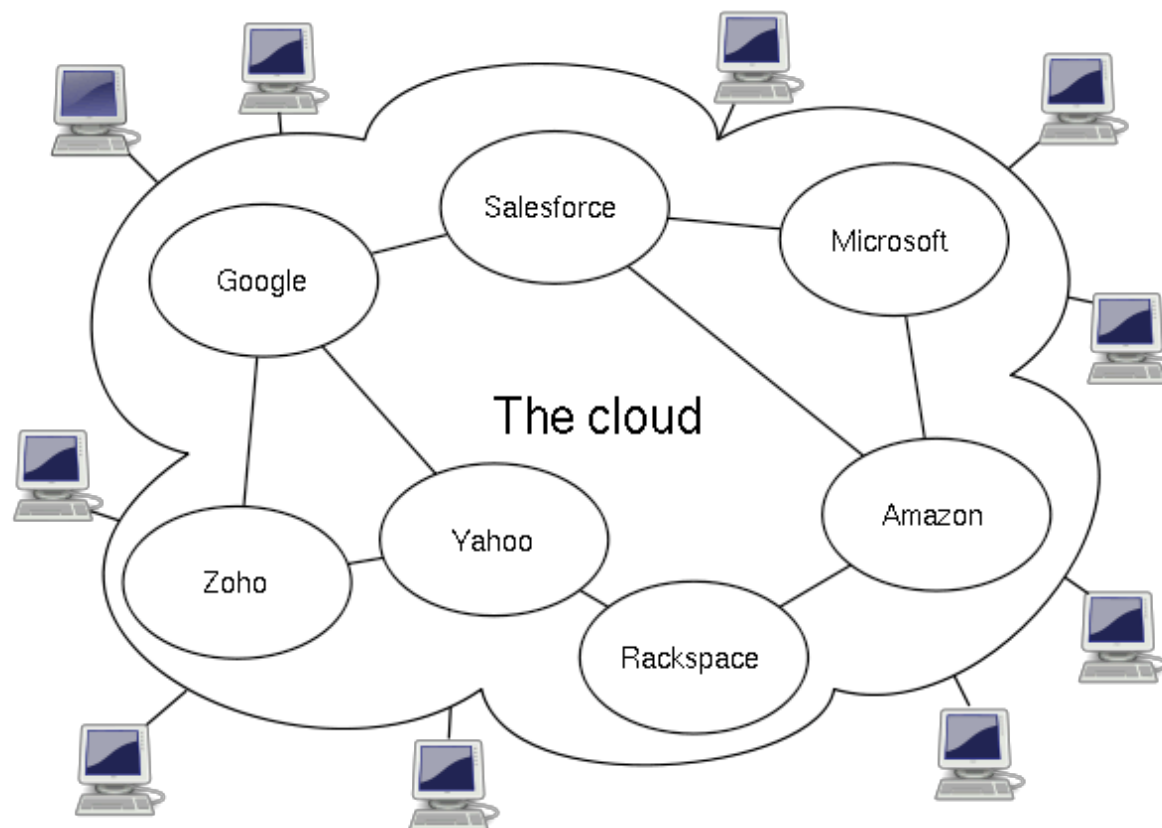# Cloud Computing Overview

CSC 501

# Cloud Computing

- Paradigm shift in computing model
  - Provision of dynamically *scalable* and *virtualized* resources as service over Internet
  - Computing and storage details abstracted from users

# History

- 1960s: concept of utility computing
  - Computing resources offered as public utility
- Late 1990s: SaaS
  - Software as a Service
  - Application licensed for use as service on demand
  - Extended by Microsoft through web services
- 2005: Amazon
  - Modernized datacenters
  - Started AWS (Amazon Web Services)
    - EC2: Amazon Elastic Compute Cloud
- 2007: Google, IBM, and universities: cloud computing initiative
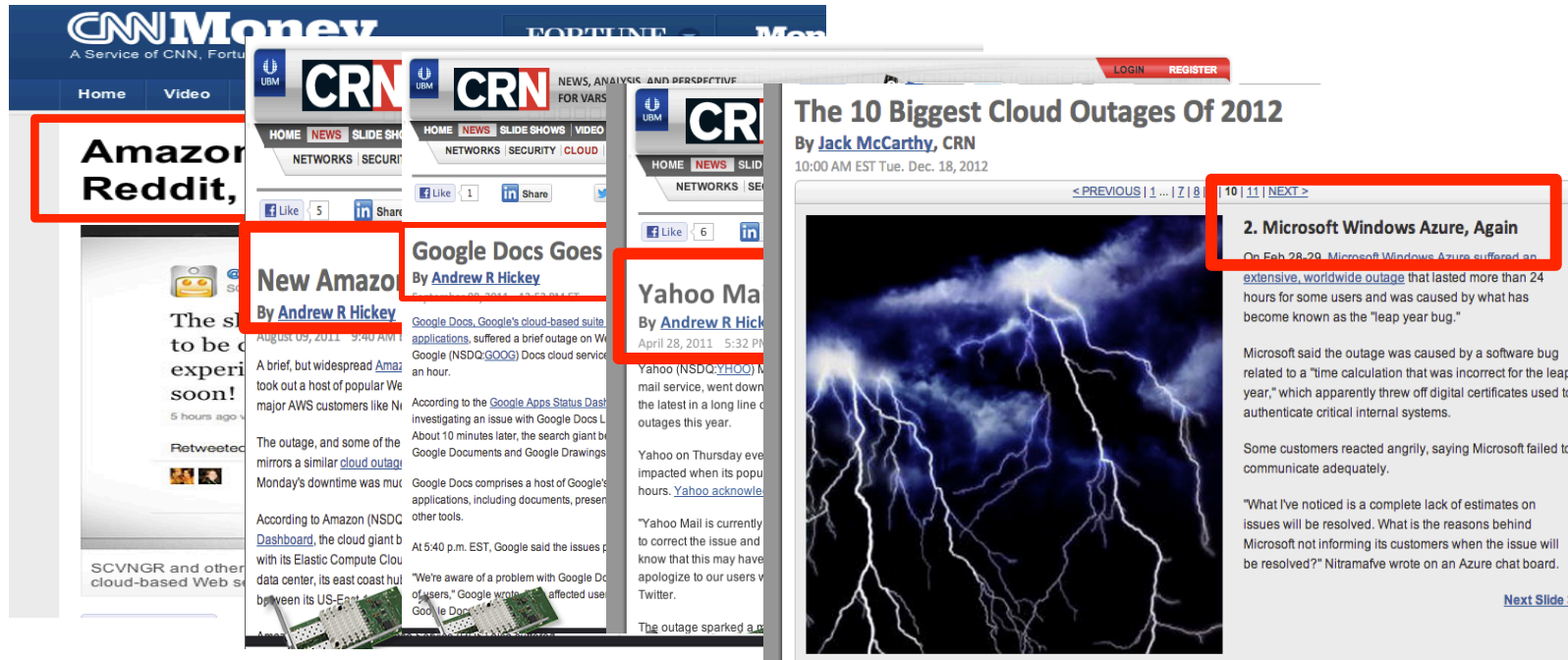
# What is Cloud Computing?

