



Duke Systems

Cool Clouds ICAC 2009

Jeff Chase
Duke University

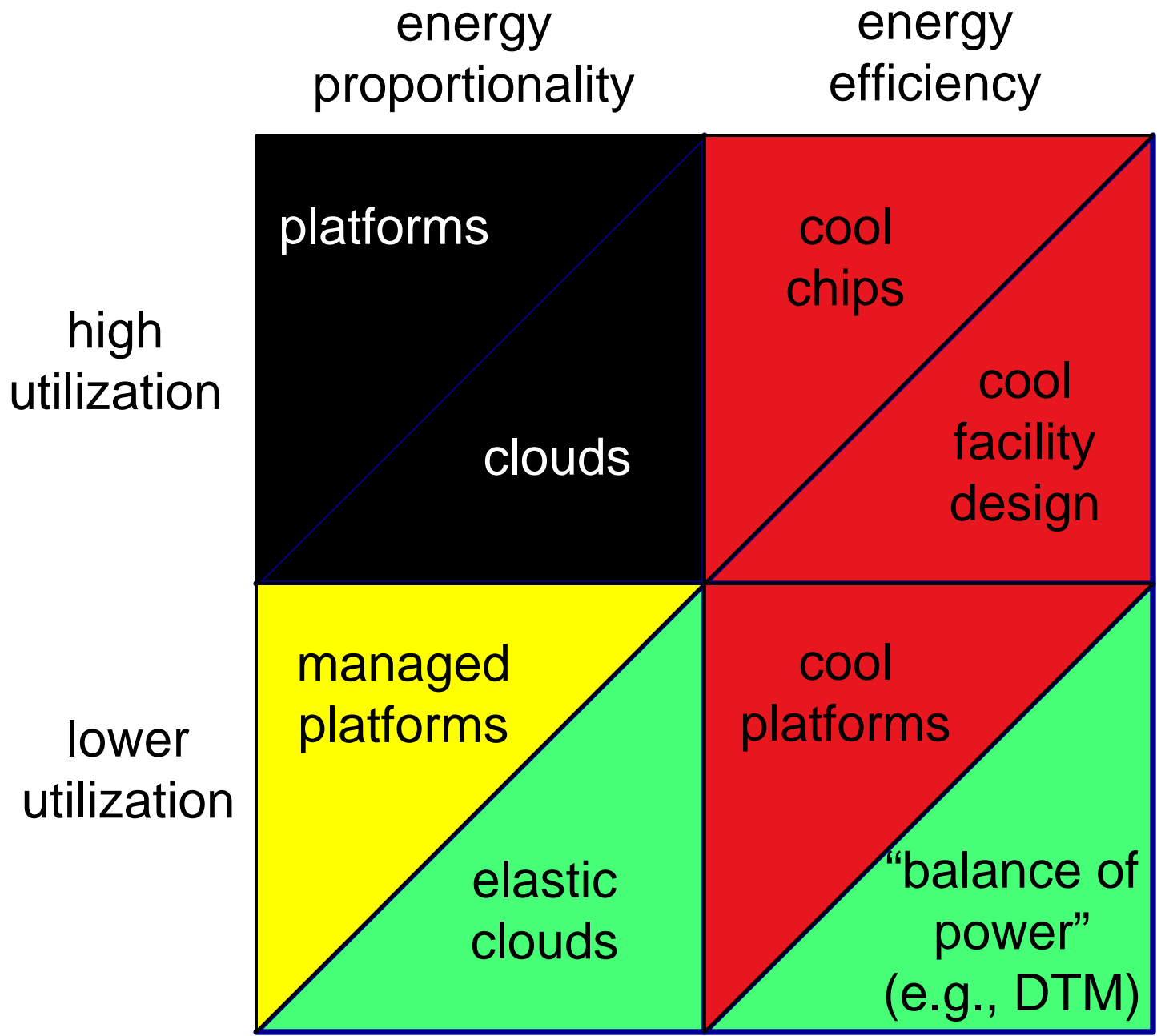


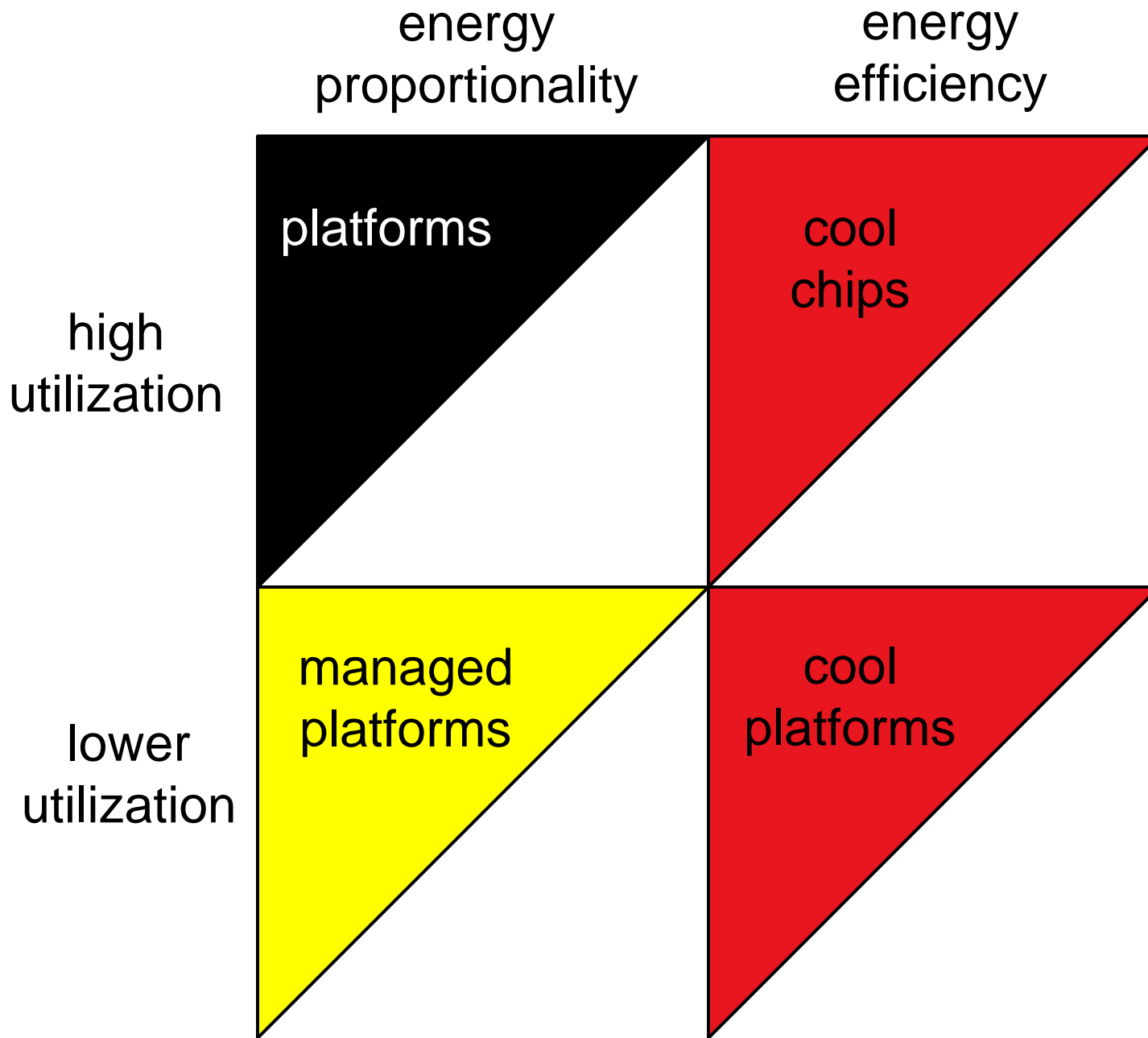
Basics

- **What is the energy cost for a DC/WHC?**
 - “20% and rising”
 - Electricity: \$70/MWH or \$700?
 - Carbon tax/offsets: 30% or 300%?
- **How much can we save with autonomies?**
 - “It depends”
 - On the platform
 - On the workload
 - “0% to 30%”

Where is the sweet spot?

- **Energy proportionality vs. efficiency**
- **High utilization vs. lower utilization**
- **Smart platforms vs. smart clouds**
- **Load shifting: flatten the demand curve**
 - **Demand-side management**
 - **Reflective applications**





