPC cluster can run at higher inlet temperatures (some liquid-cooled systems allow higher coolant temps), free cooling windows widen and chiller use drops, potentially *lowering* PUE. Also, HPC clusters often run near full capacity continuously (like training jobs run 24x7), meaning a data center's power draw will be very steady

at peak – this is different from enterprise IT which might idle often. Steady high load can sometimes allow more optimized cooling plant operation (no light load inefficiencies). But it definitely stresses backup systems and grid supply due to sustained draw.

Liquid Cooling Adoption: Tied to the above, there's a clear movement toward **liquid cooling solutions** for high-density racks. Several mainstream colocation providers have started offering liquid-cooled cabinet options or liquid-ready environments. For example, Equinix has tested rear-door cooling units and even immersion in some sites; Digital Realty has a venture with Chillydne (liquid cooling firm) to be able to cool chips directly. By 2025, it's expected perhaps ~20-30% of data center operators will have some form of liquid cooling in operation (though exact uptake depends on customer hardware deployments). Liquid cooling is not entirely new (mainframes had water cooling decades ago), but its reintroduction is significant for open compute hardware. The challenges revolve around managing water (or coolant) piping to racks, avoiding leaks, and retraining staff, but many see it as inevitable for future loads. ASHRAE in 2021 even added new liquid-cooled classes (W classes) for data centers to standardize approach ¹²¹ ¹²². The benefits: much higher heat capture per rack, potential reuse of heat (since fluid comes out hot – e.g., some European data centers pump waste heat to municipal heating, and with water cooling that becomes more feasible as the heat is already in water form). **Immersion cooling** remains niche but is being tried for certain blockchain or HPC use-cases; two-phase immersion particularly can achieve extreme densities but has overhead in managing the dielectric fluid environment.

Power Grid Constraints: As data centers proliferate, especially in certain hot markets, the **electric utility infrastructure is becoming a bottleneck**. This is an emerging trend seen in Northern Virginia, Dublin (Ireland), London, Singapore (which even had a moratorium on new DC builds 2019-2022). In Northern Virginia, Dominion Energy projected multi-year delays in fulfilling new data center power requests due to the volume – at one point, the wait for new connections in some parts (Loudoun/Prince William County) was said to be 5-7 years if immediate action isn't taken (Business Insider, 2022) ¹²³. Similarly, Silicon Valley Power had to pause new load acceptances.

This has several implications: **rising rental rates** in constrained markets (scarcity drives up cost of existing supply), pushes **development to new markets** (hence surge in Atlanta, Phoenix as alternatives), and innovation in how to get power. Some data center companies are exploring **on-site power generation** to bypass grid limits – e.g., deploying large gas turbines or fuel cell farms on-site (which still need fuel but lighten grid draw). In an opinion piece, it's argued to consider modular nuclear or off-grid solutions for AI clusters (DCD, 2025) ¹²⁴ ¹²⁵. Already, some data centers have done creative things like **temporarily using batteries to support grid** or planning to run generators during peak grid times (though emissions rules make this tricky). The bottom line: securing power has become the number one factor in data center site planning now, even above space. Power constraints are leading to **regional caps**: e.g., the Amsterdam area and Frankfurt both instituted limits on new data center developments until grid upgrades catch up, which is changing how companies plan European capacity.

Sustainability and Energy Efficiency: Under growing pressure from governments, customers, and investors, data center operators are striving for **sustainable operations**. Several facets here:

- **Renewable Energy:** Almost all major operators have commitments to use 100% renewable energy (either via direct power purchase agreements (PPAs) from wind/solar farms or via renewable energy credits). Google claims it has matched 100% of its energy with renewables since 2017, Equinix since 2020 across many regions, etc. The trend now is to go further: "24x7 carbon-free energy" – meaning

not just annual matching, but hour-by-hour localized renewable usage (Google is pushing this, trying to run each data center on carbon-free sources all the time by 2030). Data centers are among the largest corporate buyers of renewable power (they often sign long-term PPAs, e.g., Amazon is #1 corporate buyer of renewable energy globally as of 2023).