- **Computing resource leasing infrastructure**
  - Rent computing resources (hardware, infrastructure, virtual machines, software platforms, etc.) for a period of time rather than own resources permanently
  - Pay-as-you-go service model
  - Illusion of infinite resources
  - Resources on-demand
  - No up-front cost
  - Fine-grained billing (e.g. hourly)

# Why Cloud Computing?

- Pay-as-you-go resource leasing
  - Datacenter ownership is expensive
  - Avoid infrastructure maintenance hassle

- Elastic scaling
  - Instant resource on-demand
  - Illusion of infinite resources
  - No up-front cost

- Ubiquitous accessibility
  - Data aggregation

# Cloud Computing Challenges

- **Robustness**



- **Resource and energy efficiency**
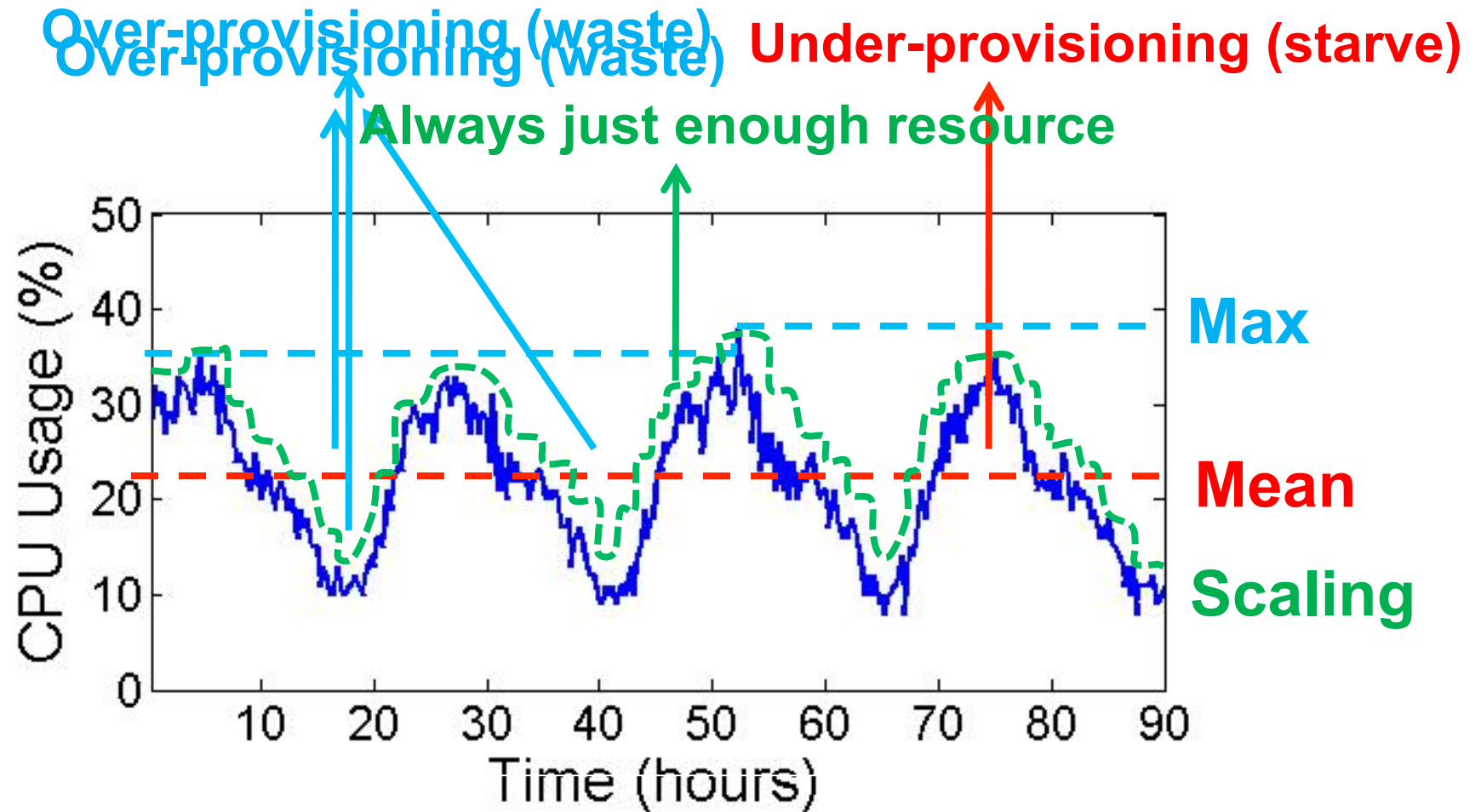  - Over-provision resources for performance

- **Accountability**
  - How do I trust the results from the cloud?

