

Data Mining for Sustainable Data Centers

Manish Marwah
Senior Research Scientist
Hewlett Packard Laboratories
manish.marwah@hp.com



Overview

- Sustainability and Data Centers
- Data mining applications
 - Chiller operation characterization
 - PV prediction
 - Anomaly detection



Motivation

Industry challenge:

Create technologies, IT infrastructure and business models for the low-carbon economy

IT industry

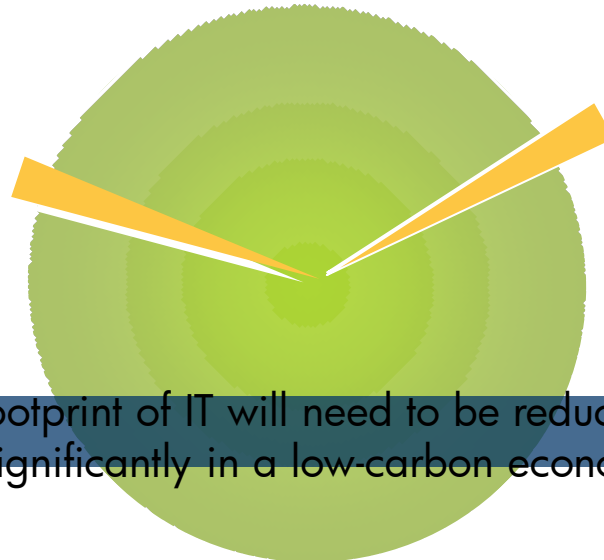
2%

Total carbon emissions



Aviation

2%



The footprint of IT will need to be reduced quite significantly in a low-carbon economy.

Sustainability

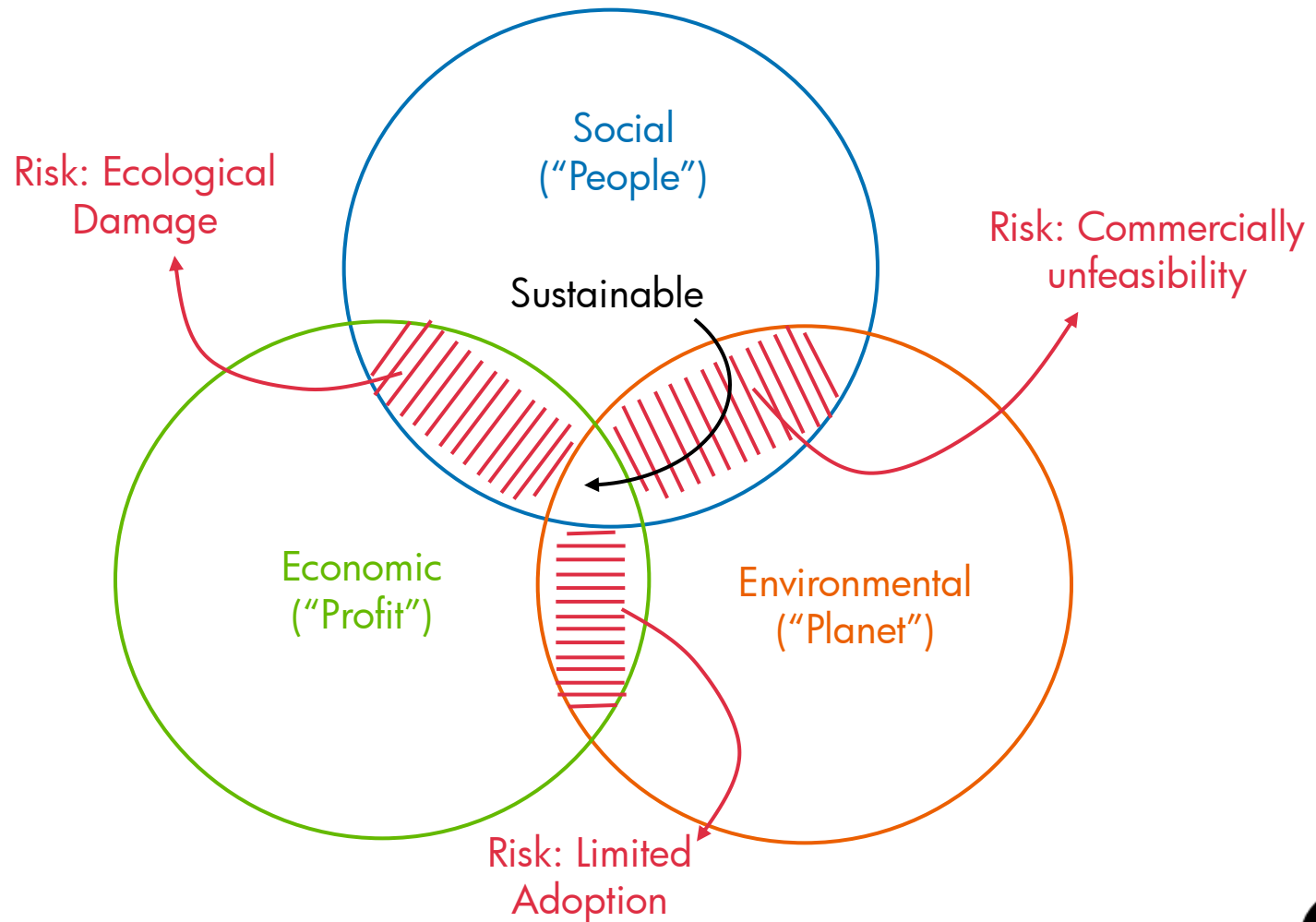
“sustainable development is development that meets the needs of the present without compromising the ability of future generations to meet their own needs”

