# Internet-Scale Service Infrastructure Efficiency

## International Symposium on System Architecture

James Hamilton, 2009/6/23

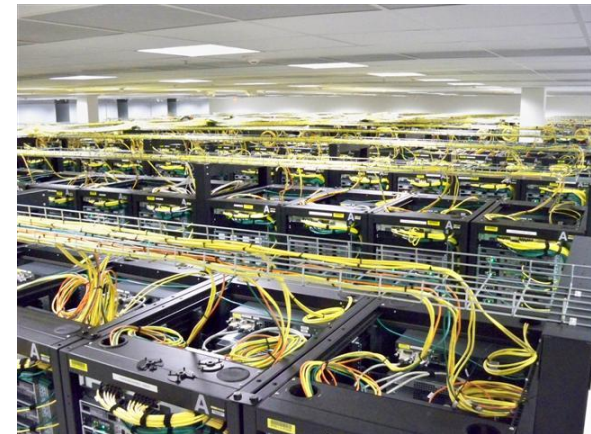**VP & Distinguished Engineer, Amazon Web Services**

e: James@amazon.com

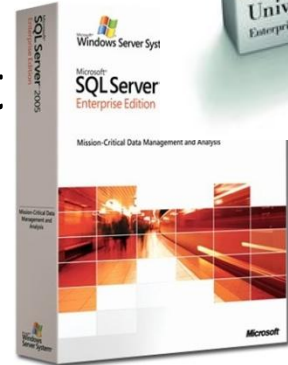w: mvdirona.com/jrh/work

b: perspectives.mvdirona.com

# Agenda



- High Scale Services
  - Infrastructure cost breakdown
  - Where does the power go?

- Power Distribution Efficiency

- Mechanical System Efficiency

- Server & Applications Efficiency
  - Hot I/O workloads & NAND flash
  - Resource consumption shaping
  - Work done per joule & per dollar

# Background & Biases

- 15 years in database engine development
  - Lead architect on IBM DB2
  - Architect on SQL Server
- Past 5 years in services
  - Led Exchange Hosted Services Team
  - Architect on the Windows Live Platform
  - Architect on Amazon Web Services
- Talk does not necessarily represent positions of current or past employers

# Services Different from Enterprises

- **Enterprise Approach:**
  - Largest cost is people -- scales roughly with servers (~100:1 common)
  - Enterprise interests center around consolidation & utilization
    - Consolidate workload onto fewer, larger systems
    - Large SANs for storage & large routers for networking

- **Internet-Scale Services Approach:**
  - Largest costs is server & storage H/W
    - Typically followed by cooling, power distribution, power
    - Networking varies from very low to dominant depending upon service
    - People costs under 10% & often under 5% (>1000+:1 server:admin)
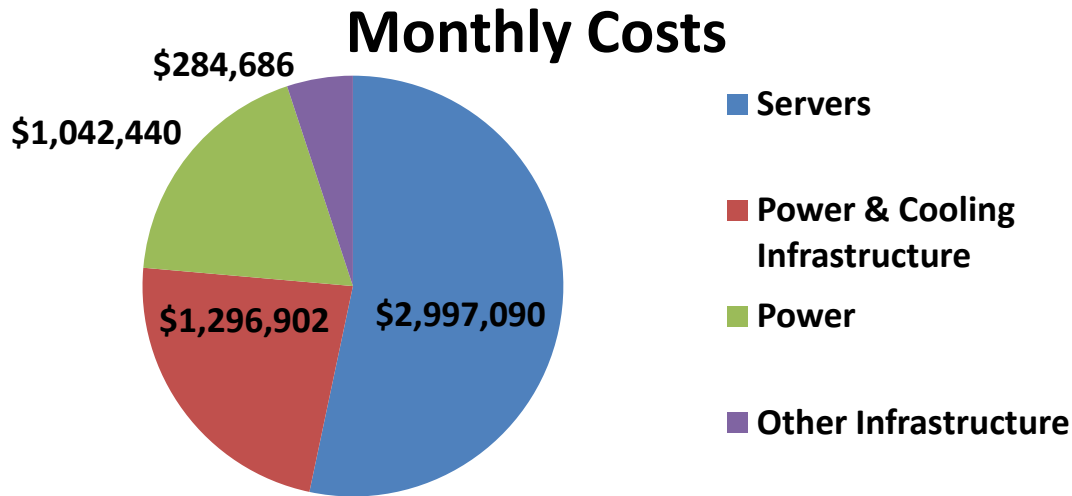  - Services interests center around work-done-per-$ (or joule)

- **Observations:**
  - People costs shift from top to nearly irrelevant.
  - Expect high-scale service techniques to spread to enterprise
  - Focus instead on work done/$ & work done/joule

# Power & Related Costs Dominate

- **Assumptions:**
  - Facility: ~$200M for 15MW facility (15-year amort.)
  - Servers: ~$2k/each, roughly 50,000 (3-year amort.)
  - Average server power draw at 30% utilization: 80%
  - Commercial Power: ~$0.07/kWhr