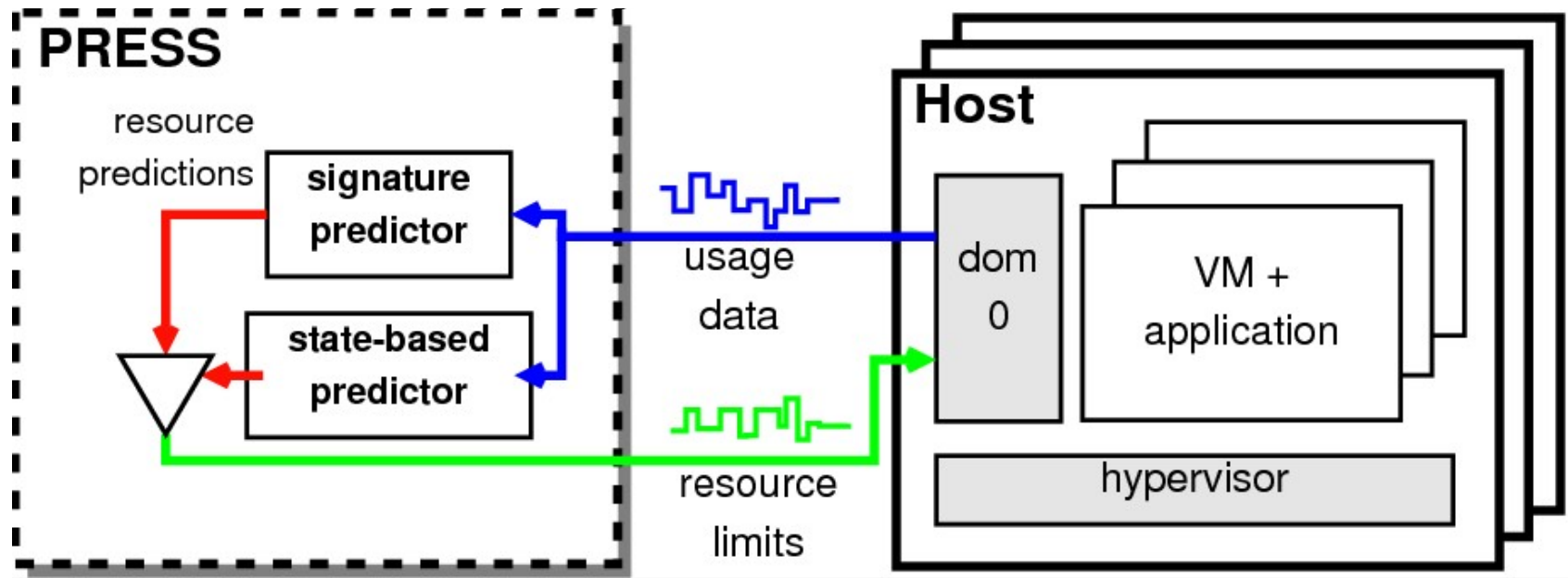# Efficient Cloud Resource Management

- **Elastic resource scaling**
  - Reduce resource waste
  - Minimize SLO violations
  - Reduce energy consumption by slowing down CPU frequencies

- **Server consolidation**
  - Reduce the number of physical hosts
  - Reduce energy consumption by shutting down machines

# Motivation

- RUBiS Web server driven by real workload trace:
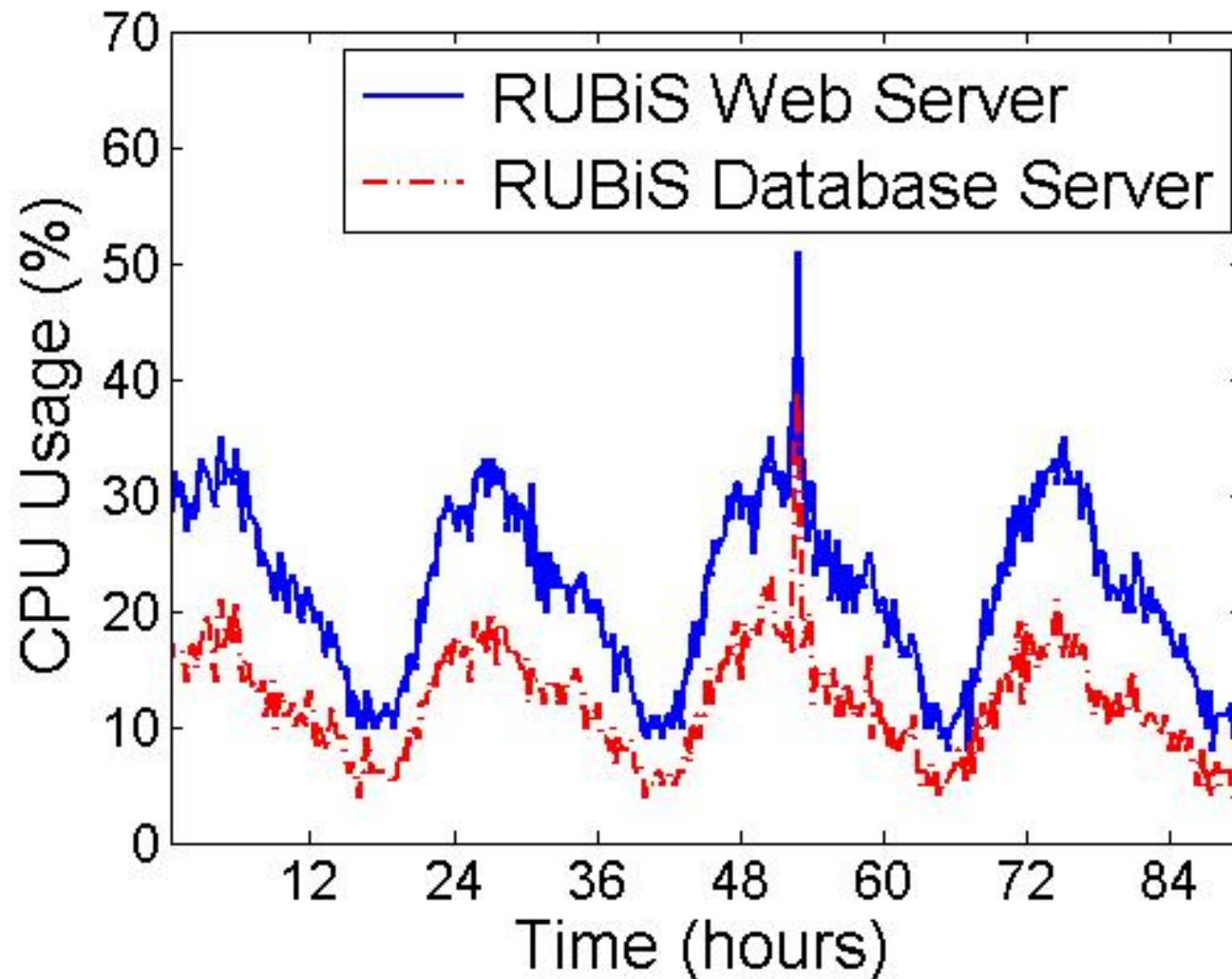
# System Architecture



