# AXIS 6: Power & Cooling (Comparative Analysis)

**Presentation Format**: 30-45 min presentation (through Google Meet)

- **Mission**:
  - Explain the power and thermal challenges of modern AI infrastructure
  - Compare cooling technologies and analyze datacenter design tradeoffs
  - YOU CAN USE CHATGPT but always add source to every statements you say
    - If unsure of statements (because no sources found), put in italic in a different colors
  - Some data in tables below were not checked. If you decide to use them, feel free but it needs sources OR a formula that explains how to get the value

## Part A: The Power Trajectory of AI Systems

### GPU Power Evolution

- How has GPU power consumption evolved?
  - What is driving the power increases?
  - Is there a ceiling to GPU power?
  - How does power scale with performance?

| GPU | Year | TDP (W) | FP16 TFLOPs | W/TFLOP | Process Node |
|-----|------|---------|-------------|---------|--------------|
| V100 | 2017 | ? | ? | ? | ? |
| A100 | 2020 | ? | ? | ? | ? |
| H100 SXM | 2022 | ? | ? | ? | ? |
| H200 SXM | 2024 | ? | ? | ? | ? |
| B200 SXM | 2024 | ? | ? | ? | ? |
| B300 | 2025 | ? | ? | ? | ? |
| R100 (Rubin) | 2026 | ? | ? | ? | ? |

- System-level power (full node):

| System | GPUs | GPU Power | Total System Power | Year |
|--------|------|-----------|--------------------|------|
| DGX-1 | 8× V100 | ? | ? | 2017 |
| DGX A100 | 8× A100 | ? | ? | 2020 |
| DGX H100 | 8× H100 | ? | ? | 2022 |
| DGX B200 | 8× B200 | ? | ? | 2024 |
| GB200 NVL72 | 72× B200 | ? | ? | 2024 |

- Rack power density evolution:

| Era | Typical Rack Power | kW/rack | Cooling Method |
|---|---|---|---|
| Traditional IT (2010) | ? | ? | Air |
| GPU compute (2018) | ? | ? | Air |
| AI training (2022) | ? | ? | Air/DLC |
| AI training (2024) | ? | ? | DLC required |
| AI training (2026) | ? | ? | DLC/Immersion |

- What drives power increases?

---

## Part B: GPU Power Delivery

### How does GPU power delivery work?

- Power delivery hierarchy
    - From grid to chip: what are the conversion stages?
    - Where are the efficiency losses?
    - How is power distributed within a GPU?

```
1    Grid (AC) → Transformer → UPS → PDU → PSU → VRM → GPU Die
```

| Stage | Input | Output | Typical Efficiency | Loss |
|---|---|---|---|---|
| Utility transformer | ? | ? | ? | ? |
| UPS | ? | ? | ? | ? |
| PDU | ? | ? | ? | ? |
| PSU (AC-DC) | ? | ? | ? | ? |
| VRM (DC-DC) | ? | ? | ? | ? |

- Voltage Regulator Modules (VRMs)
    - What is a VRM and why is it critical?
    - What are the phases and how do they work?
    - Why is VRM efficiency important?
    - What are the thermal challenges?

| VRM Aspect | Description | Typical Value |
|---|---|---|
| Input voltage | ? | ? |
| Output voltage | ? | ? |
| Phase count | ? | ? |
| Efficiency | ? | ? |

| VRM Aspect | Description | Typical Value |
|---|---|---|
| Power loss (1kW GPU) | ? | ? |

- 12V vs 48V power distribution
  - Why is 48V better for high-power systems?
  - What is the current reduction benefit?
  - Who is pushing 48V adoption?
  - What are the challenges?

| Aspect | 12V Distribution | 48V Distribution |
|---|---|---|
| Current for 1kW | ? | ? |
| Cable thickness | ? | ? |
| I²R losses | ? | ? |
| Connector size | ? | ? |
| Industry adoption | ? | ? |

- 12VHPWR connector issues
  - What is the 12VHPWR connector?
  - What problems occurred?
  - What is 12V-2x6 (the replacement)?
  - How much power can these connectors handle?

| Connector | Max Power | Pins | Issues |
|---|---|---|---|
| 8-pin PCIe | ? | ? | ? |
| 12VHPWR | ? | ? | ? |
| 12V-2x6 | 600a?W | ? | ? |

- 800V DC for AI factories
  - Why is NVIDIA pushing 800V DC?
  - What efficiency gains are possible?
  - How does this change datacenter design?
  - What are the safety considerations?