energy
proportionality

energy
efficiency

high
utilization

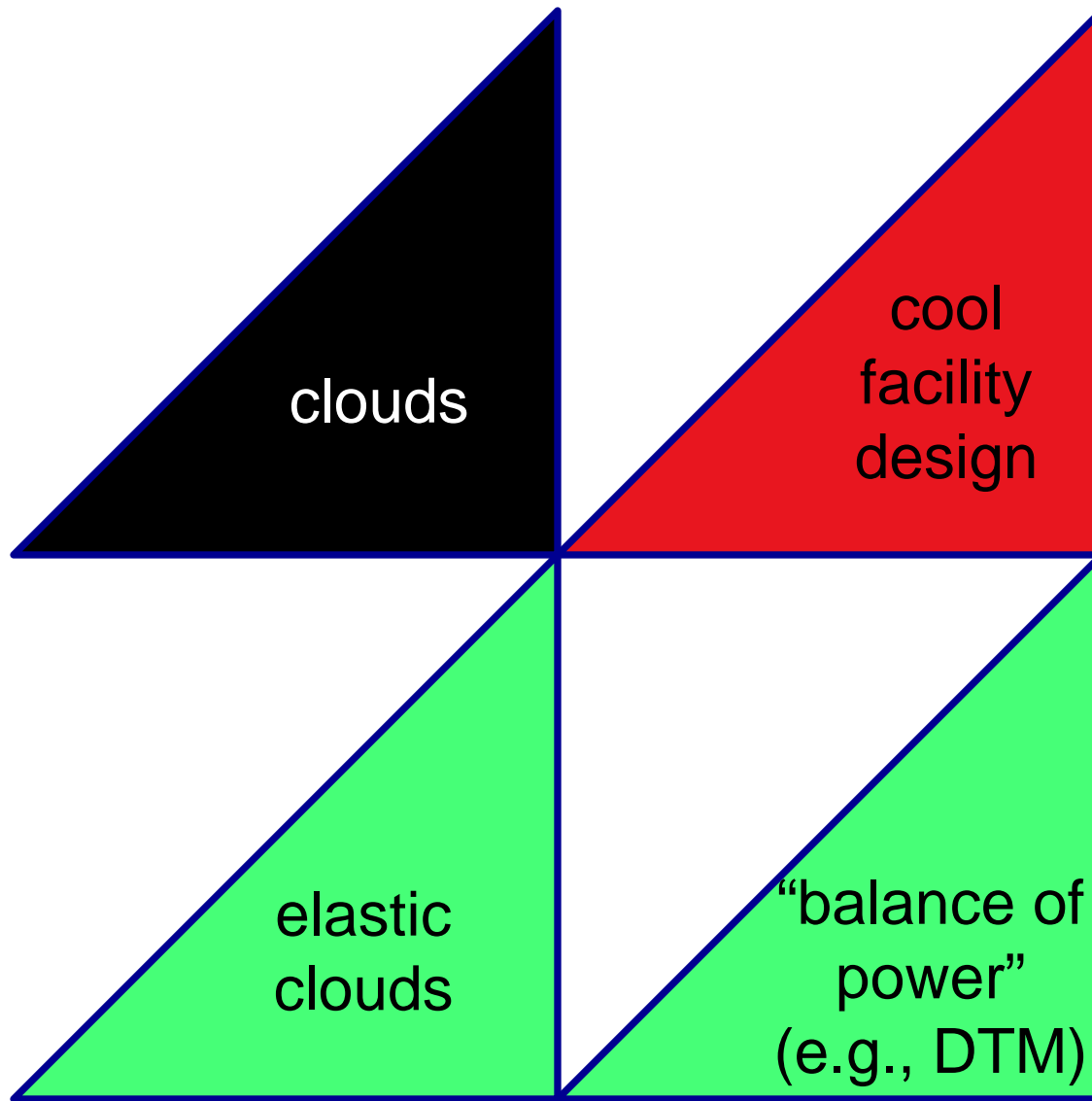
clouds

cool
facility
design

lower
utilization

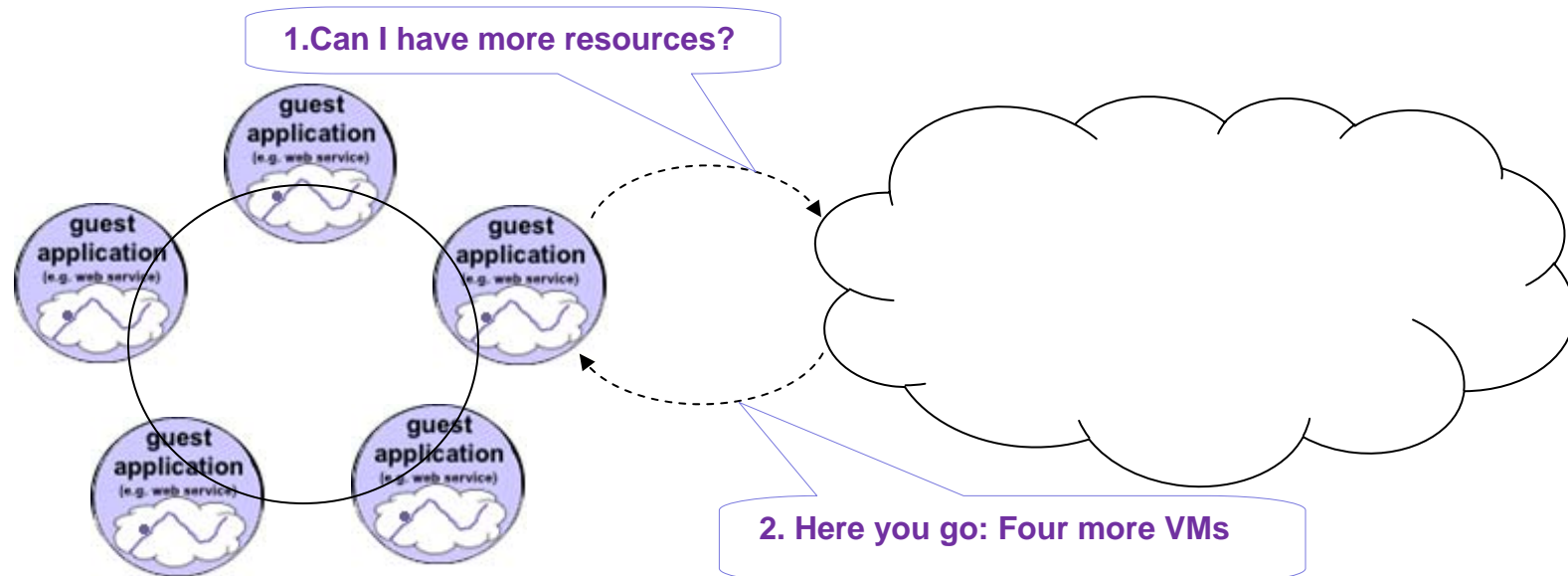
elastic
clouds