the [Brundtland Commission](#) of the [United Nations](#), 1987



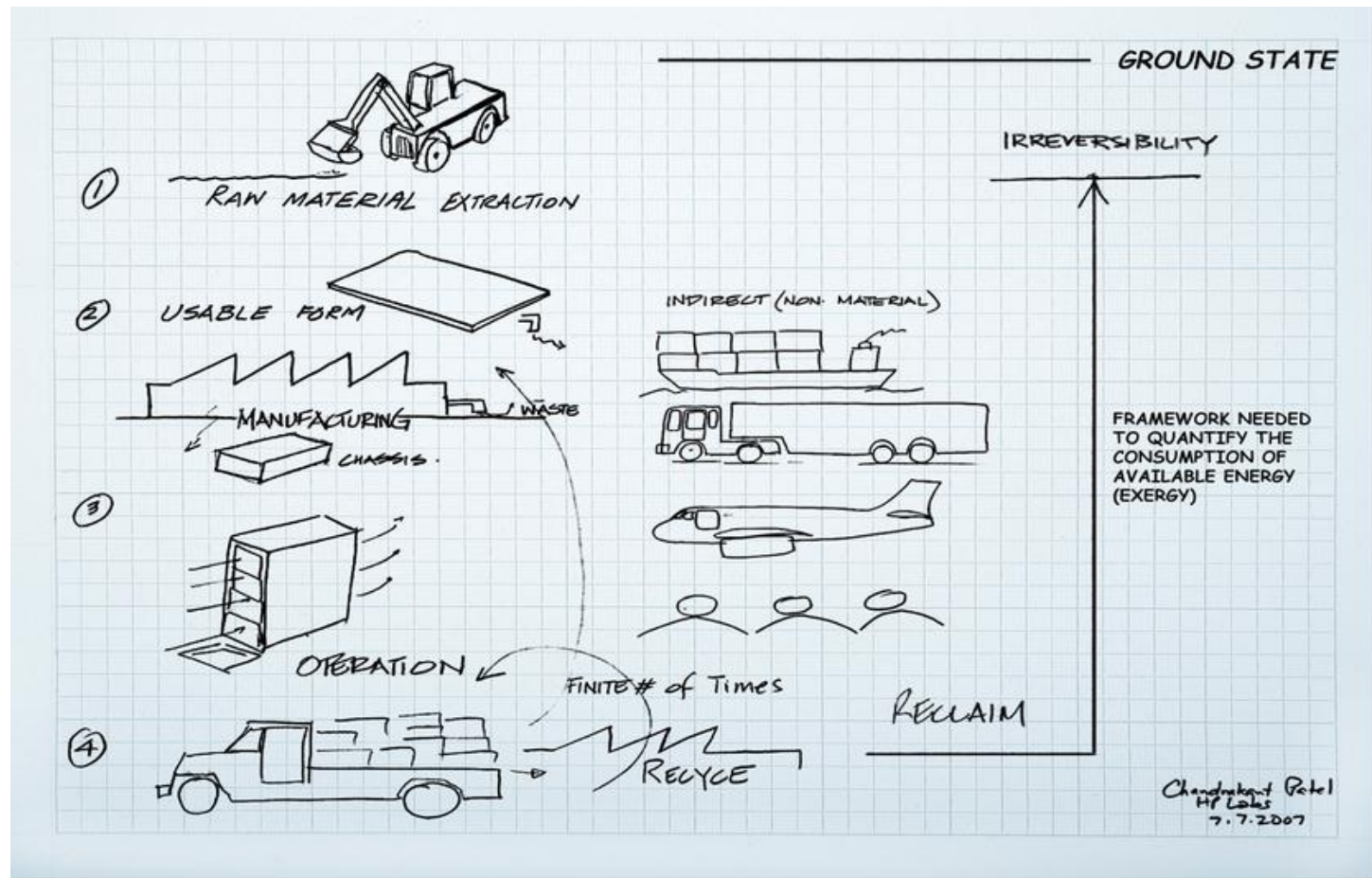
Sustainability

What do I mean by "sustainability"?