- **Energy Efficiency and Cooling Innovations:** As mentioned, average PUE has plateaued ~1.57, but new builds still push the envelope lower in favorable climates (some hyperscalers achieved PUE ~1.1 in cool regions using advanced cooling). For HPC, new metrics are considered like **Compute Efficiency** (to incorporate IT efficiency as well). Adopting liquid cooling also is a sustainability move – it can cut cooling energy by reducing fan power and allowing higher coolant temps meaning more economizer use ⁶⁰ ⁷¹. Another metric: **WUE (water usage)**. There's increasing scrutiny on water consumption – in some water-stressed regions, data centers may be asked to limit or eliminate water cooling. This is driving interest in **waterless cooling** solutions (like refrigerant-based chillers or direct-to-chip liquid that then uses dry coolers, etc.). Some operators commit to using greywater or investing in community water projects to offset usage.
- **Heat Reuse:** Especially in colder climates (Nordic Europe, Canada), waste heat reuse is an emerging practice – connecting data centers to district heating networks. For instance, in Denmark and Sweden there are data centers warming thousands of homes. In the U.S., this is less common due to distribution distances and lower need in many areas, but some projects exist (e.g., in Seattle, an Amazon campus attempted heat reuse from a nearby data center for heating office buildings). Heat reuse improves overall energy efficiency (reduces community's need to generate separate heating). Some jurisdictions may encourage or even mandate exploring it (the EU has proposed making large data centers assess feasibility of heat reuse).
- **Carbon and Reporting:** Data centers are being included in corporate ESG reporting. Some are adopting **carbon neutrality goals**, using carbon offsets for any emissions they can't eliminate (like diesel generator tests – some buy offsets for generator fuel emissions). Others focus on **Scope 3 emissions** – e.g., encouraging use of low-embodied-carbon building materials, recycling of servers, etc. Circular economy efforts (like IT hardware recycling, batteries recycling as second-life) are trending.
- **Regulation:** While few direct data center-specific laws exist in U.S., places like **EU** have Code of Conduct for Data Centres (voluntary but might become baseline), and certain states have high efficiency standards in building codes (e.g., ASHRAE 90.4 energy standard is adopted by some states requiring certain efficiency levels for cooling and power systems). If grids become more stressed, governments might impose limits on new data centers (as Singapore did temporarily). Also climate policies could affect data centers – e.g., carbon pricing or limits on diesel generator running hours beyond emergencies.

Edge Computing & 5G: Another trend is the move toward smaller edge deployments due to 5G and IoT. While not as massive in impact as AI or sustainability, edge computing is gradually expanding. Companies like EdgeConneX, Vapor IO, American Tower are deploying micro-data centers at cell tower sites or regional locations to shorten latency for certain applications (augmented reality, autonomous systems, gaming, etc.). This trend will likely result in a more distributed layer of infrastructure complementing big central clouds. One concrete example: AWS has "Wavelength" edge zones collocated in telco networks to provide <10 ms latency for mobile apps. These edge sites won't be huge power consumers individually, but

collectively could be significant. They also might use more **autonomous operations** (since unmanned sites).

Automation and AI Ops: Data center operations themselves are adopting AI – for instance, using machine learning to optimize cooling (Google's DeepMind system reduced cooling energy by 30% in their DCs by AI recommendations), or predictive maintenance (AI analyzing sensor data to predict equipment failures). Also, more automation in monitoring and management (DCIM tools getting smarter). This trend, sometimes termed "AI for AI", will help manage complexity as infrastructure grows. However, an Uptime Institute survey found many operators are cautious about fully autonomous data centers – trust in AI operations was actually declining, with fewer than half feeling confident in AI-driven management ¹²⁶ ¹²⁷. So humans are still very much in the loop.

Security and Edge Regulation: Another emerging aspect is physical and cyber security as national security issues. There's discussion in some countries about treating large data centers as critical infrastructure (with corresponding regulations for security and resilience). Already in U.S., there's CFIUS oversight on foreign investment in data centers, and some prohibitions (like Chinese companies limited in buying near sensitive locations). As data centers house more critical data (including government cloud, etc.), expect more compliance requirements.