“balance of
power”
(e.g., DTM)

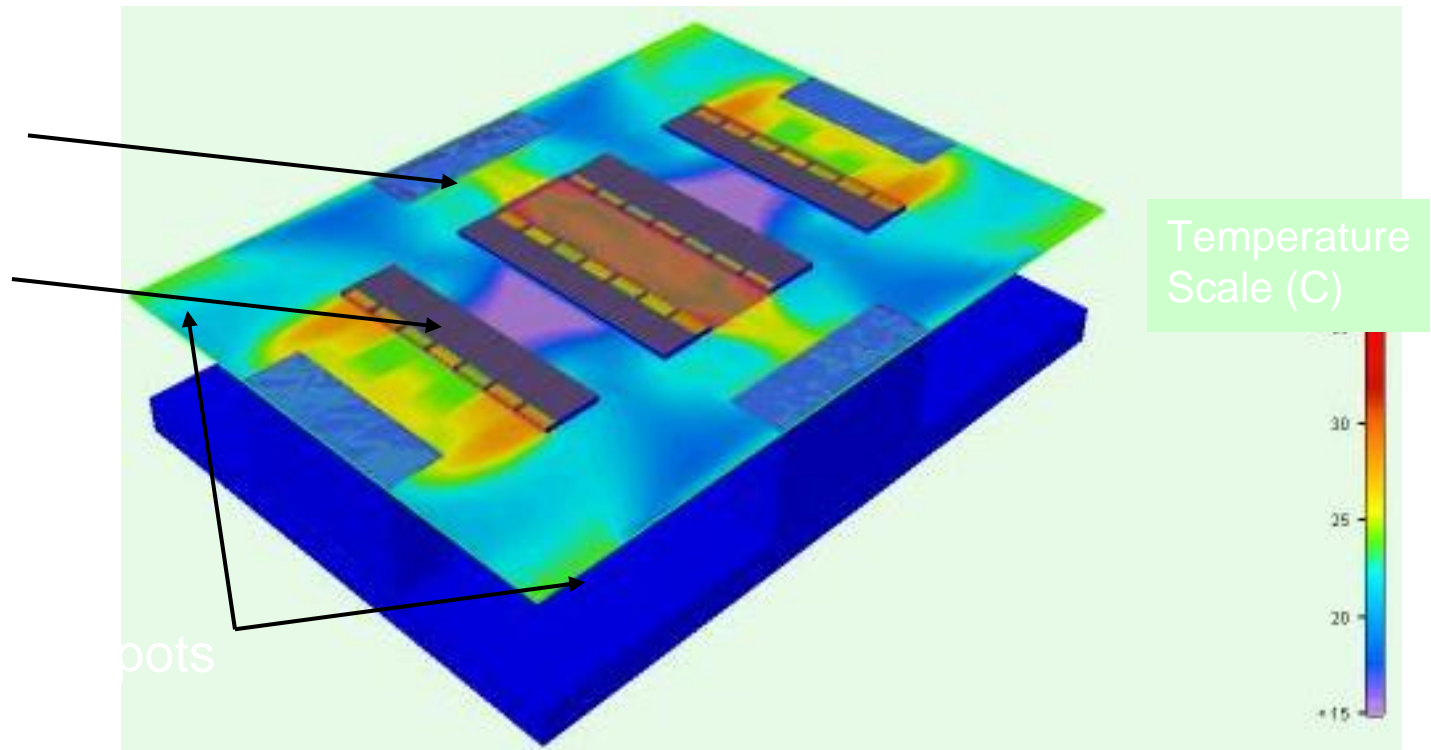


The Elasticity of Power

- ▶ Clouds: “boundless infrastructure on demand”
- ▶ Elasticity: Grow/shrink resource slices as required