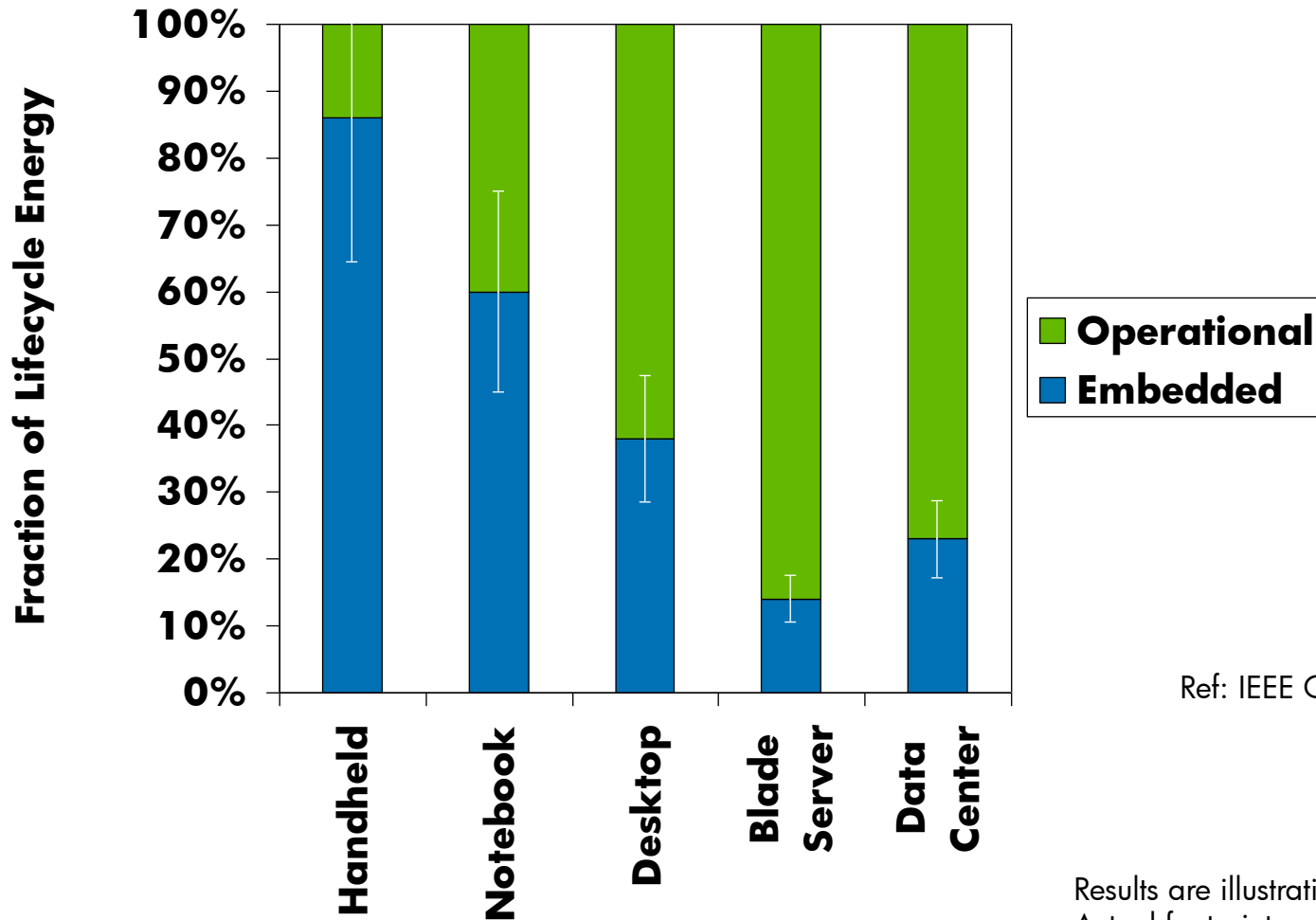
Environmental Sustainability

- Impact factors (e.g., carbon, water, toxicity, etc.)
- Life Cycle View



Sustainable Data Centers

Lifecycle Assessment



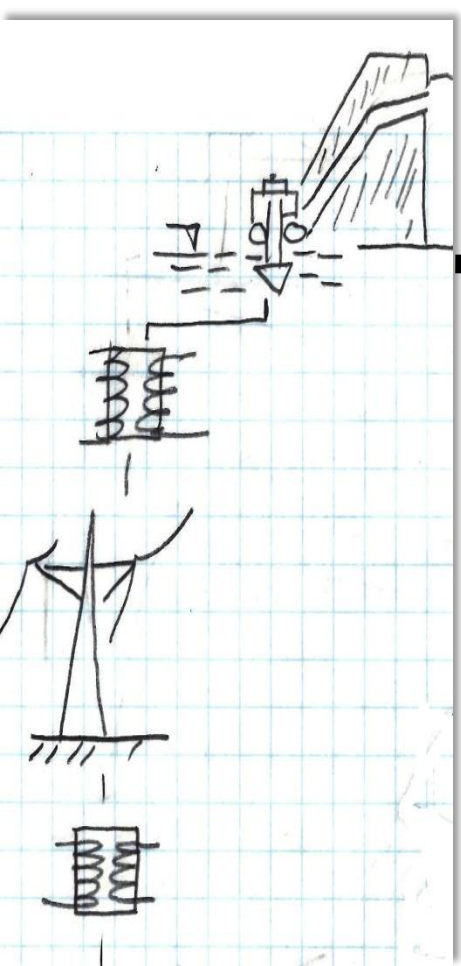
Ref: IEEE Computer 2009

Results are illustrative only.
Actual footprint may differ.