## Monthly Costs



$284,686

$1,042,440

$1,296,902

$2,997,090

- Servers
- Power & Cooling Infrastructure
- Power
- Other Infrastructure

3yr server & 15 yr infrastructure amortization
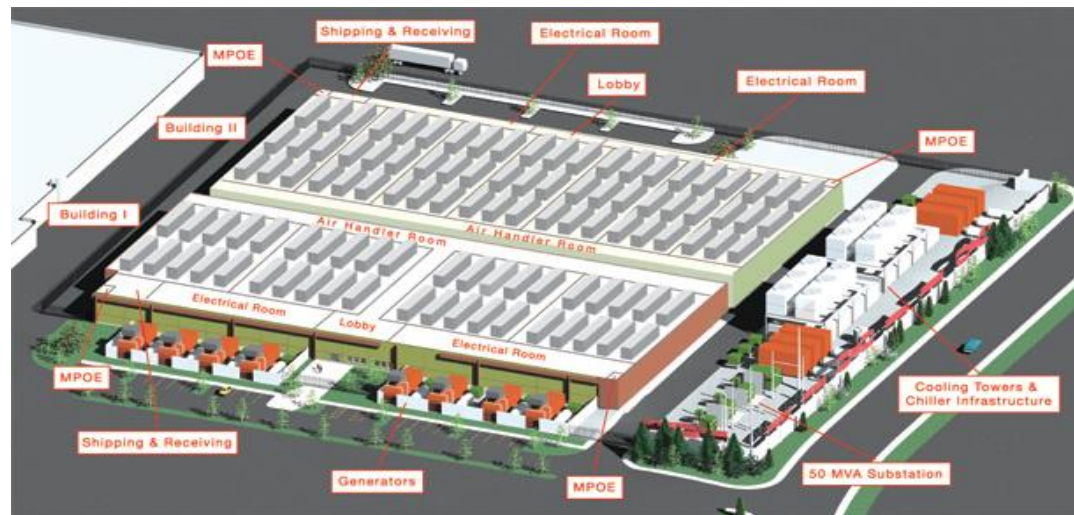


- **Observations:**
  - $2.3M/month from charges functionally related to power
  - Power related costs trending flat or up while server costs trending down

    Details at: http://perspectives.mvdirona.com/2008/11/28/CostOfPowerInLargeScaleDataCenters.aspx

# PUE & DCiE

- Measure of data center infrastructure efficiency
- Power Usage Effectiveness
  - PUE = (Total Facility Power)/(IT Equipment Power)
- Data Center Infrastructure Efficiency
  - DCiE = (IT Equipment Power)/(Total Facility Power) * 100%
- Help evangelize **tPUE** (power to server components)
  - http://perspectives.mvdirona.com/2009/06/15/PUEAndTotalPowerUsageEfficiencyTPUE.aspx



http://www.thegreengrid.org/en/Global/Content/white-papers/The-Green-Grid-Data-Center-Power-Efficiency-Metrics-PUE-and-DCiE

# Where Does the Power Go?

- **Assuming a pretty good data center with PUE ~1.7**
  - Each watt to server loses ~0.7W to power distribution losses & cooling
  - IT load (servers): 1/1.7=> 59%

- **Power losses are easier to track than cooling:**
  - Power transmission & switching losses: 8%
    - Detailed power distribution losses on next slide
  - Cooling losses remainder:100-(59+8) => 33%

- **Observations:**
  - **Server  efficiency & utilization improvements highly leveraged**
  - **Cooling costs unreasonably high**



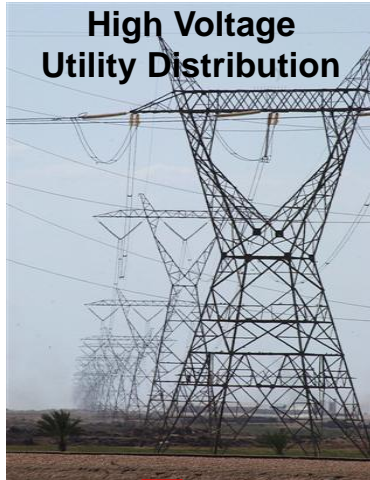http://perspectives.mvdirona.com

# Agenda



- High Scale Services
  - Infrastructure cost breakdown
  - Where does the power go?
- Power Distribution Efficiency
- Mechanical System Efficiency
- Server & Applications Efficiency
  - Hot I/O workloads & NAND flash
  - Resource consumption shaping
  - Work done per joule & per dollar

# Power Distribution

**High Voltage Utility Distribution**

8% distribution loss
.997^3*.94*.99 = 92.2%

**IT Load (servers, storage, Net, …)**

**2.5MW Generator (180 gal/hr)**

115kv

13.2kv

208V

~1% loss in switch gear & conductors

**Transformers**

**UPS: Rotary or Battery**

13.2kv

**Transformers**

13.2kv

480V

**Transformers**

0.3% loss
99.7% efficient

6% loss
94% efficient, ~97% available

0.3% loss
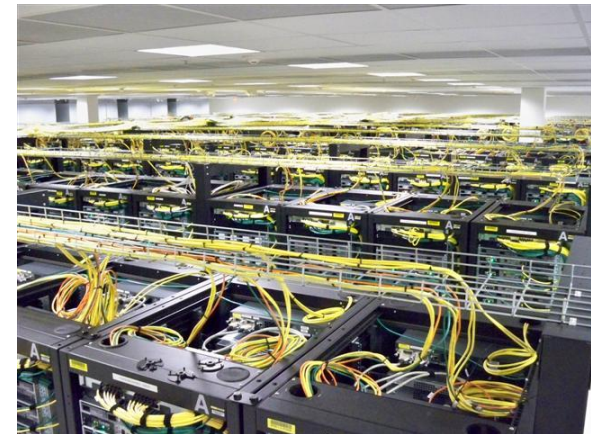99.7% efficient

0.3% loss
99.7% efficient

# Power Distribution Efficiency Summary



- Two additional conversions in server:
    1. Power Supply: often <80% at typical load
    2. On board step-down (VRM/VRD): <80% common
        - ~95% efficient both available & affordable