Conclusion of Trends: In summary, the industry is pivoting to handle *much denser and more power-hungry computing*, while contending with external constraints like power grid limits and environmental accountability. Solutions include technical innovation (liquid cooling, new architectures), strategic geographic shifts (finding places with available power and land), and operational changes (sustainability commitments, advanced automation). Providers that adapt to these trends – enabling AI hardware, demonstrating green credentials, collaborating with utilities – will lead in the next phase of growth. Meanwhile, ignoring these could leave facilities obsolete or non-competitive (e.g., a data center that can't handle >10 kW/rack might struggle to attract new customers in a few years).

The macro context: demand for data center capacity is still increasing briskly (digital transformation, AI, streaming, etc.), but how that capacity is delivered is evolving. We might see more **modular deployments** (to accelerate builds), maybe resurgence of **on-site power generation** (like fuel cells, as Microsoft tested 3MW hydrogen fuel cell system in 2022), and **collaboration with governments** to ensure sustainable growth (like Virginia continuing to invest in grid due to data center tax revenue incentive). The next 5-10 years will likely bring greener, higher-density, and more widely distributed data center infrastructure to underpin our connected world.

Glossary

- **Colocation (Colo):** The practice of renting space in a third-party data center to host your servers or IT equipment. The provider supplies power, cooling, physical security, and network access, while the customer provides and manages the IT hardware.
- **Hyperscale Data Center:** An extremely large facility (tens of megawatts of IT load or more) designed for a single hyperscale company's services (like cloud providers or social media giants). Characterized by massive scale, high uniformity, and automation. Hyperscale also implies the ability to rapidly expand capacity.

- **Edge Data Center:** A small data center located close to end-users or devices to provide low-latency services. It often has limited capacity compared to core sites and may be in unconventional locations (cell tower base, on-premises at factories, etc.). Supports edge computing for applications like IoT, 5G, AR/VR.
- **PUE (Power Usage Effectiveness):** A metric for data center energy efficiency. Defined as *Total Facility Power / IT Equipment Power*. A PUE of 1.5 means for every 1 W used by IT gear, 0.5 W is used for cooling, power overhead, etc. Ideal PUE is 1.0 (no overhead). Lower PUE = more efficient.
- **WUE (Water Usage Effectiveness):** A metric for water efficiency, measuring *liters of water used per kWh of IT energy*. Tracks how much water cooling consumes. Lower WUE is better (zero means no water used in cooling).
- **Tier III / Tier IV:** Levels from Uptime Institute's tier classification. *Tier III* data centers have N+1 redundancy and allow concurrent maintenance (99.982% uptime goal). *Tier IV* adds fault tolerance with 2N redundancy (99.995% uptime). Higher tiers incur greater cost for extra backup systems.
- **UPS (Uninterruptible Power Supply):** A battery-backed power unit that provides short-term power during an outage and conditions incoming power. In data centers, typically a large system that instantaneously takes over if utility power fails, bridging until generators start.
- **Generator (Genset):** On-site engine (usually diesel) generator used for backup power. Kicks in during a utility outage to supply electricity to the data center, usually within seconds, and can run the facility as long as fuel is available.
- **Meet-Me Room (MMR):** A secure room in a data center where telecom carriers and customers interconnect. Houses fiber patch panels and network gear for cross-connects. It's the nexus of connectivity, enabling different networks and tenants to "meet" and exchange traffic directly.
- **Cross-Connect:** A physical cable connecting two parties within the data center (e.g., from a customer's rack to a carrier's port in the MMR). Facilitates private network connections. Typically managed by the data center staff and subject to a monthly fee.
- **Availability Zone (AZ):** A term from cloud computing; a distinct location comprised of one or more data centers, with independent power/cooling/network, within a cloud **Region**. AZs are isolated from each other to prevent one failure impacting others, allowing cloud applications to be deployed with high availability by spanning AZs.
- **SOC 2:** A security and availability audit framework (under AICPA) for service organizations. Data center operators undergo SOC 2 Type II audits to verify they have effective controls for security, availability, processing integrity, confidentiality, and privacy. Many enterprise customers require a SOC 2 report from their colo provider.
- **Liquid Cooling:** Cooling methods that use liquids (water, coolant, etc.) to remove heat directly from IT equipment, as opposed to traditional air cooling. Examples: direct-to-chip cold plates, immersion

cooling, rear-door cooling. Liquid's higher thermal conductivity allows cooling very high-density loads more effectively.