Cloud Data Center

Supply and Demand Side



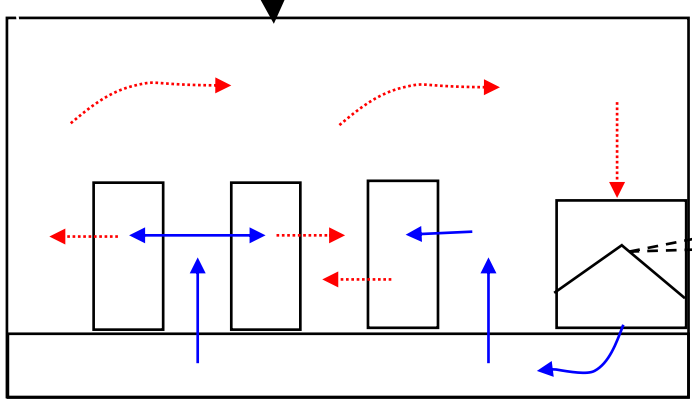
Power

Switch Gear

UPS

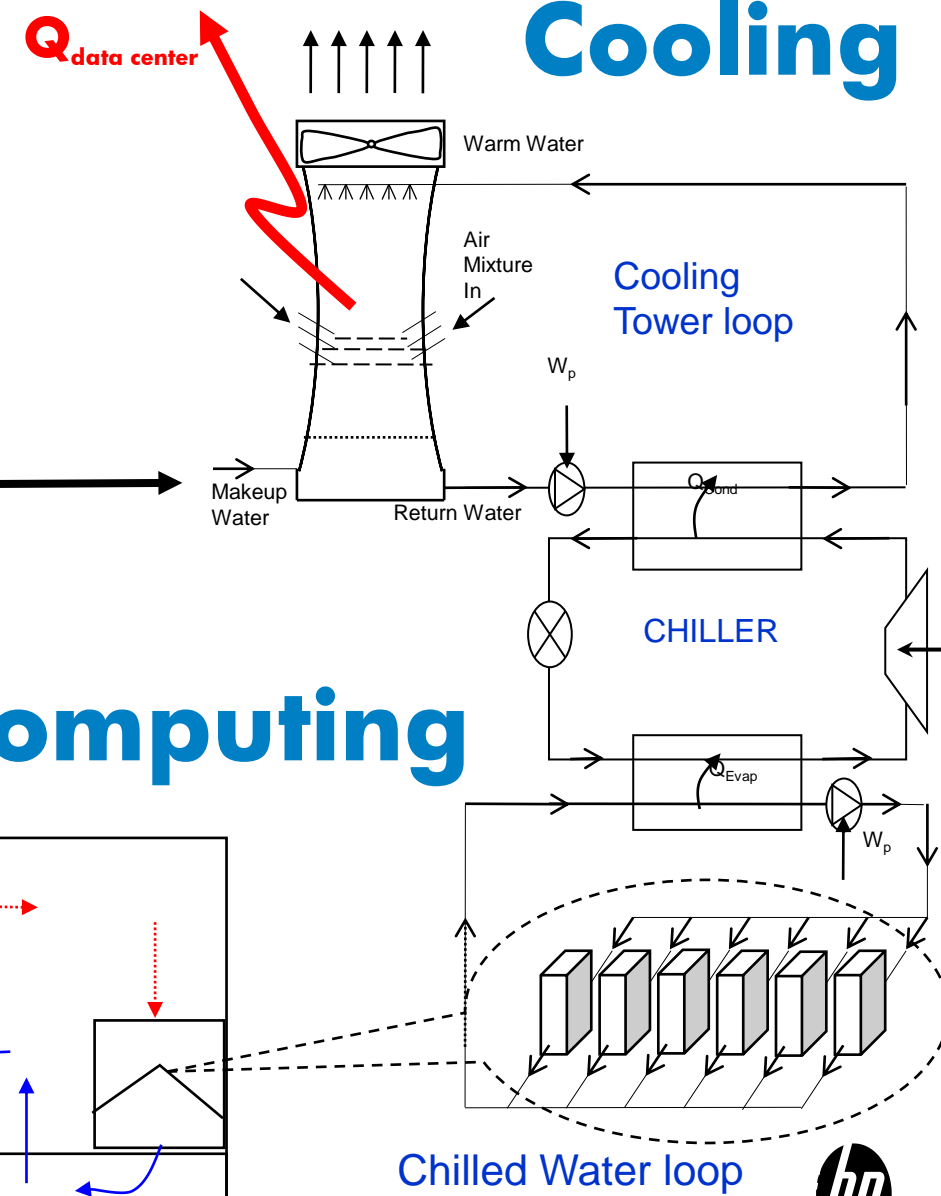
PDU

Computing



Data Center

Cooling

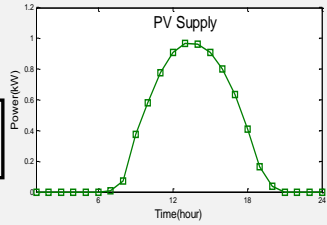
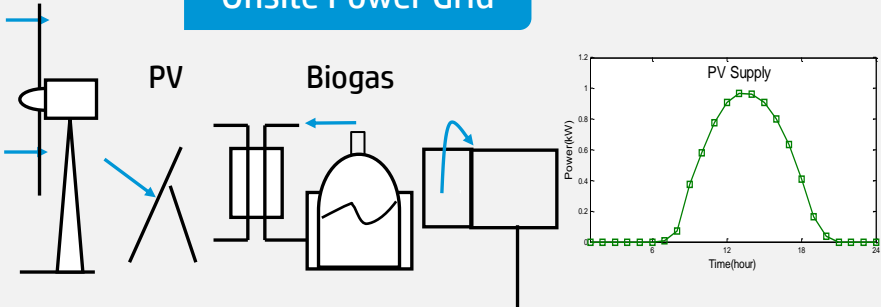


Chilled Water loop

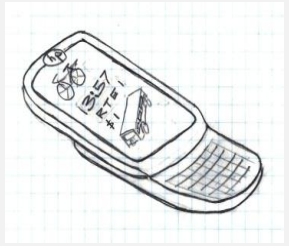
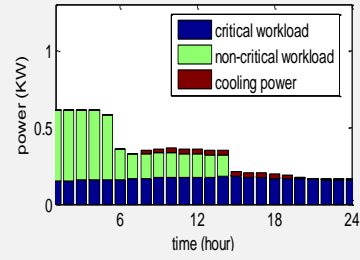