- Rules to minimize power distribution losses:
    1. Oversell power (more theoretic load that power)
    2. Avoid conversions (fewer transformer steps & efficient UPS)
    3. Increase efficiency of conversions
    4. High voltage as close to load as possible
    5. Size VRMs & VRDs to load & use efficient parts
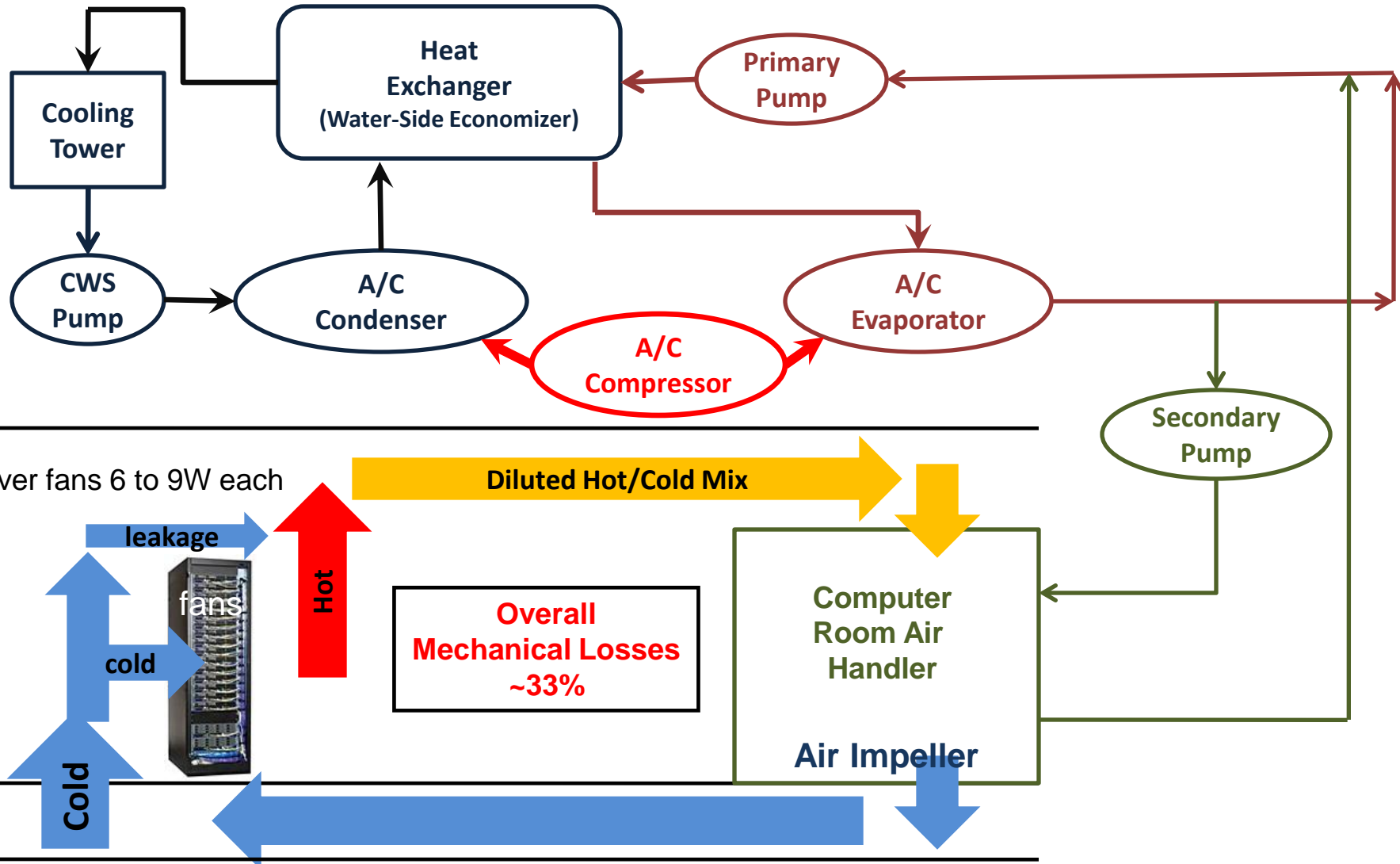    6. DC distribution potentially a small win (regulatory issues)