---

# Part C: Cooling Technologies

## Air Cooling

- How does air cooling work?
  - Heat sink design principles

- Airflow management (hot aisle/cold aisle)
- Fan power consumption
- What are the limits?

| Air Cooling Aspect | Typical Value | Limit |
|---|---|---|
| Max heat dissipation per GPU | ? | ? |
| Max rack density | ? | ? |
| Airflow per rack | ? | ? |
| Fan power overhead | ? | ? |

- When does air cooling fail?
  - At what TDP does air become impractical?
  - What are the acoustic limits?
  - How does altitude affect air cooling?

# Direct Liquid Cooling (DLC)

- How does DLC work?
  - Cold plate design and attachment
  - Manifold and distribution systems
  - Coolant types and flow rates
  - Heat rejection (CDU, dry coolers, cooling towers)

| DLC Component | Function | Key Specifications |
|---|---|---|
| Cold plate | ? | ? |
| Manifold | ? | ? |
| Quick disconnects | ? | ? |
| CDU (Coolant Distribution Unit) | ? | ? |
| Facility water loop | ? | ? |

- DLC performance characteristics:

| Aspect | Air Cooling | Direct Liquid Cooling |
|---|---|---|
| Max GPU TDP supported | ? | ? |
| Max rack density | ? | ? |
| PUE impact | ? | ? |
| Maintenance complexity | ? | ? |
| Capital cost | ? | ? |
| Operating cost | ? | ? |

- Coolant types:

| Coolant | Thermal Properties | Cost | Safety | Use Case |
|---|---|---|---|---|
| Water/glycol | ? | ? | ? | ? |
| Propylene glycol | ? | ? | ? | ? |
| Dielectric fluids | ? | ? | ? | ? |

## Immersion Cooling

- How does immersion cooling work?
    - Single-phase vs two-phase
    - Tank design and fluid management
    - Heat rejection methods
    - Maintenance considerations

| Aspect | Single-Phase Immersion | Two-Phase Immersion |
|---|---|---|
| Fluid type | ? | ? |
| Operating principle | ? | ? |
| Max heat flux | ? | ? |
| Fluid cost | ? | ? |
| Complexity | ? | ? |
| Maturity | ? | ? |

- Immersion cooling advantages and challenges:

| Advantage | Challenge |
|---|---|
| Highest heat density | ? |
| No fans required | ? |
| Reduced PUE | ? |
| Component longevity | ? |

## Rear-Door Heat Exchangers (RDHx)

- What is RDHx?
    - How does it supplement air cooling?
    - What densities can it support?
    - When is it the right choice?

| RDHx Type | Cooling Capacity | Best For |
|---|---|---|
| Passive RDHx | ? | ? |
| Active RDHx | ? | ? |

---

# Part D: Power Usage Effectiveness (PUE)

## Understanding PUE

- What is PUE?
    - Definition: Total Facility Power / IT Equipment Power
    - What does PUE measure and not measure?
    - What are the components of overhead?

```
1    PUE = (IT Load + Cooling + Power Distribution + Lighting + Other) / IT Load
```

| PUE Component | Typical % of Overhead | Reduction Strategies |
|---|---|---|
| Cooling | ? | ? |
| Power distribution losses | ? | ? |
| Lighting and other | ? | ? |

- PUE benchmarks:

| Datacenter Type | Typical PUE | Best-in-Class PUE |
|---|---|---|
| Legacy enterprise | ? | ? |
| Modern enterprise | ? | ? |
| Hyperscale (air) | ? | ? |
| Hyperscale (DLC) | ? | ? |
| AI-optimized | ? | ? |

- How does cooling choice affect PUE?

| Cooling Method | Typical PUE | Why |
|---|---|---|
| Traditional air (CRAC) | ? | ? |
| Hot/cold aisle containment | ? | ? |
| Free air cooling | ? | ? |
| Direct liquid cooling | ? | ? |

| Cooling Method | Typical PUE | Why |
|---|---|---|
| Immersion cooling | ? | ? |

- Best-in-class examples:

| Company | Facility | PUE | How Achieved |
|---|---|---|---|
| Google | ? | ? | ? |
| Meta | ? | ? | ? |
| Microsoft | ? | ? | ? |
| NVIDIA DGX Cloud | ? | ? | ? |