Reading List

1. **Uptime Institute (2024) - Global Data Center Survey 2024 Report.** Insights into industry performance, efficiency (average PUE ~1.56), rack density trends, outages analysis, and sustainability metrics [65](#) [21](#).
2. **Uptime Institute - Tier Standard Overview.** Explains Tier I-IV infrastructure criteria and expected availabilities [86](#) [41](#), helpful for understanding reliability design levels in data centers.
3. **U.S. Chamber of Commerce (2017) - "Data Centers: Jobs and Opportunities" Report.** Industry landscape and definitions (enterprise vs colocation), growth drivers, plus foundational definition of data centers [4](#) [7](#).
4. **PhoenixNAP (2021) - "Data Center Tiers Explained".** Plain-language explanation of Tier 1-4 characteristics and comparison table of redundancies and downtime per year [86](#) [41](#).
5. **Equinix (2024) - Blog: "What Is Water Usage Effectiveness (WUE) in Data Centers?"** Discusses the water-energy trade-off in cooling and why measuring WUE alongside PUE is important [60](#) [71](#).
6. **DCD (2025) - "CBRE: Vacancy rates in top data center markets hit record low".** News on data center market capacity and demand, noting NoVA ~2,930 MW capacity and rapid growth in Atlanta and Phoenix [100](#) [107](#).
7. **Energy EESI (2025) - "Data Centers and Water Consumption".** Highlights water use statistics: large facilities using up to 5 million gallons/day and U.S. data centers ~449 million gallons/day (2021) [128](#) [74](#); also links AI-driven demand to rising water and energy use.
8. **CoreSite (2022) - "Is the time right for lithium-ion batteries in data centers?"** Covers Li-ion vs VRLA battery cost and performance, noting Li-ion's 15-year lifespan vs 5-year VRLA, smaller footprint, faster recharge [53](#) [48](#), and prediction that 35% of data center UPS batteries could be Li-ion by 2025 [51](#).
9. **Flexential (2024) - "Essential considerations for effective data center site selection".** Comprehensive blog on site selection factors: power grid reliability, fiber proximity, climate, natural disaster risks, and expansion room [26](#) [24](#).
10. **Digital Realty (2025) - Cross Connect Product Brief.** Describes cross-connects as physical connections in meet-me-rooms enabling fast, low-latency links between customers and carriers [22](#) [83](#), and emphasizes interconnection benefits (security, performance).
11. **Dgtl Infra (2024) - "Top 250 Data Center Companies" List.** Ranks major operators by footprint. E.g., Equinix operates 251 data centers (30.2 million sq ft) [113](#), Digital Realty 312 data centers globally [114](#), and NTT GDC 95 data centers (1.1+ GW) [115](#), providing context on leading players.

12. **Dgtl Infra (2024) – “Data Center Compliance: Standards and Key Requirements”.** Reviews compliance regimes like ISO 27001, SOC 2, PCI DSS, HIPAA and why they matter [95](#) [94](#), plus practices for security and continuity (e.g., NFPA 75 for fire protection).
13. **DCD (2025) – Opinion: “Bridging the power gap: Why data centers can’t wait for the grid”.** Commentary on power delivery delays for AI data centers and suggestion of off-grid modular power solutions; notes even NoVA and Santa Clara reaching capacity limits, with projects queued behind years of utility upgrades [3](#) [124](#).
14. **ASHRAE Journal (2021) – “Thermal Guidelines for Data Processing Environments, 5th Ed.”** (Summary/Excerpt). Details recommended temperature/humidity ranges for IT equipment classes (A1 to A4) [78](#) [81](#) and the introduction of liquid cooling classes, guiding modern cooling designs.
15. **EPA Report to Congress (2007) – “Server and Data Center Energy Efficiency”.** (Historical reference) Landmark study (EPA 2007) quantifying data center energy use at ~1.5% of U.S. electricity then, spurring industry efficiency efforts and initiatives like The Green Grid’s PUE metric development [129](#) [69](#). (Provides background on why metrics like PUE were adopted).