Supply and Demand in a Data Center

Onsite Power Grid

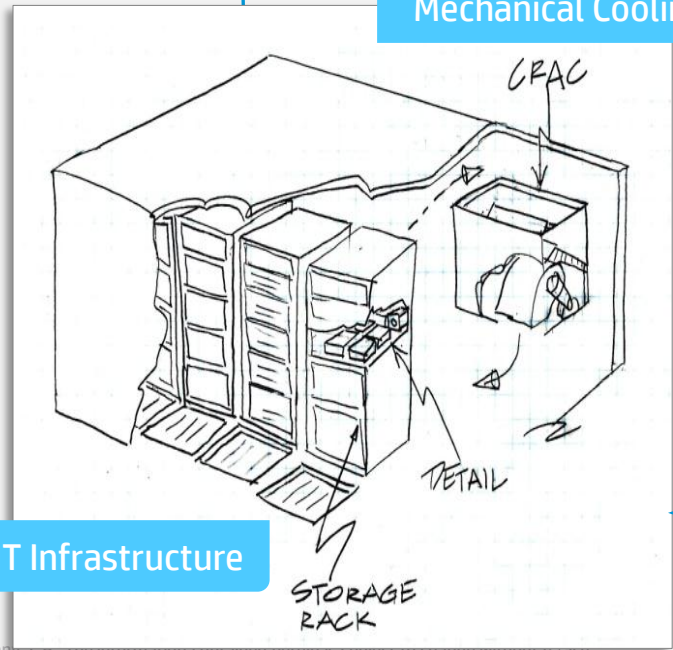


IT Services



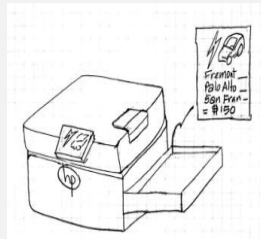
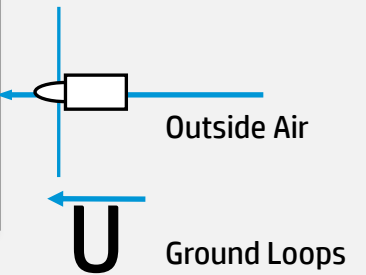
Ecosystem of Clients

Mechanical Cooling



IT Infrastructure

Local Cooling Grid



Supply

Demand



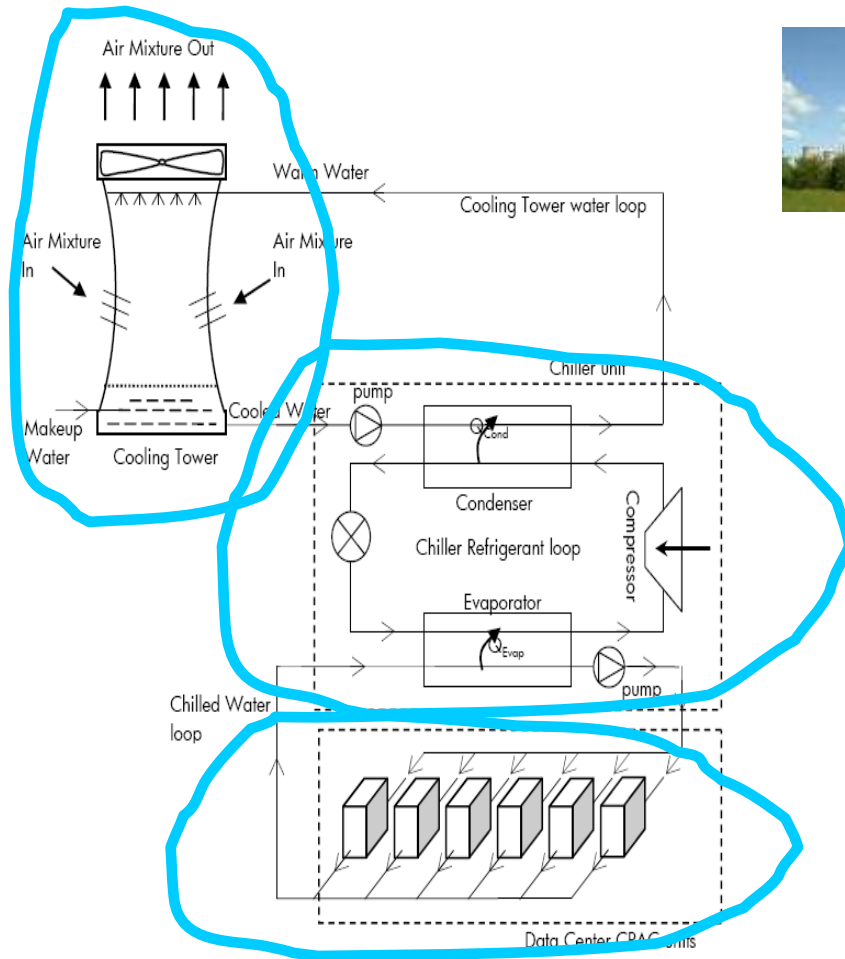
Sustainable Operation and Management of Chillers using Temporal Data Mining (KDD '09)

- Data Centers
 - Cooling Infrastructure
- Problem Statement
- Prior Work
- Our Approach
 - Symbolic representation
 - Event encoding
 - Motif mining
 - Sustainability characterization
- Experimental Results
- Summary



Data Center Cooling Infrastructure

Consumes from 1/3 up to 1/2 of total power consumption



Cooling Towers

Chiller Unit

Water Return (T_{in})



Water Supply (T_{out})



Computer room air-conditioner (CRAC)



Ensemble of Chillers

- Challenging to operate efficiently
 - Complex physical system
 - Dynamic
 - Heterogeneous
 - Inter-dependencies
 - Many constraints
 - Accurate models not available
 - Rapid cycles undesirable – reduce lifespan
- Domain experts determine settings based on heuristics
- Can it be automated through a data-driven approach?



Chiller Ensemble

- Which unit to turn ON/OFF?
- At what utilization?
- How to handle increase/decrease in cooling load?