- Beyond PUE: other efficiency metrics
    - What is WUE (Water Usage Effectiveness)?
    - What is CUE (Carbon Usage Effectiveness)?
    - Why do these matter for AI datacenters?

| Metric | Definition | Typical Values | Best-in-Class |
|---|---|---|---|
| PUE | ? | ? | ? |
| WUE | ? | ? | ? |
| CUE | ? | ? | ? |

---

# Part E: Infrastructure Requirements

## Power Density Planning

- Rack power density considerations
    - How do you plan for increasing density?
    - What infrastructure upgrades are needed?
    - How do you handle mixed densities?

| Density Tier | kW/Rack | Infrastructure Requirements |
|---|---|---|
| Low density | ? | ? |
| Medium density | ? | ? |
| High density | ? | ? |
| Ultra-high density | ? | ? |

- Power distribution architecture:

| Component | Function | Sizing Consideration |
|---|---|---|
| Utility feed | ? | ? |
| Main switchgear | ? | ? |
| UPS systems | ? | ? |
| PDUs | ? | ? |
| RPPs (Remote Power Panels) | ? | ? |

## Cooling Capacity Planning

- Cooling load calculations
  - How do you size cooling for AI racks?
  - What is N+1 redundancy?
  - How much water is needed for DLC?

| Cooling Capacity Unit | Conversion | Context |
|---|---|---|
| 1 ton of cooling | ? BTU/hr | ? kW |
| 1 kW IT load (air) | ? tons | Includes overhead |
| 1 kW IT load (DLC) | ? tons | Direct rejection |

- Water requirements for liquid cooling:

| Cooling Method | Water Usage | GPM per MW |
|---|---|---|
| Evaporative (cooling tower) | ? | ? |
| Dry cooler | ? | ? |
| DLC (closed loop) | ? | ? |

## Backup Power Systems

- UPS sizing and architecture
  - What UPS topologies exist?
  - How long does UPS need to last?
  - What is rotary vs battery UPS?

| UPS Type | Efficiency | Runtime | Best For |
|---|---|---|---|
| Double conversion | ? | ? | ? |
| Line interactive | ? | ? | ? |
| Rotary UPS | ? | ? | ? |
| Battery + flywheel | ? | ? | ? |

- Generator requirements:

| Facility Size | Generator Capacity | Fuel Storage | Startup Time |
|---|---|---|---|
| 10 MW | ? | ? | ? |
| 100 MW | ? | ? | ? |
| 500 MW | ? | ? | ? |
| 1 GW | ? | ? | ? |

## Grid Connection for GW-Scale Facilities

- What does a GW-scale AI datacenter need?
  - Substation requirements
  - Transmission line upgrades
  - Grid stability considerations
  - Timeline for new connections

| Scale | Grid Requirements | Typical Lead Time |
|---|---|---|
| 50 MW | ? | ? |
| 200 MW | ? | ? |
| 500 MW | ? | ? |
| 1 GW+ | ? | ? |

- Power sourcing strategies:

| Strategy | Description | Pros | Cons |
|---|---|---|---|
| Grid connection | ? | ? | ? |
| On-site generation | ? | ? | ? |
| PPA (Power Purchase Agreement) | ? | ? | ? |
| Behind-the-meter solar/wind | ? | ? | ? |
| Nuclear (SMR) | ? | ? | ? |

---

# Part F: Deep Dive Topics

## Chip-Level Power Management

- Dynamic Voltage and Frequency Scaling (DVFS)
  - How does DVFS work?
  - What is the power/frequency relationship?

- How do GPUs implement DVFS?

| Power State | Voltage | Frequency | Power | Use Case |
|---|---|---|---|---|
| Max boost | ? | ? | ? | Peak compute |
| Base clock | ? | ? | ? | Sustained |
| Idle | ? | ? | ? | Low utilization |
| Sleep | ? | ? | ? | Inactive |

- Power gating
    - What is power gating?
    - What can be gated (cores, memory, I/O)?
    - What are the wake-up latencies?