# Agenda

- High Scale Services
  - Infrastructure cost breakdown
  - Where does the power go?
- Power Distribution Efficiency
- Mechanical System Efficiency
- Server & Applications Efficiency
  - Hot I/O workloads & NAND flash
  - Resource consumption shaping
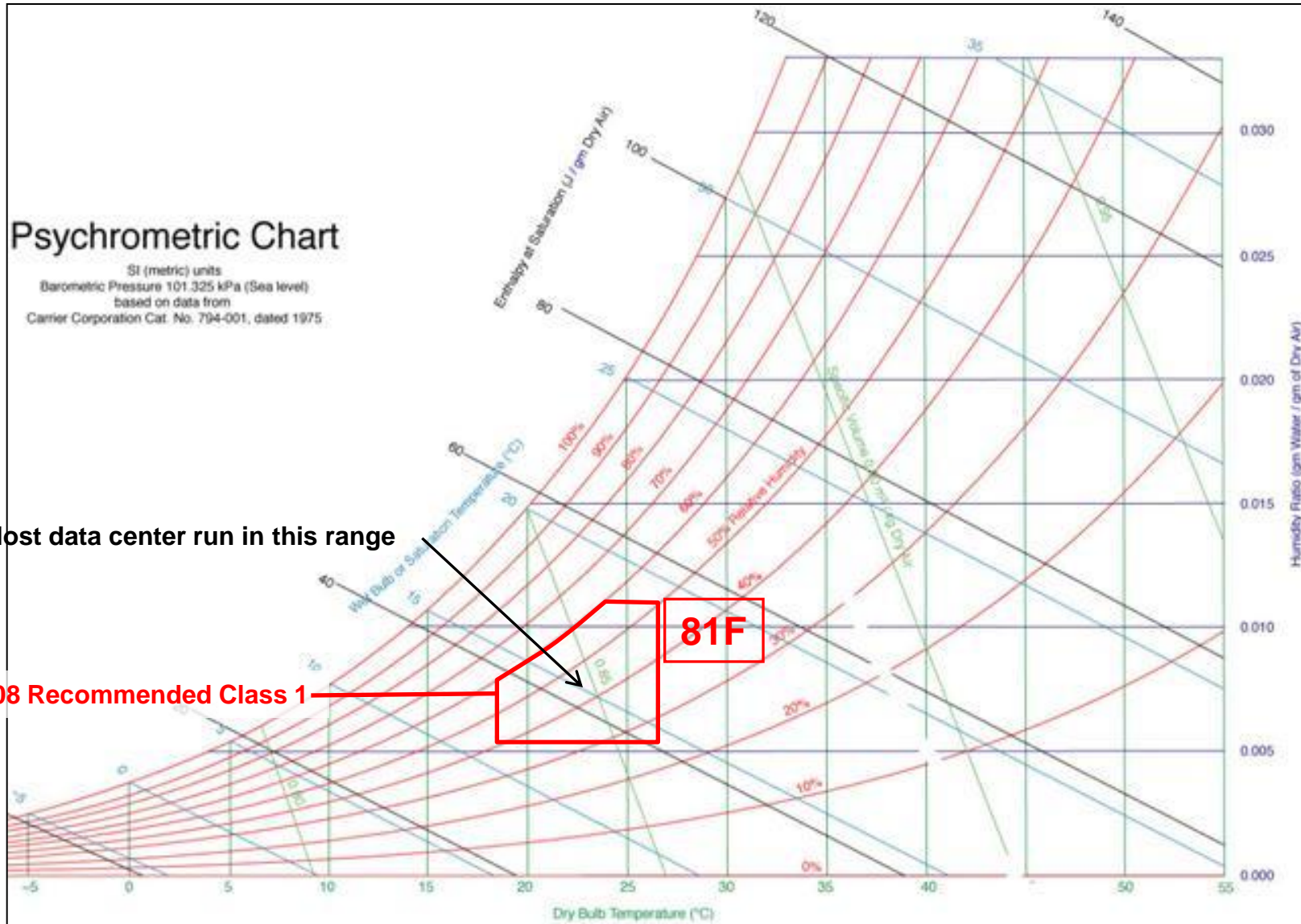  - Work done per joule & per dollar

# Conventional Mechanical Design

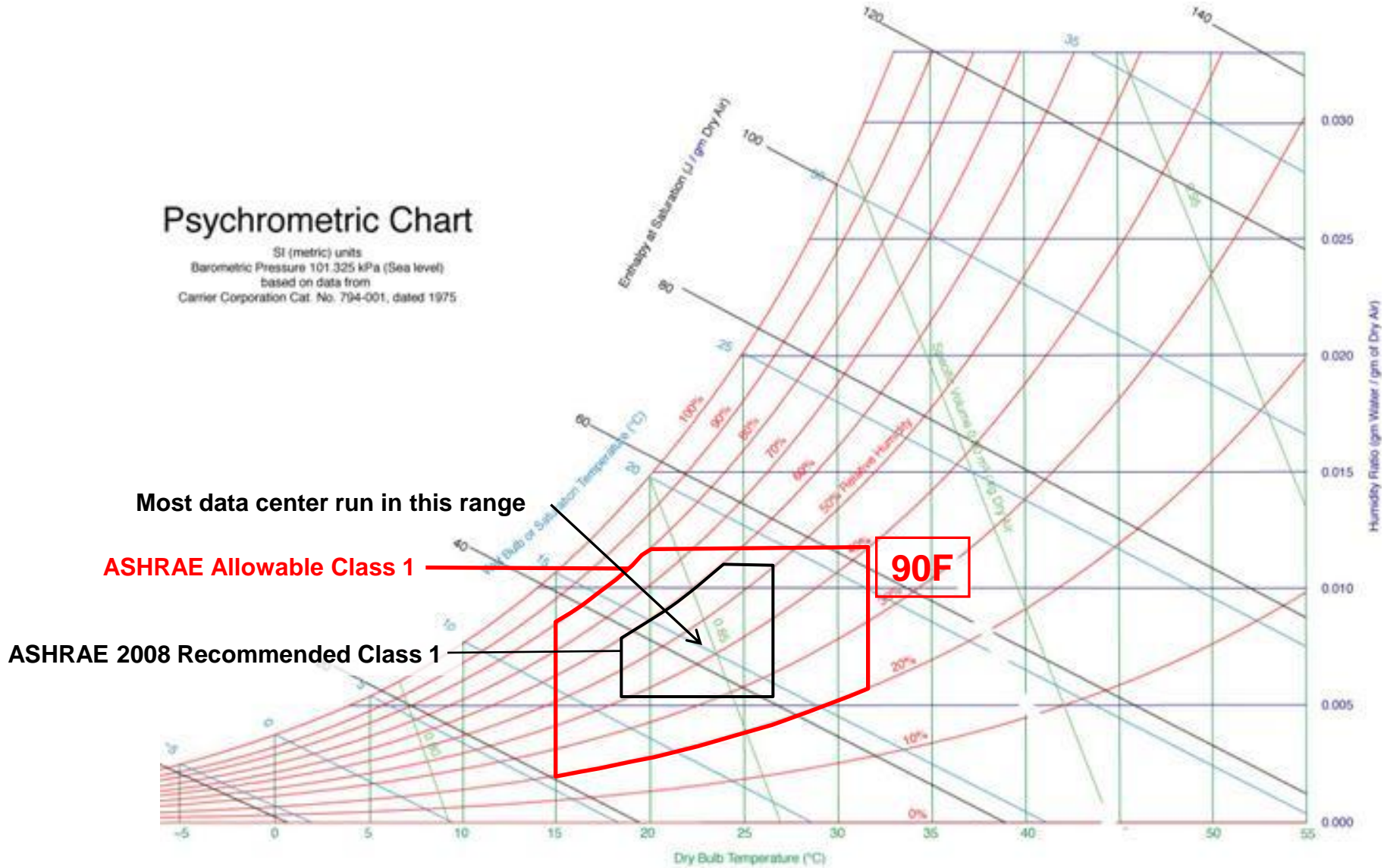Blow down & Evaporative Loss for 15MW facility: ~360,000 gal/day