Problem Statement

- Given the following chiller time series
 - utilization levels
 - power consumption
 - cooling loads
- Is it possible to determine which operational settings are more energy efficient?
- And then use this information to advise data center facility operators



Some Terminology

- IT cooling load
- Chiller utilization
- Chiller power consumption
- Coefficient of performance (COP)

$$\frac{\text{Cooling Load}}{\text{Power consumption}}$$



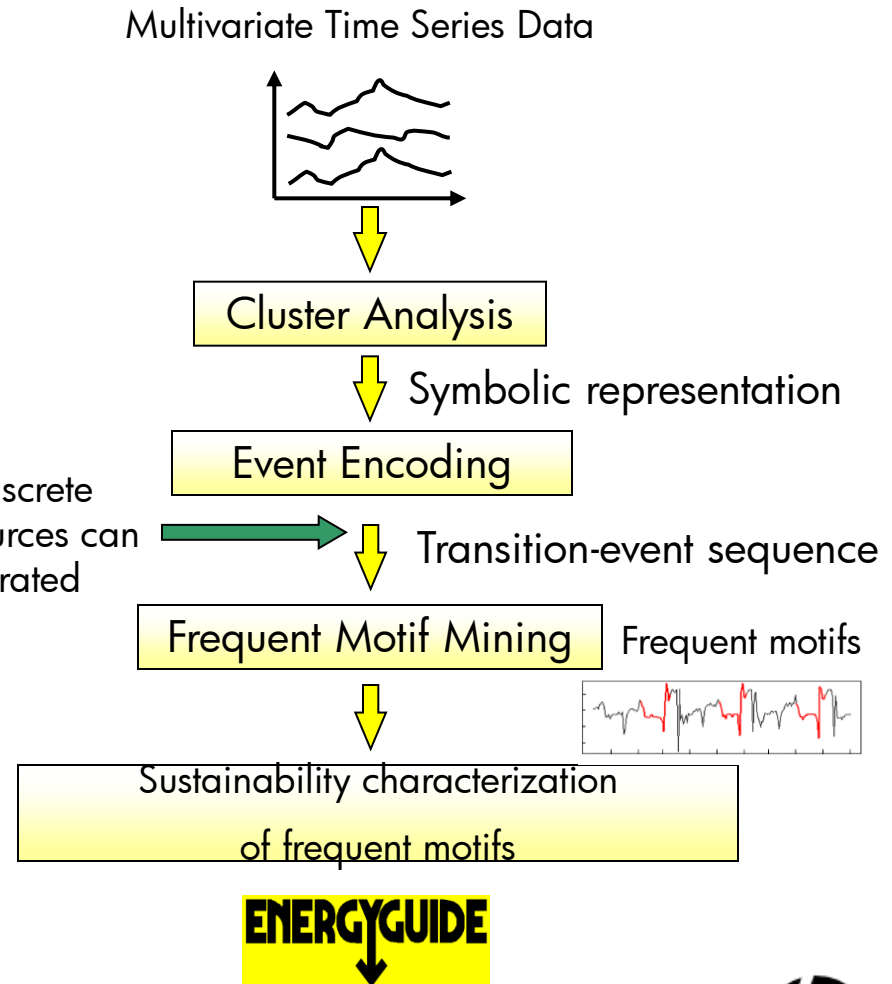
Prior Work

- Classical approaches to model time series data
 - Principal component analysis
 - Discrete Fourier transforms
- Discrete representations: SAX [Keogh et al.]
- Motifs: Repeating subsequences [Yankov et al.]



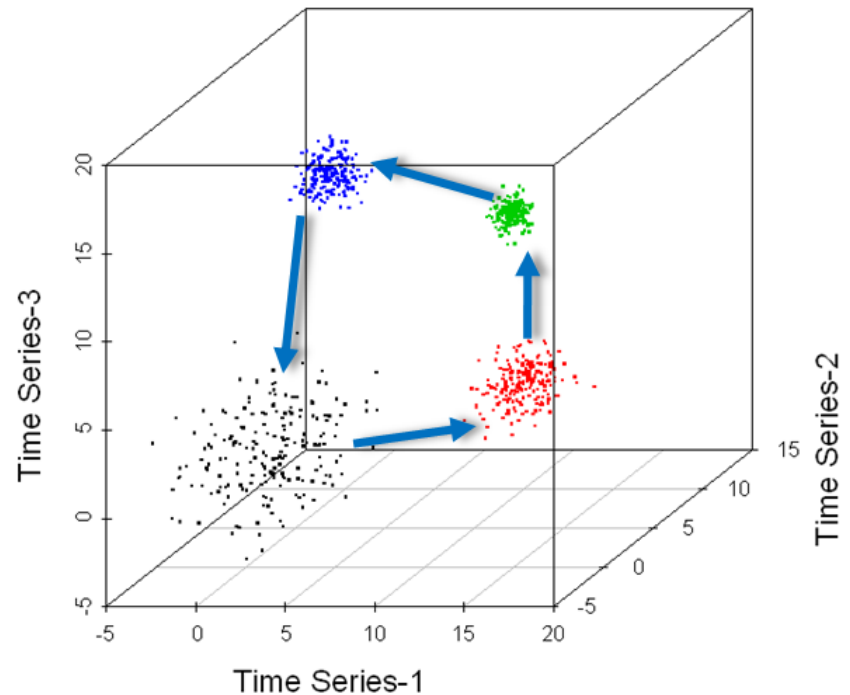
Our approach

- Goal: Sustainability characterization of multi-variate time series data
 - Chiller utilization data
- Four Main Steps
 - Symbolic representation
 - Event encoding
 - Motif mining
 - Sustainability Characterization



Clustering

- Individual vector:
Utilization across all chiller units
- Raw Data: Sequence of such vectors
- Perform k-means clustering
- Use cluster labels to encode multi-variate time series



Some Definitions

- Event Sequence

$$\langle (E_1, t_1), (E_2, t_2), \dots, (E_N, t_N) \rangle$$

E_i = Event type t_i = Time of occurrence

$$\langle (A,1), (B,3), (D,4), (C,6), (A,12), (E,14), (B,15), (D,17), (C,20), (A,21) \rangle$$