- Collect historical resource usage
- Online prediction for future usage
- Scaling: dynamically adjust resource caps based on prediction
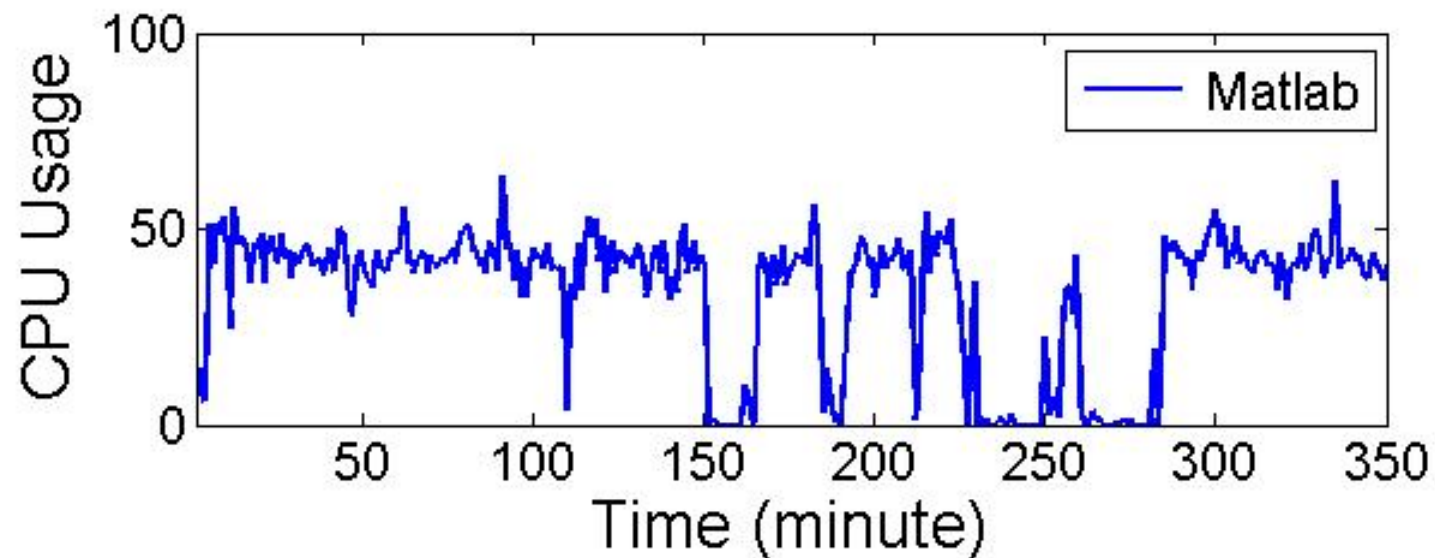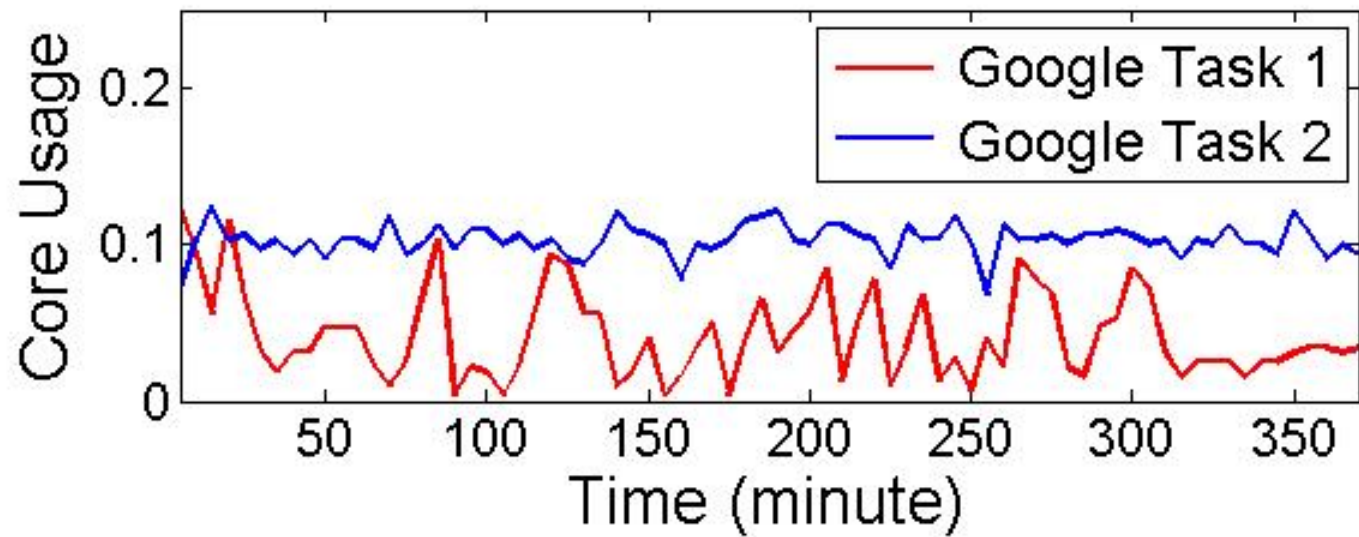- Black-box: no need for application knowledge