**Cooling Tower**

**Heat Exchanger (Water-Side Economizer)**

**Primary Pump**

**CWS Pump**

**A/C Condenser**

**A/C Compressor**

**A/C Evaporator**

**Secondary Pump**

Server fans 6 to 9W each

**Diluted Hot/Cold Mix**

leakage

fans

**Hot**

**cold**

**Overall Mechanical Losses ~33%**

**Cold**

**Computer Room Air Handler**

**Air Impeller**

# ASHRAE 2008 Recommended



Psychrometric Chart

SI (metric) units
Barometric Pressure 101.325 kPa (Sea level)
based on data from
Carrier Corporation Cat. No. 794-001, dated 1975

**Most data center run in this range**

**81F**

**ASHRAE 2008 Recommended Class 1**
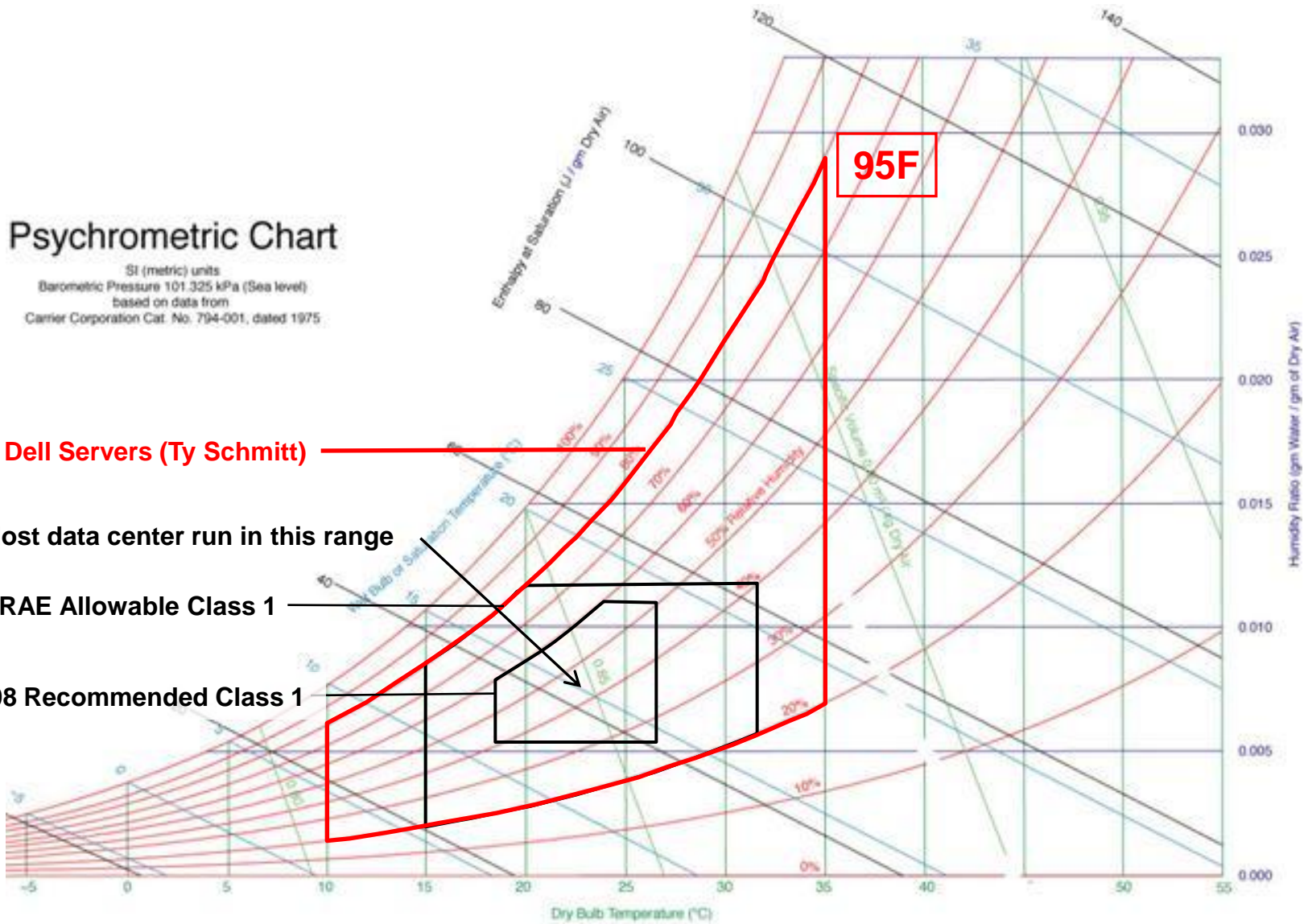
Dry Bulb Temperature (°C)

# ASHRAE Allowable



**Psychrometric Chart**

SI (metric) units
Barometric Pressure 101.325 kPa (Sea level)
based on data from
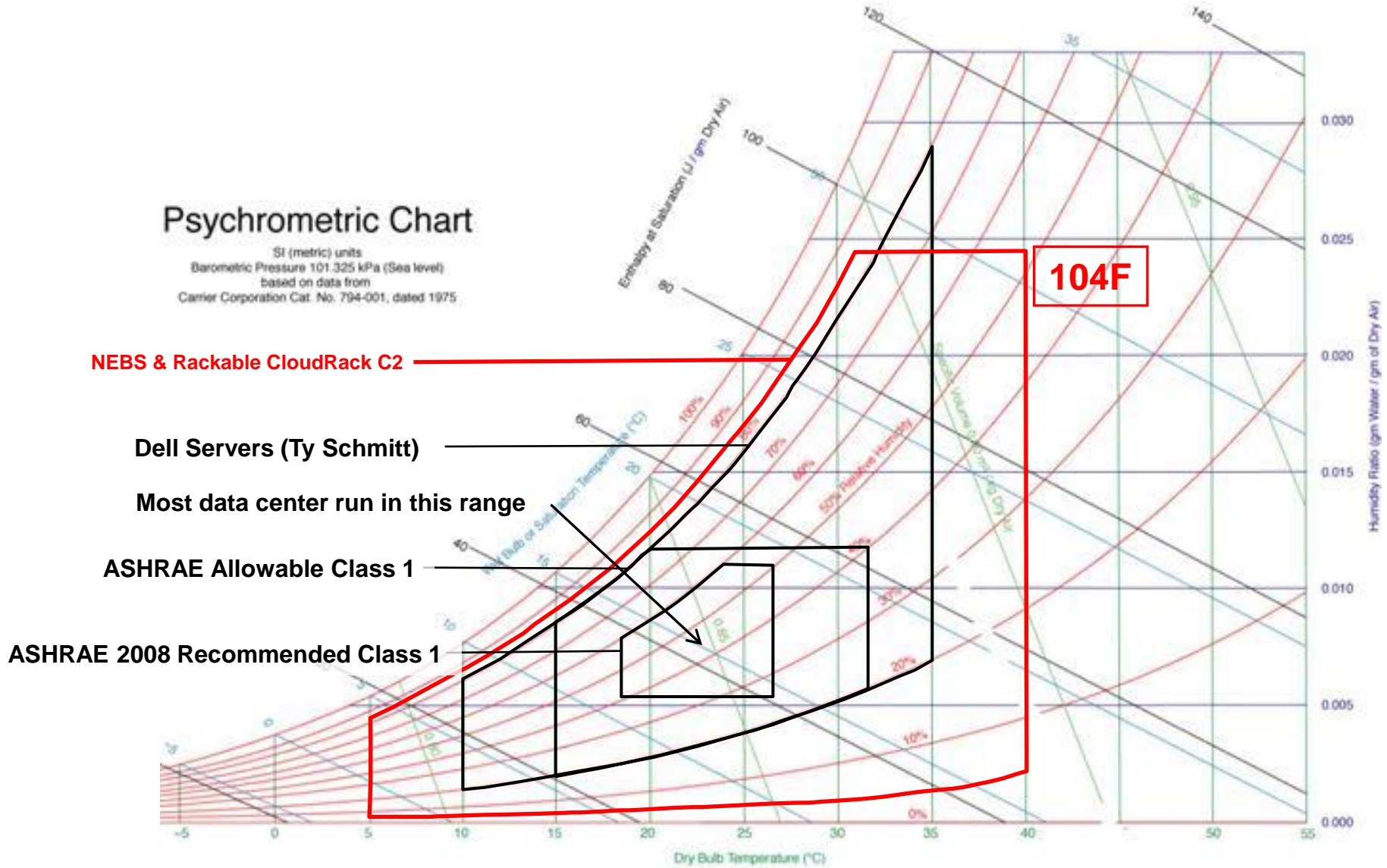Carrier Corporation Cat. No. 794-001, dated 1975