## Thermal Throttling Behavior

- How do GPUs throttle under thermal stress?
    - Temperature thresholds
    - Throttling mechanisms
    - Performance impact

| Threshold | Temperature | Action |
|---|---|---|
| Target | ~83°C | ? |
| Throttle start | ~85°C | ? |
| Max operating | ~90°C | ? |
| Shutdown | ~95°C | ? |

## Heat Sink and Cold Plate Design

- Heat sink design principles:

| Parameter | Impact | Tradeoff |
|---|---|---|
| Fin density | ? | ? |
| Base thickness | ? | ? |
| Heat pipe count | ? | ? |
| Material (Cu vs Al) | ? | ? |

- Cold plate design for GPUs:

| Design Aspect | Consideration | Best Practice |
|---|---|---|
| Contact area | ? | ? |

| Design Aspect | Consideration | Best Practice |
|---|---|---|
| Channel design | ? | ? |
| Flow rate | ? | ? |
| Pressure drop | ? | ? |

## Stranded Power in Datacenters

- What is stranded power?
    - Why does it occur?
    - How much power is typically stranded?
    - How do you minimize stranded capacity?

---

# Part G: Companies & Industry Landscape

## System Vendors

| Company | Products | Cooling Approach | Market Position |
|---|---|---|---|
| NVIDIA | DGX, MGX, HGX | Air + DLC ready | ? |
| Dell | PowerEdge XE | ? | ? |
| HPE | Cray EX | ? | ? |
| Supermicro | GPU servers | ? | ? |
| Lenovo | ThinkSystem | ? | ? |

## Cooling Infrastructure Vendors

| Company | Products | Technology Focus |
|---|---|---|
| Vertiv | Liebert, CDUs | Full stack cooling |
| Schneider Electric | APC, cooling | Power + thermal |
| Asetek | Cold plates, CDUs | DLC pioneer |
| CoolIT | DLC systems | Rack-level DLC |
| GRC | ICEraQ | Single-phase immersion |
| LiquidCool | Immersion tanks | Two-phase immersion |
| Submer | SmartPod | Immersion systems |

## Power Infrastructure Vendors

| Company | Products | Specialty |
|---|---|---|
| Schneider Electric | UPS, PDUs, switchgear | End-to-end power |
| Vertiv | UPS, PDUs | Critical power |
| Eaton | UPS, PDUs | Power distribution |
| ABB | Transformers, switchgear | Utility-scale |
| Caterpillar | Generators | Backup power |

---

# Part H: Summary Comparison

## Cooling Technology Comparison

| Aspect | Air Cooling | Rear-Door HX | Direct Liquid | Single-Phase Immersion | Two-Phase Immersion |
|---|---|---|---|---|---|
| Max kW/rack | ? | ? | ? | ? | ? |
| PUE achievable | ? | ? | ? | ? | ? |
| Capital cost | ? | ? | ? | ? | ? |
| Operating cost | ? | ? | ? | ? | ? |
| Maintenance | ? | ? | ? | ? | ? |
| Maturity | ? | ? | ? | ? | ? |
| GPU compatibility | ? | ? | ? | ? | ? |

## AI System Power Summary

| System | GPUs | Total Power | Cooling Method | Rack Density |
|---|---|---|---|---|
| DGX A100 | 8× A100 | ~6.5 kW | Air | ? |
| DGX H100 | 8× H100 | ~10.2 kW | Air/DLC | ? |
| DGX B200 | 8× B200 | ~14.3 kW | DLC | ? |
| GB200 NVL72 | 72× B200 | ~120 kW | DLC | ? |
| AMD MI300X (8-way) | 8× MI300X | ? | ? | ? |
| Google TPU v5p pod | ? | ? | ? | ? |

## Datacenter Efficiency Comparison

| Operator | Facility Type | PUE | WUE | Cooling Method |
|---|---|---|---|---|
| Google | Hyperscale | ? | ? | ? |

| Operator | Facility Type | PUE | WUE | Cooling Method |
|---|---|---|---|---|
| Meta | Hyperscale | ? | ? | ? |
| Microsoft | Azure | ? | ? | ? |
| AWS | Cloud | ? | ? | ? |
| CoreWeave | AI-focused | ? | ? | ? |
| Lambda Labs | AI-focused | ? | ? | ? |

## TCO Impact of Cooling Choice

| Cost Component | Air Cooling | DLC | Immersion |
|---|---|---|---|
| Capital ($/kW IT) | ? | ? | ? |
| Power cost ($/kW-yr) | ? | ? | ? |
| Maintenance ($/kW-yr) | ? | ? | ? |
| Floor space ($/kW-yr) | ? | ? | ? |
| 5-year TCO ($/kW) | ? | ? | ? |