# Research Challenges

- We need prediction to do scaling

  - Fast reaction to dynamic workload changes

- A good prediction algorithm is hard

  - Resource demands are dynamic

  - Need to avoid both over- and under-estimation

  - Accommodate different workloads

  - Low overhead to manage large cloud system
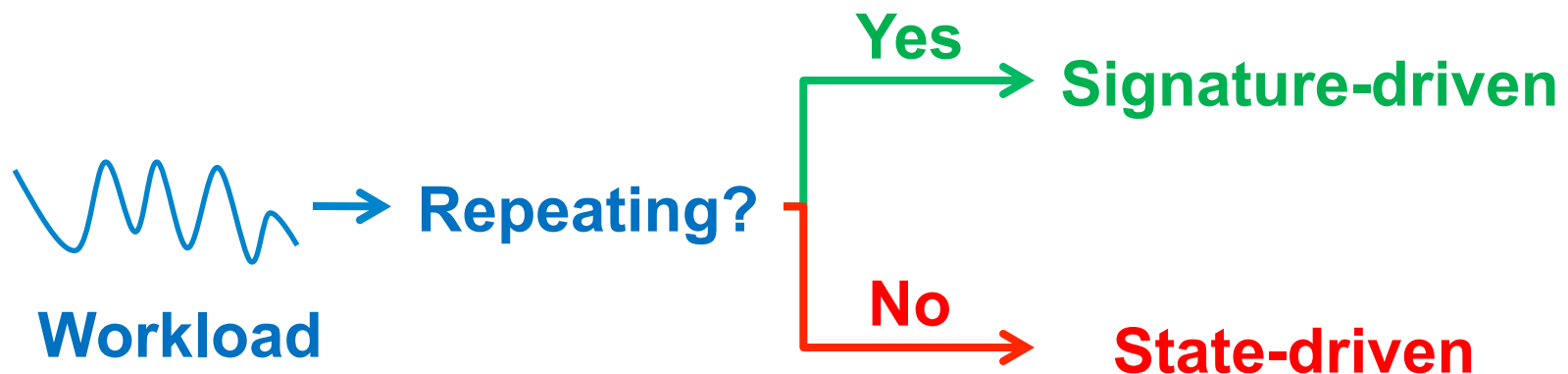
# Repeating Workload

# Non-Repeating Workload
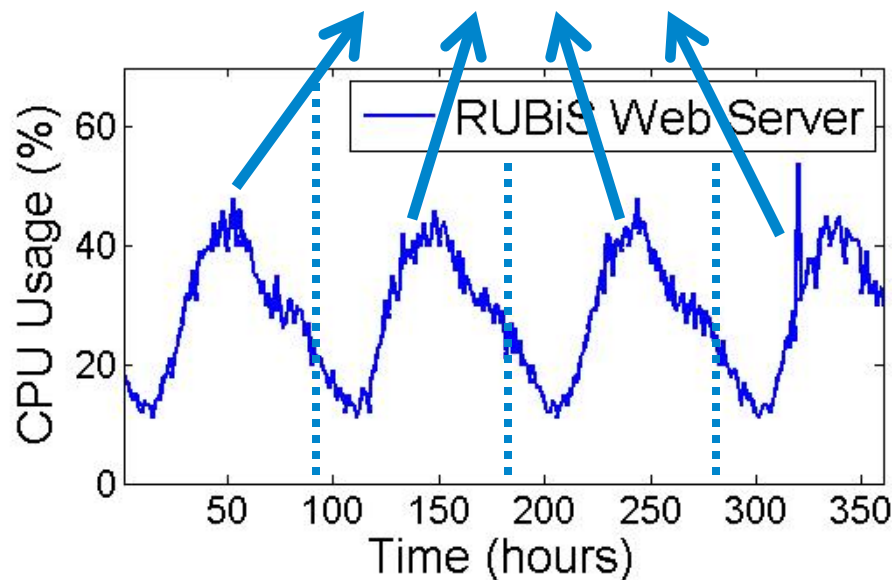
# Our approach: Hybrid Prediction Model

- For repeating workload: signature-driven
  - Use fast Fourier transform (FFT) to find repeating patterns for prediction
  - Lower overhead

- For non-repeating workload: state-driven
  - Use Markov chain model for prediction
  - Higher overhead, accommodate different workloads