**Most data center run in this range**

**ASHRAE Allowable Class 1**

**90F**

**ASHRAE 2008 Recommended Class 1**

Dry Bulb Temperature (°C)

# Dell PowerEdge 2950 Warranty



Psychrometric Chart
SI (metric) units
Barometric Pressure 101.325 kPa (Sea level)
based on data from
Carrier Corporation Cat. No. 794-001, dated 1975

**95F**

**Dell Servers (Ty Schmitt)**

**Most data center run in this range**

**ASHRAE Allowable Class 1**

**ASHRAE 2008 Recommended Class 1**

# NEBS (Telco) & Rackable Systems



Psychrometric Chart
SI (metric) units
Barometric Pressure 101.325 kPa (Sea level)
based on data from
Carrier Corporation Cat. No. 794-001, dated 1975

**104F**

**NEBS & Rackable CloudRack C2**

**Dell Servers (Ty Schmitt)**

**Most data center run in this range**

**ASHRAE Allowable Class 1**

**ASHRAE 2008 Recommended Class 1**

# Air Cooling

- Allowable component temperatures higher than hottest place on earth
  - Al Aziziyah, Libya: 136F/58C (1922)
- It's only a mechanical engineering problem
  - More air & better mechanical designs
  - Tradeoff: power to move air vs cooling savings & semi-conductor leakage current
  - Partial recirculation when external air too cold
- Currently available equipment:
  - 40C: Rackable CloudRack C2
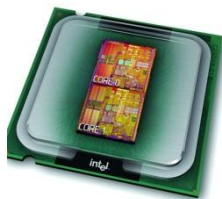  - 35C: Dell Servers

**Memory: 3W - 20W**
**Temp Spec: 85C-105C**

**Hard Drives: 7W- 25W**
**Temp Spec: 50C-60C**

**Rackable CloudRack C2**
**Temp Spec: 40C**

Thanks for data & discussions:
Ty Schmitt, Dell Principle Thermal/Mechanical Arch.
& Giovanni Coglitore, Rackable Systems CTO

**I/O: 5W - 25W**
**Temp Spec: 50C-60C**

**Processors/Chipset: 40W - 200W**
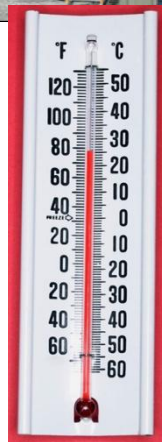**Temp Spec: 60C-70C**

# Air-Side Economization & Evaporative Cooling





- Avoid direct expansion cooling entirely
- Ingredients for success:
  - Higher data center temperatures
  - Air side economization
  - Direct evaporative cooling
- Particulate concerns:
  - Usage of outside air during wildfires or datacenter generator operation
  - Solution: filtration & filter admin or heat wheel & related techniques
- Others: higher fan power consumption, more leakage current, higher failure rate

# Mechanical Efficiency Summary

- Mechanical System Optimizations:
    1. Tight airflow control, short paths & large impellers
    2. Raise data center temperatures
    3. Cooling towers rather than A/C
    4. Air side economization & evaporative cooling
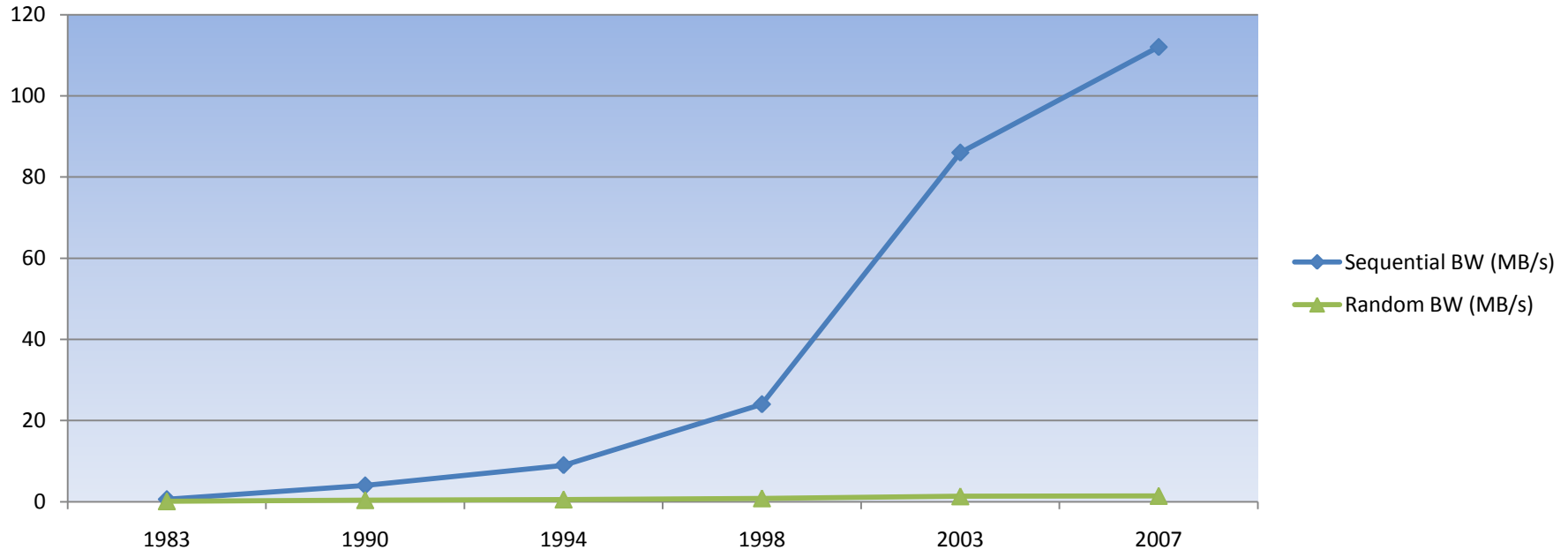        - outside air rather than A/C & towers