-
- [1](#) [2](#) [18](#) [21](#) [38](#) [39](#) [64](#) [65](#) [66](#) [67](#) [68](#) [69](#) [126](#) [127](#) datacenter.uptimeinstitute.com
<https://datacenter.uptimeinstitute.com/rs/711-RIA-145/images/2024.GlobalDataCenterSurvey.Report.pdf?version=0>
- [3](#) [124](#) [125](#) [Bridging the power gap: Why data centers can’t wait for the grid - DCD](#)
<https://www.datacenterdynamics.com/en/opinions/bridging-the-power-gap-why-data-centers-cant-wait-for-the-grid/>
- [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) uschamber.com
https://www.uschamber.com/assets/documents/ctec_datacenterrrpt_lowres.pdf
- [11](#) [Availability Zones - AWS Fault Isolation Boundaries](#)
<https://docs.aws.amazon.com/whitepapers/latest/aws-fault-isolation-boundaries/availability-zones.html>
- [12](#) [Question from CLF-C02 Practice Exam - Availability Zone - Reddit](#)
https://www.reddit.com/r/AWSCertifications/comments/191wql/question_from_clfc02_practice_exam_availability/
- [13](#) [Understanding the Differences Between 5 Common Types of Data ...](#)
<https://www.datacenterfrontier.com/sponsored/article/11427373/belden-understanding-the-differences-between-5-common-types-of-data-centers>
- [14](#) [Types of Data Centers: Enterprise, Colocation, Hyperscale - Dgtl Infra](#)
<https://dgtlinfra.com/types-of-data-centers/>
- [15](#) [16](#) [19](#) [Difference Between Retail and Wholesale Colocation | Volico](#)
<https://www.volico.com/difference-between-retail-and-wholesale-colocation/>
- [17](#) [Hyperscale vs. Colocation - Interconnections - The Equinix Blog](#)
<https://blog.equinix.com/blog/2020/08/27/hyperscale-vs-colocation/>
- [20](#) [It's Time To Learn About Latency - TeleGeography Blog](#)
<https://blog.telegeography.com/its-time-to-learn-about-latency>
- [22](#) [23](#) [83](#) [Cross Connect | Digital Realty](#)
<https://www.digitalrealty.com/platform-digital/connectivity/interconnection/cross-connect>