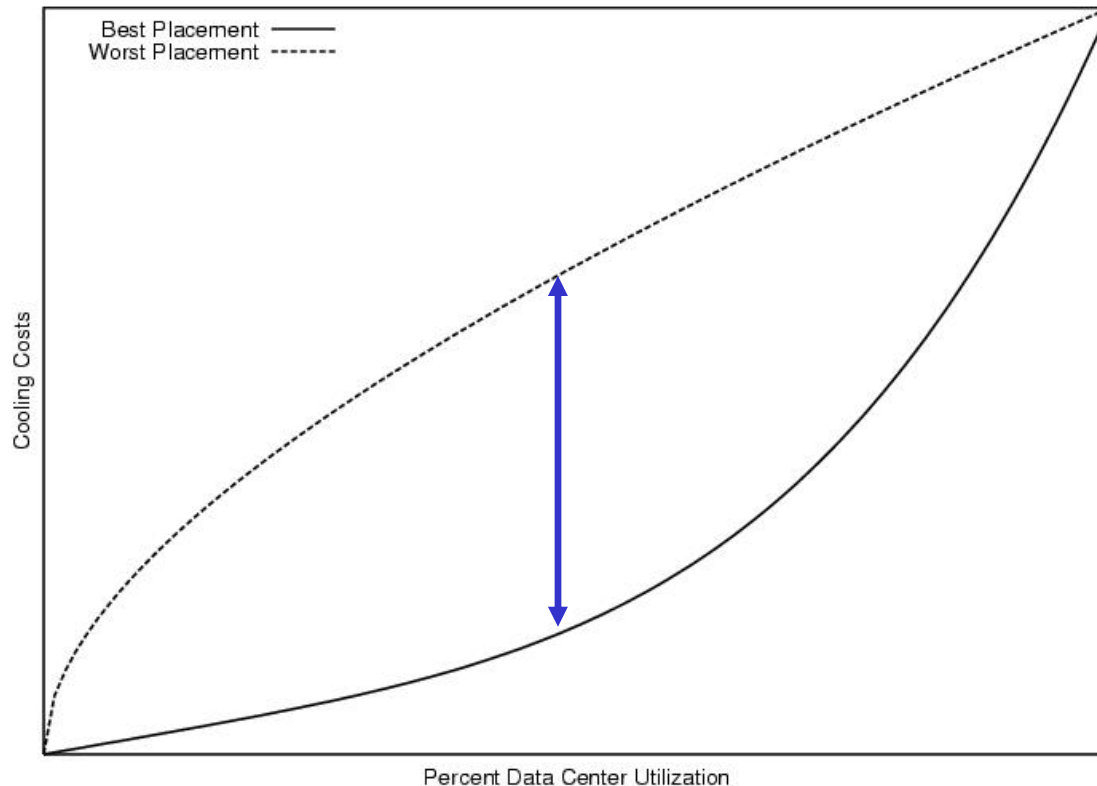


“Balance of Power”



- Continuous thermal sensors
- Infer “thermal topology”
- Place workload to optimize cooling
- Dynamic thermal management for DC/WHC

The Importance of Being Idle

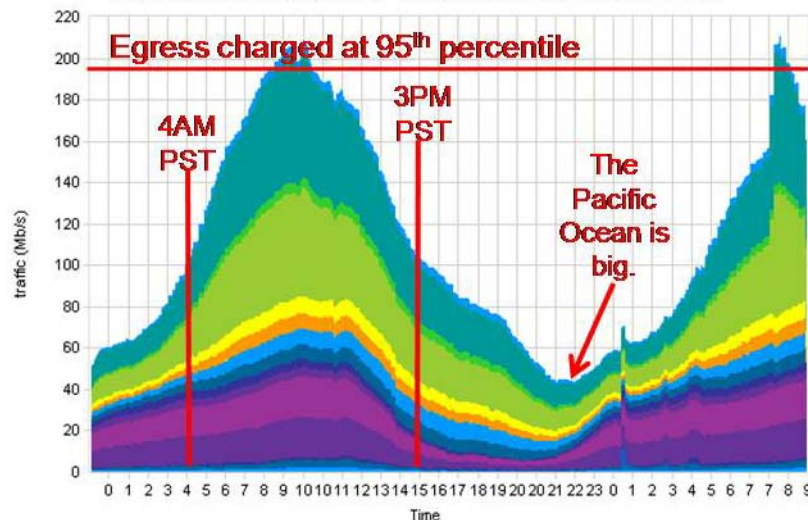


— 0%: No choices; 100%: No choices

*Only midrange has a useful spread
between good choices and bad choices.*

Flatten the Demand Curve?