# Agenda

- **High Scale Services**
  - Infrastructure cost breakdown
  - Where does the power go?
- **Power Distribution Efficiency**
- **Mechanical System Efficiency**
- **Server & Applications Efficiency**
  - Hot I/O workloads & NAND flash
  - Resource consumption shaping
  - Work done per joule & per dollar

# Disk Random BW vs Sequential BW



Source: Dave Patterson with James Hamilton updates

- Disk sequential BW lagging DRAM and CPU
- Disk random access BW growth ~10% of sequential
- **Conclusion**: Storage Chasm widening requiring larger memories & more disks

# Memory to Disk Chasm

- Disk I/O rates grow slowly while CPU data consumption grows near Moore pace
  - Random read 1TB disk: 15 to 150 days*
- Sequentialize workloads
  - Essentially the storage version of cache conscious algorithms
    - e.g. map/reduce
  - Disks arrays can produce acceptable aggregate sequential bandwidth
- Redundant data: materialized views & indexes
  - Asynchronous maintenance
  - Delta or stacked indexes (from IR world)
- Distributed memory cache (remote memory "closer" than disk)
- I/O Cooling: Blend hot & cold data (using HDD)
- I/O concentration: partition hot & cold (SSD & HDD mix)

*** Tape is Dead, Disk is Tape, Flash is Disk, Ram Locality is King (Jim Gray)**

# Case Study: TPC-C with SSD

| Slot | Controller | Disks | | | Capacity | | Usage | |
|---|---|---|---|---|---|---|---|---|
| 0 | Dell PERC5i | 8x73GB,15K,SAS | | RAID10 | Disk 6 | 15GB | OS | **O/S & Log** |
| | | | | | 279.99GB | 260GB | Logs | |
| 3 | Dell PERC6/E | 15x36GB,15K,SAS | | RAID0 | Disk 2 | 488.92GB | DB data | |
| | | 15x36GB,15K,SAS | | RAID0 | Disk 3 | 488.92GB | DB data | |
| 4 | Dell PERC6/E | 15x36GB,15K,SAS | | RAID0 | Disk 4 | 488.92GB | DB data | **Data** |
| | | 15x36GB,15K,SAS | | RAID0 | Disk 5 | 488.92GB | DB data | |
| 6 | Dell PERC6/E | 15x73GB,15K,SAS | | RAID0 | Disk 0 | 1016.23GB | DB data | |
| | | 15x73GB,15K,SAS | | RAID0 | Disk 1 | 1016.23GB | DB data | |



- 98 HDD total
  - 90 data disks (primarily random access)
  - 8 log & O/S disks (primarily sequential access)
- Compute HDD/SSD cross-over using fictitious SSD
  - 128GB SSD @ 7k IOPS
- 90 HDD to store 2,464GB (short stroked)
  - 106GB static & 2,357GB dynamic (60 day rule)
  - 90 disk HDD budget: $26,910 (disks $299 each)
  - ***Requires 20 SSDs to support @ up to $1,346 each***
- Static content only (drop 60 day rule)
  - Conservatively estimate 45k IOPS
    - Used 90 short stroked disks at 500 IOPS each
  - ***Requires 7 SSDs at up to $3,844 (easy)***
- **Very hot I/O workloads a win on SSD**

http://www.tpc.org/results/FDR/TPCC/Dell_2900_061608_fdr.pdf

# Summary

- CPU optimizations are always welcome but the biggest design & optimization problems today are at the datacenter level

- In work at all levels, focus on:
  - Work done per dollar
  - Work done per joule

- Single dimensional performance measurements are not interesting at scale unless balanced against cost

# More Information



- **This Slide Deck:**
  - I will post these slides to http://mvdirona.com/jrh/work later this week
- **Power and Total Power Usage Effectiveness (tPUE)**
  - http://perspectives.mvdirona.com/2009/06/15/PUEAndTotalPowerUsageEfficiencyTPUE.aspx
- **Berkeley Above the Clouds**
  - http://perspectives.mvdirona.com/2009/02/13/BerkeleyAboveTheClouds.aspx
- **Degraded Operations Mode**
  - http://perspectives.mvdirona.com/2008/08/31/DegradedOperationsMode.aspx
- **Cost of Power**
  - http://perspectives.mvdirona.com/2008/11/28/CostOfPowerInLargeScaleDataCenters.aspx
  - http://perspectives.mvdirona.com/2008/12/06/AnnualFullyBurdenedCostOfPower.aspx
- **Power Optimization:**
  - http://labs.google.com/papers/power_provisioning.pdf
- **Cooperative, Expendable, Microslice Servers**
  - http://perspectives.mvdirona.com/2009/01/15/TheCaseForLowCostLowPowerServers.aspx
- **Power Proportionality**
  - http://www.barroso.org/publications/ieee_computer07.pdf
- **Resource Consumption Shaping:**
  - http://perspectives.mvdirona.com/2008/12/17/ResourceConsumptionShaping.aspx
- **Email**
  - James@amazon.com