- [24 25 26 27 28 29 30 Essential considerations for effective data center site selection | Flexential](https://www.flexential.com/resources/blog/essential-considerations-effective-data-center-site-selection)
<https://www.flexential.com/resources/blog/essential-considerations-effective-data-center-site-selection>
- [31 Data Center Retail Sales & Use Tax Exemption | Virginia Economic Development Partnership](https://www.vedp.org/incentive/data-center-retail-sales-use-tax-exemption)
<https://www.vedp.org/incentive/data-center-retail-sales-use-tax-exemption>
- [32 45 46 47 59 Understanding Key Elements of Data Center Power Distribution](https://www.tierpoint.com/blog/data-center-power-distribution)
<https://www.tierpoint.com/blog/data-center-power-distribution/>
- [33 123 Data Centers Face Seven-Year Wait for Dominion Power Hookups](https://www.energyconnects.com/news/utilities/2024/august/data-centers-face-seven-year-wait-for-dominion-power-hookups/)
<https://www.energyconnects.com/news/utilities/2024/august/data-centers-face-seven-year-wait-for-dominion-power-hookups/>
- [34 US Data Center Trends & Why They Require New Dark Fiber](https://www.7x24exchange.org/us-data-center-trends-why-they-require-new-dark-fiber/)
<https://www.7x24exchange.org/us-data-center-trends-why-they-require-new-dark-fiber/>
- [35 36 37 94 95 96 98 99 Data Center Compliance: Standards and Key Requirements - Dgtl Infra](https://dgtlinfra.com/data-center-compliance-standards/)
<https://dgtlinfra.com/data-center-compliance-standards/>
- [40 41 42 43 44 86 87 88 89 90 91 92 93 Data Center Tiers Classification Explained: \(Tier 1, 2, 3, 4\)](https://phoenixnap.com/blog/data-center-tiers-classification)
<https://phoenixnap.com/blog/data-center-tiers-classification>
- [48 49 50 51 53 58 Is the Time Right for Lithium Ion Batteries in Data Centers?](https://www.coresite.com/blog/is-the-time-right-for-lithium-ion-batteries-in-data-centers)
<https://www.coresite.com/blog/is-the-time-right-for-lithium-ion-batteries-in-data-centers>
- [52 54 55 56 57 Lithium Ion Batteries for the data center. Are they ready for production yet? - Uptime Institute Blog](https://journal.uptimeinstitute.com/lithium-ion-batteries-in-the-data-center)
<https://journal.uptimeinstitute.com/lithium-ion-batteries-in-the-data-center/>
- [60 61 62 63 71 77 What Is Water Usage Effectiveness \(WUE\) in Data Centers? - Interconnections - The Equinix Blog](https://blog.equinix.com/blog/2024/11/13/what-is-water-usage-effectiveness-wue-in-data-centers)
[https://blog.equinix.com/blog/2024/11/13/what-is-water-usage-effectiveness-wue-in-data-centers/](https://blog.equinix.com/blog/2024/11/13/what-is-water-usage-effectiveness-wue-in-data-centers)
- [70 Data Center Water Usage: A Comprehensive Guide - Dgtl Infra](https://dgtlinfra.com/data-center-water-usage)
[https://dgtlinfra.com/data-center-water-usage/](https://dgtlinfra.com/data-center-water-usage)
- [72 73 74 75 76 128 Data Centers and Water Consumption | Article | EESI](https://www.eesi.org/articles/view/data-centers-and-water-consumption)
<https://www.eesi.org/articles/view/data-centers-and-water-consumption>
- [78 79 80 81 129 Understanding Environmental Conditions for Data Center Cooling - Therma](https://www.therma.com/data-center-cooling-overview-and-guide)
[https://www.therma.com/data-center-cooling-overview-and-guide/](https://www.therma.com/data-center-cooling-overview-and-guide)
- [82 What Is a Meet-Me Room?](https://www.sunbirddcim.com/glossary/meet-me-room)
<https://www.sunbirddcim.com/glossary/meet-me-room>
- [84 85 113 114 115 116 117 Top 250 Data Center Companies in the World as of 2024 - Dgtl Infra](https://dgtlinfra.com/top-data-center-companies)
[https://dgtlinfra.com/top-data-center-companies/](https://dgtlinfra.com/top-data-center-companies)
- [97 Data Center Operations | IBX Standards and Compliance - Equinix](https://www.equinix.com/data-centers/design/standards-compliance)
<https://www.equinix.com/data-centers/design/standards-compliance>
- [100 107 108 109 110 CBRE: Vacancy rates in top data center markets hit record low - DCD](https://www.datacenterdynamics.com/en/news/cbre-vacancy-rates-in-top-data-center-markets-hit-record-low)
<https://www.datacenterdynamics.com/en/news/cbre-vacancy-rates-in-top-data-center-markets-hit-record-low/>
- [101 102 103 104 111 112 Virginia Data Center Subsidy Costs Balloon by 1,051% - Good Jobs First](https://goodjobsfirst.org/virginia-data-center-subsidy-costs-balloon-by-1051)
[https://goodjobsfirst.org/virginia-data-center-subsidy-costs-balloon-by-1051/](https://goodjobsfirst.org/virginia-data-center-subsidy-costs-balloon-by-1051)

- ¹⁰⁵ Loudon Tax Incentives - SDIA Knowledge Hub
<https://knowledge.sdialliance.org/loudon-tax-incentives-47dc384b526341d0bed00685162c7bbf>
- ¹⁰⁶ Data Centers | Loudoun County Economic Development
<https://loudounpossible.com/business-sector/data-centers>
- ¹¹⁸ ¹²⁰ Q4 2024 data center colocation results: Digital Realty, Equinix, and ...
<https://www.datacenterdynamics.com/en/news/q4-2024-data-center-colocation-results-digital-realty-equinix-and-iron-mountain/>
- ¹¹⁹ Data Center Colocation Company Evaluation Report 2025 | Equinix ...
<https://www.businesswire.com/news/home/20250903529641/en/Data-Center-Colocation-Company-Evaluation-Report-2025-Equinix-Digital-Realty-and-NTT-Lead-in-Data-Center-Expansion---ResearchAndMarkets.com>
- ¹²¹ ¹²² ASHRAE publishes liquid cooling guidelines as chip power moves ...
<https://www.datacenterdynamics.com/en/news/ashrae-publishes-liquid-cooling-guidelines-as-chip-power-moves-into-uncharted-territory/>