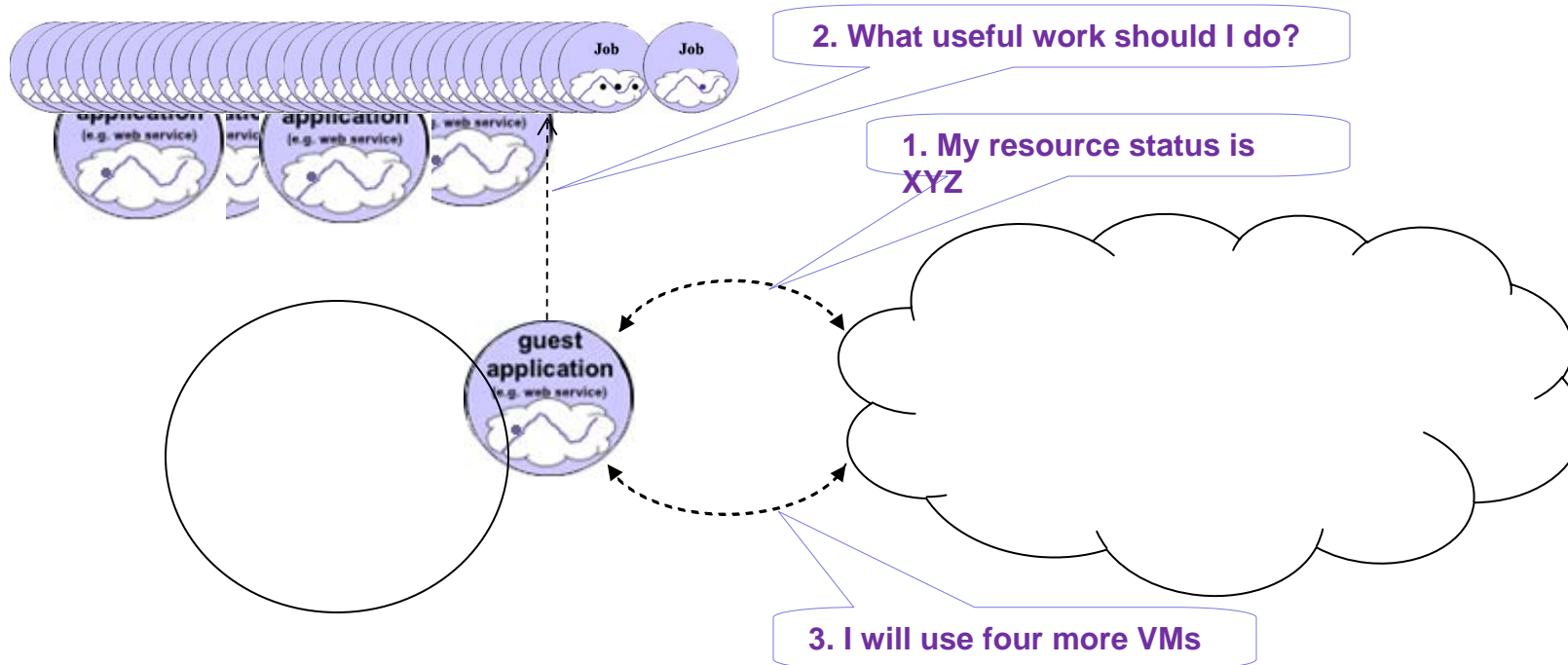
- ▶ Statistical multiplexing is not enough
 - ▶ Not for smart clouds or smart (electrical) grids
 - ▶ Wide variance in aggregate demand
 - ▶ Congestion → higher price, higher carbon footprint
- ▶ Smoother ride? **Demand side management**