**Yes**

**Signature-driven**

**Workload** → **Repeating?**

**No**

**State-driven**

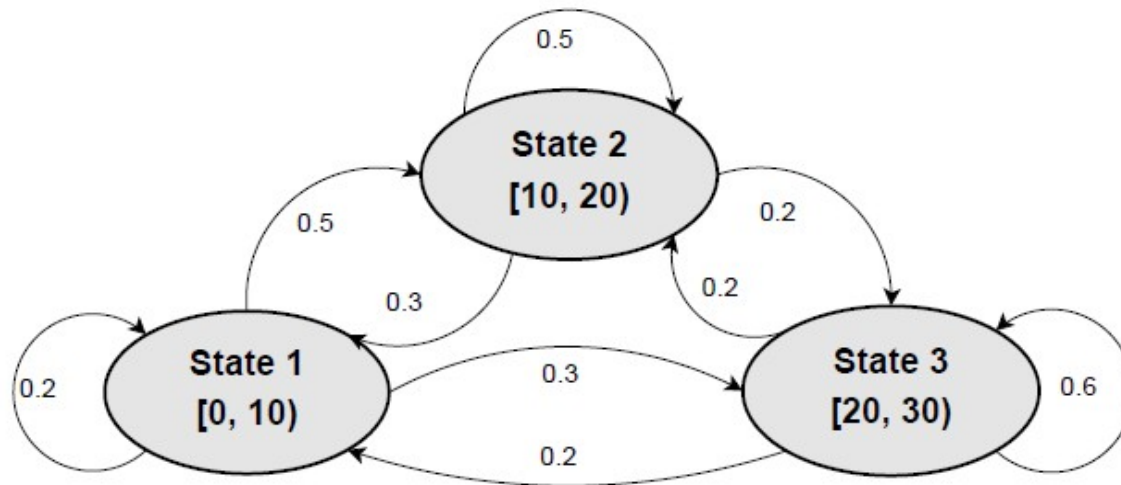13

# Signature-Driven Approach

- Signature Extraction
  - Use fast Fourier transform (FFT) to discover dominant frequency
  - Use comprehensive difference check to validate repeating pattern
  - Signature alignment to predict next value

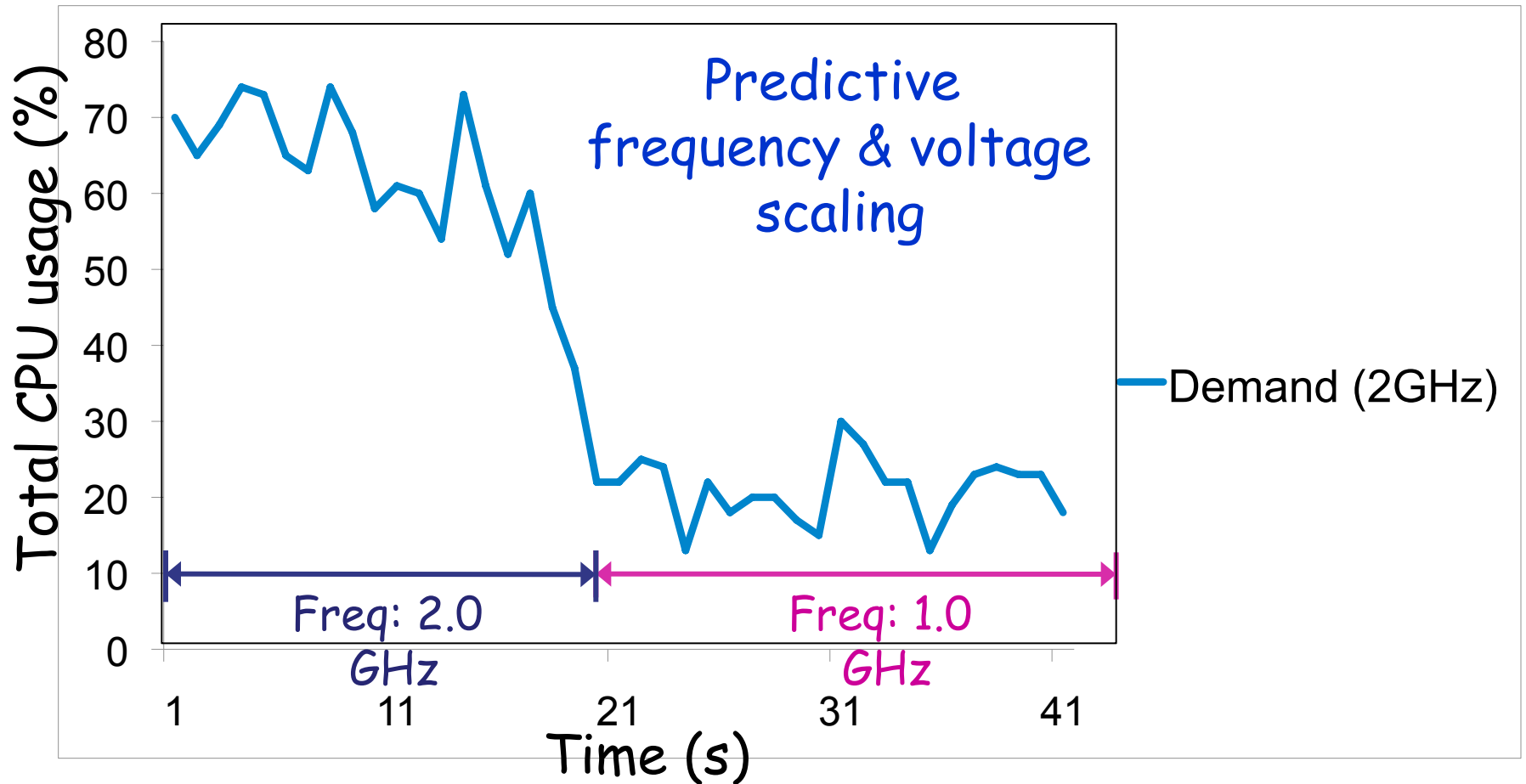**Similar? (mean value, correlation)**

# State-Driven Approach

- Repeating patterns sometimes do not exist
- Uses discrete-time Markov chain to do short-term prediction
  - Dynamically learn transition probability matrix to predict future values

state transition matrix:



$$\begin{bmatrix} 0.2 & 0.5 & 0.3 \\ 0.3 & 0.5 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}$$
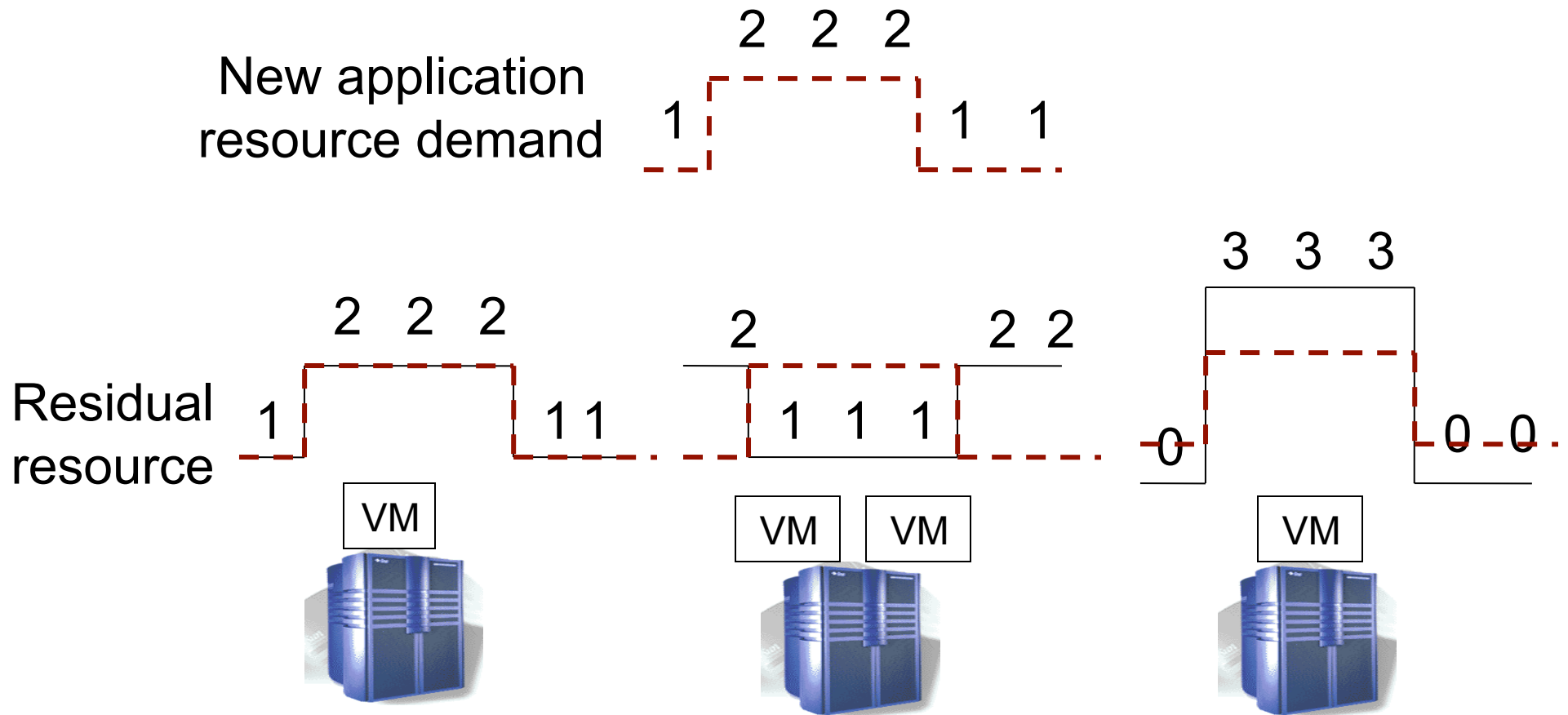
# Energy Saving via Resource Scaling

# Problem: Application Consolidation

- Application consolidation is needed for
  - Use fewest hosts to run all applications
  - Reduce resource provisioning cost
  - Reduce energy consumption

- Application consolidation needs to achieve
  - Meet dynamic resource requirements of different applications
  - Alleviate negative sharing impact
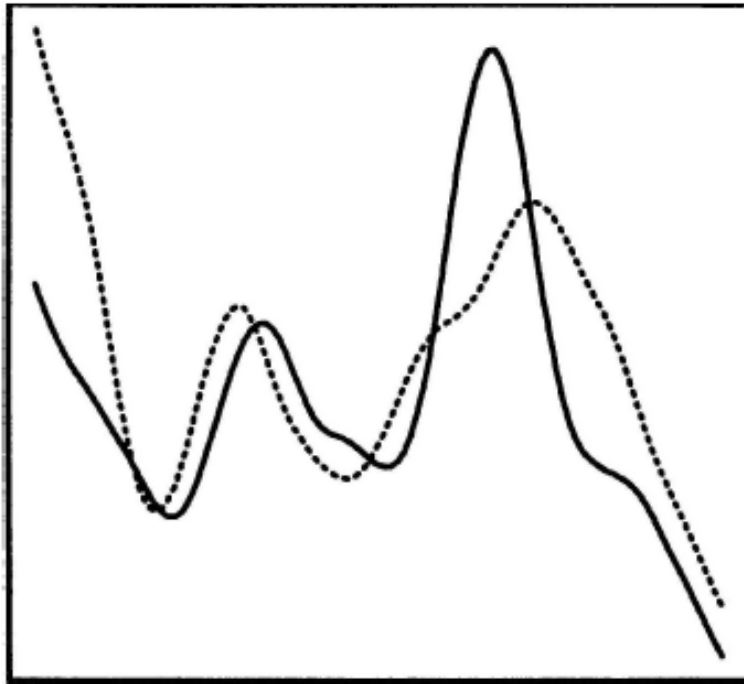  - Balance resource utilization

# A Case for Fine-Grained Patten Matching

# PAC: Pattern-driven Consolidation

- Use time series to capture fine-grained signatures of dynamic applications
  - Dynamically discover repeating patterns using signal processing techniques