- Episode
 - Ordered collection of events occurring together

$$(A \rightarrow B \rightarrow C)$$

- Episode occurrence
 - Events same ordering as episode in the **data**.

$$\langle (A,1), (B,3), (D,4), (C,6), (E,12), (A,14), (B,15), (C,17) \rangle$$

- Motifs
 - Frequently occurring episodes



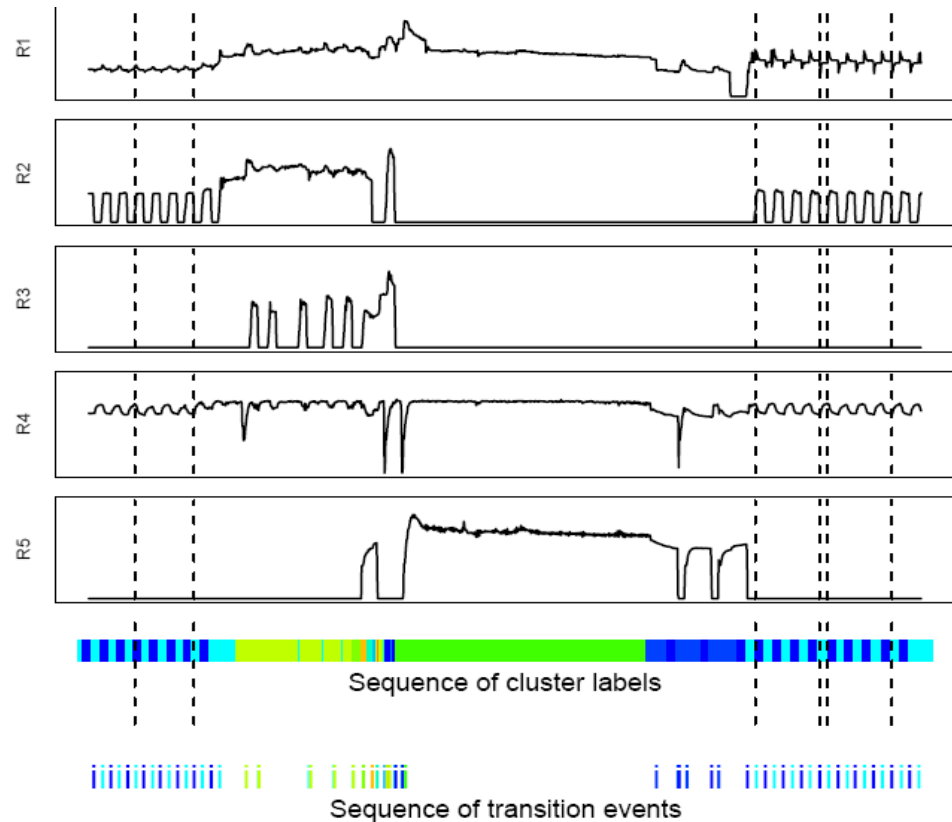
Redescribing time series data

- Perform run-length encoding:
 - Note transitions from one symbol to another
- Higher level of abstraction
 - Transition events

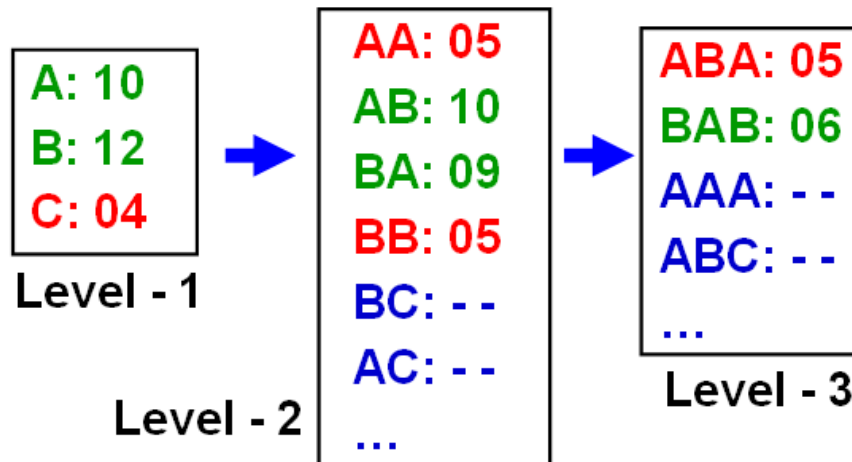
Symbol Sequence : d d d b a c c d d d d c b

↓

Event Sequence : $\langle (d-b, 4), (b-a, 5), (a-c, 6), (c-d, 8), (d-c, 12), (c-b, 13) \rangle$