James Hamilton

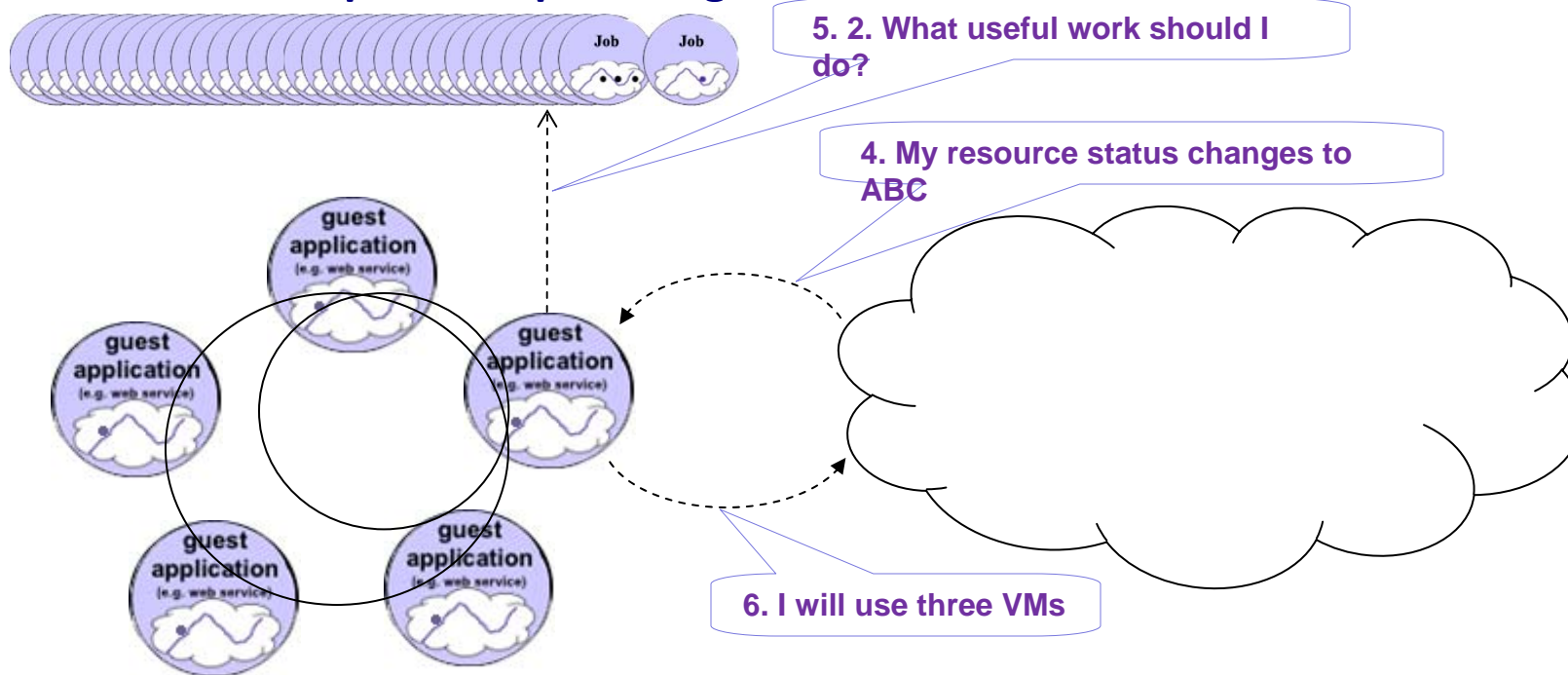
Demand Side Management

- ▶ **Reflection** in elastic cloud applications:
 - ▶ Adapt behavior based on resource availability
 - ▶ Opportunistically exploit surplus resources
 - ▶ Defer/avoid work during congestion



Reflective elastic applications

- ▶ **Reflection** in elastic cloud applications:
 - ▶ Adapt behavior based on resource availability
 - ▶ Opportunistically exploit surplus resources
 - ▶ Defer/avoid work during congestion
 - ▶ Require deeper integrated control



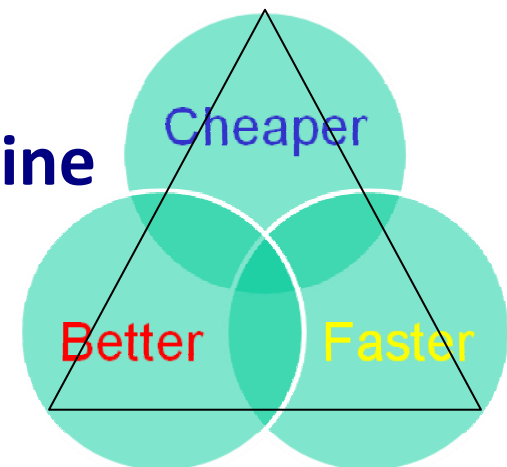
Reflective control objectives

Given:

- ▶ Elastic resource utilization
- ▶ Variable/dynamic cost, utility
- ▶ (Variable/dynamic price)

How to:

- ▶ Balance budget, accuracy, and deadline
- ▶ Perform cost-benefit analysis



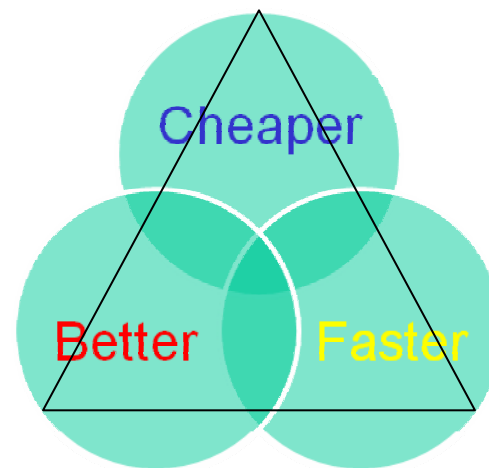
DSM/Reflection: Challenges

1. How to specify objectives?

- ▶ Multiple objectives: deadline, budget, accuracy
- ▶ How much parallelism for opportunistic/speculative use?
- ▶ Policies to adapt to variable/dynamic price?