- Perform dynamic VM consolidation using signatures
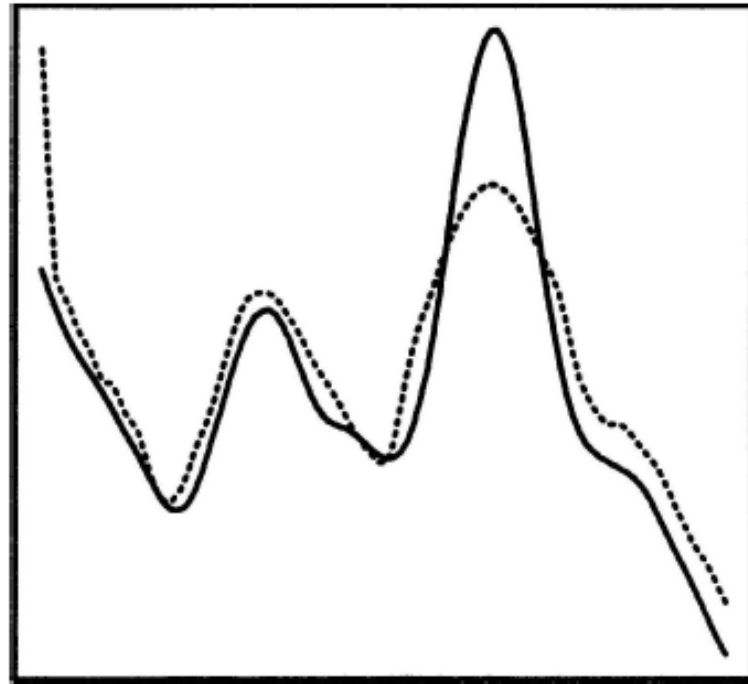  - Use time series indexing and dynamic time warping (DTW) to achieve fast and robust signature matching

# PAC: Pattern-driven Consolidation

- Challenges

  - Fine-grained pattern matching in distributed computing environments

  - Scalable application consolidation

  - Pattern heterogeneity

- Solutions

  - Use dynamic time warping (DTW) to achieve robust matching

  - Use time series indexing to achieve fast signature matching

  - Pattern alignment

# Dynamic Time Warping (DTW)



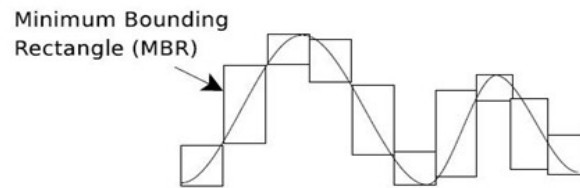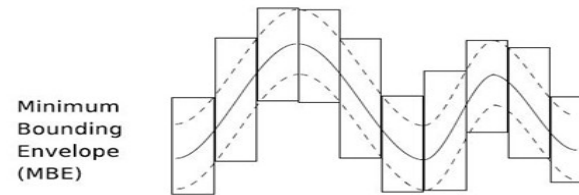(a) Two different time-series that are out of phase by little

(b) After alignment of the time-series by dynamic time warping

DTW Phase alignment example

# Signature Indexing for Pre-Filtering

- ## One-dimensional case for simplicity:
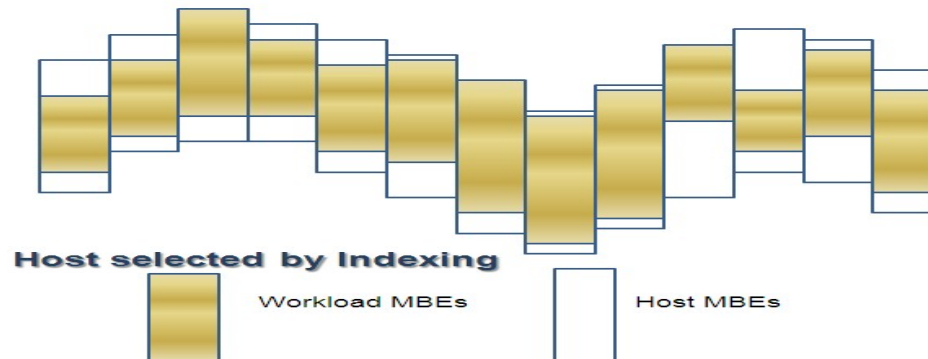  - Split the time series to MBRs expand to be MBE



Minimum Bounding Rectangle (MBR)

Minimum Bounding Envelope (MBE)

(a) Signature indexing using MBR.

(b) Signature Pre-filtering.

  - Matching lower-bound of the load MBE is lower than the upper-bound of the host resource MBE



**Host selected by Indexing**

Workload MBEs          Host MBEs

  - Admission Control: Pass if more then 65% matching in all MBEs (adjustable threshold)

# Pattern-Driven Consolidation

Collect VM resource usage in the recent past
Extract signature patterns and predict for future.
Use Indexing and DTW to match VMs and hosts
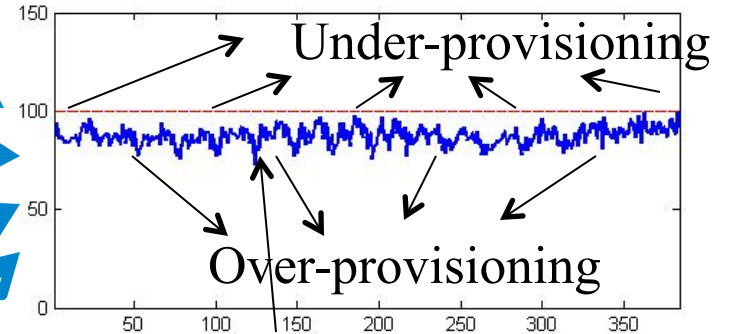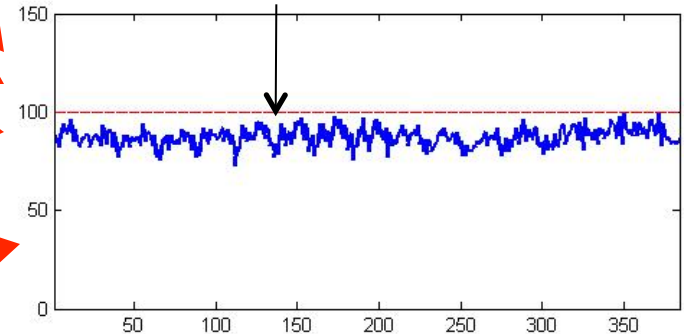Perform VM consolidation based on matching result

# Consolidation example

# System Consolidation Performance



Response time comparison

SLO violation comparison