2. Does it generalize? To what extent can we “factor out” reflective policies from applications?

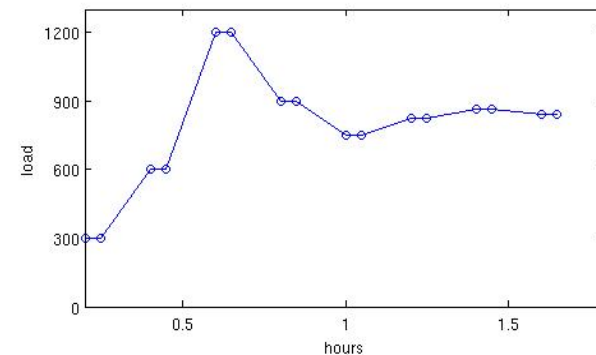
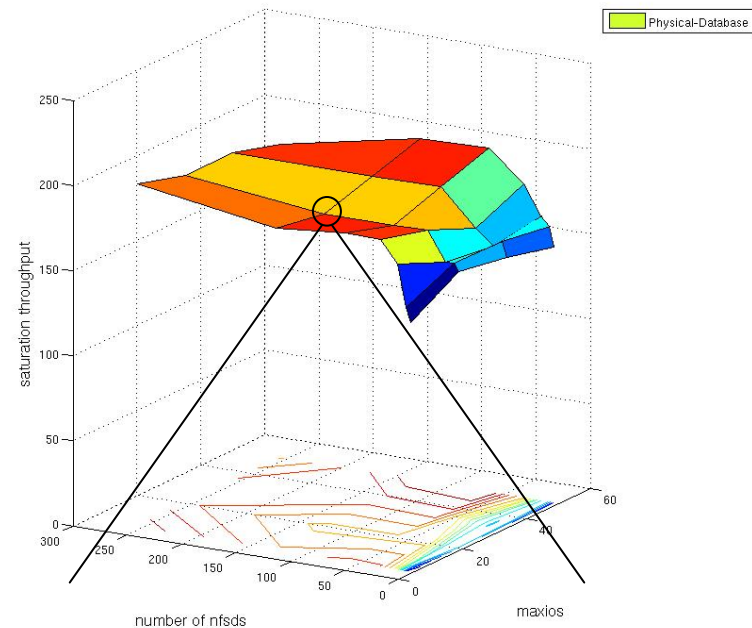
3. *What does it require from the “cool cloud”?*



(Backup slides)

Workbench-assisted benchmarking

- **Goal: response surface map:**
peak rate= $F(W, R, C)$
- **Parallelism**
 - Data dependency at each point in surface
 - Partition surface arbitrarily: embarrassingly data parallel
- **What experiments to run?**
 - Need notion of experiment utility: $u(e)$
 - Highly selective sampling

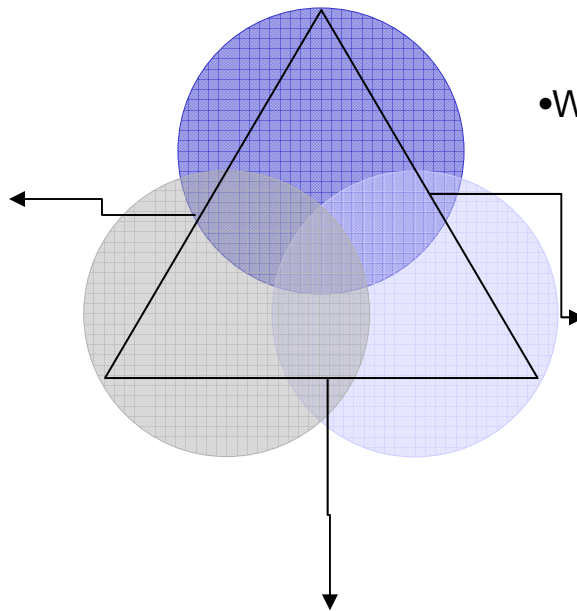


Better? Cheaper? Faster?

- How long should I run each experiment?
- Which search techniques can I use?
- How to quantify cost?

- Do I need more samples?
- Is confidence level satisfied?
- How many times should I run experiment?

- How to sample accurately?
- Can I trust result?



- What if resources cost less at 4am?
- How to handle dynamic resources?
- Do I have enough resources? Should I get more/ return some?
- Wait for more resources vs. running lower rank experiments?

- What's the optimal parallelism?

- What if I have deadline?