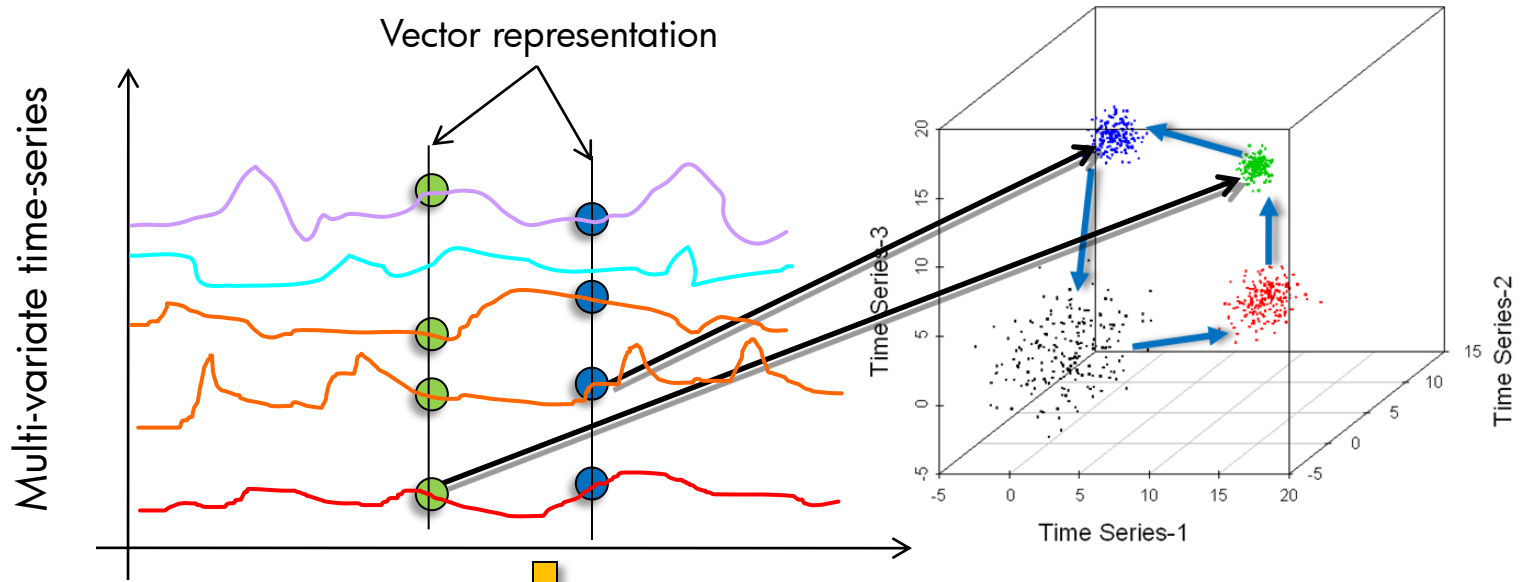


Level-wise (Apriori-based) motif mining



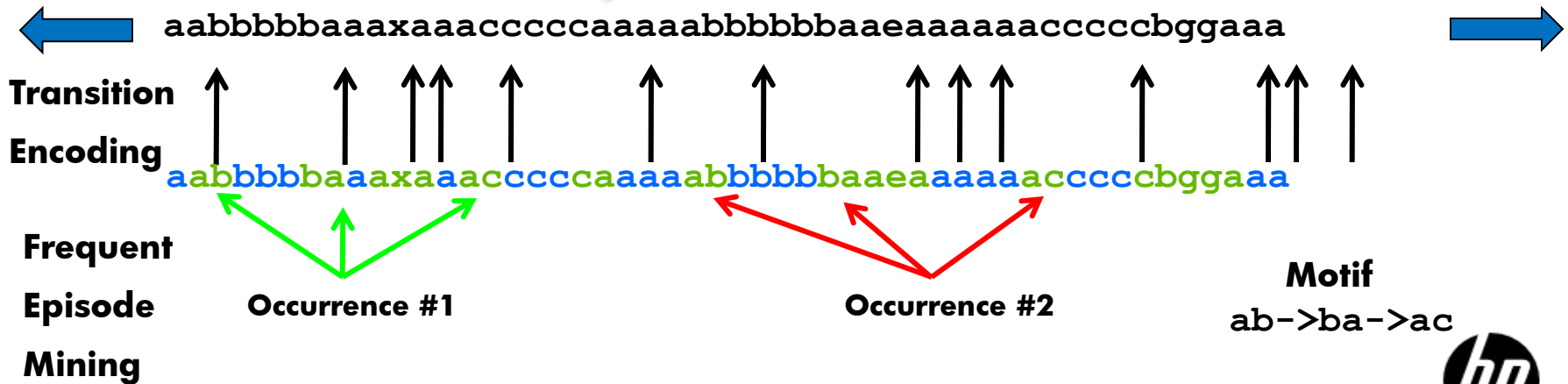
Candidate generation
followed by counting

Methodology Summary



Clustering

Discrete representation of chiller ensemble time-series



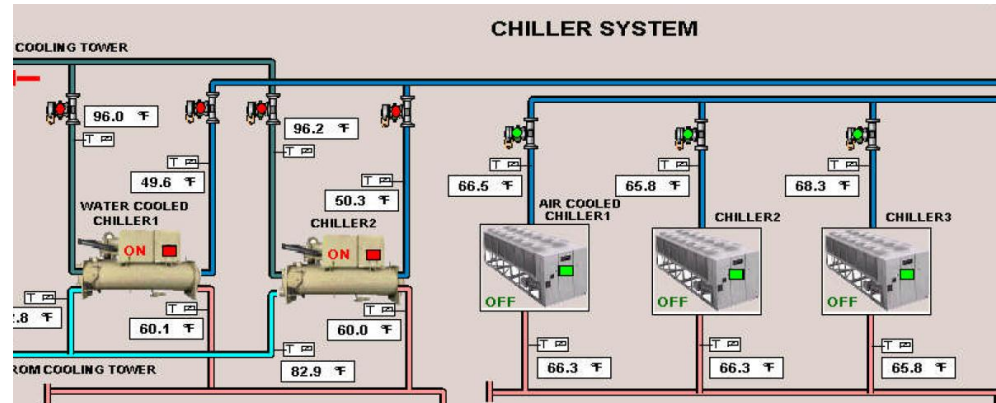
Sustainability characterization of Motifs

- Average motif COP (coefficient of performance)
 - Indicates cooling efficiency of a chiller unit
 - $$\text{COP} = \frac{\text{IT Cooling Load}}{\text{Power consumed}}$$
- Frequency of oscillations of a motif
 - Impacts chiller lifespan
 - Normalized number of mean-crossings

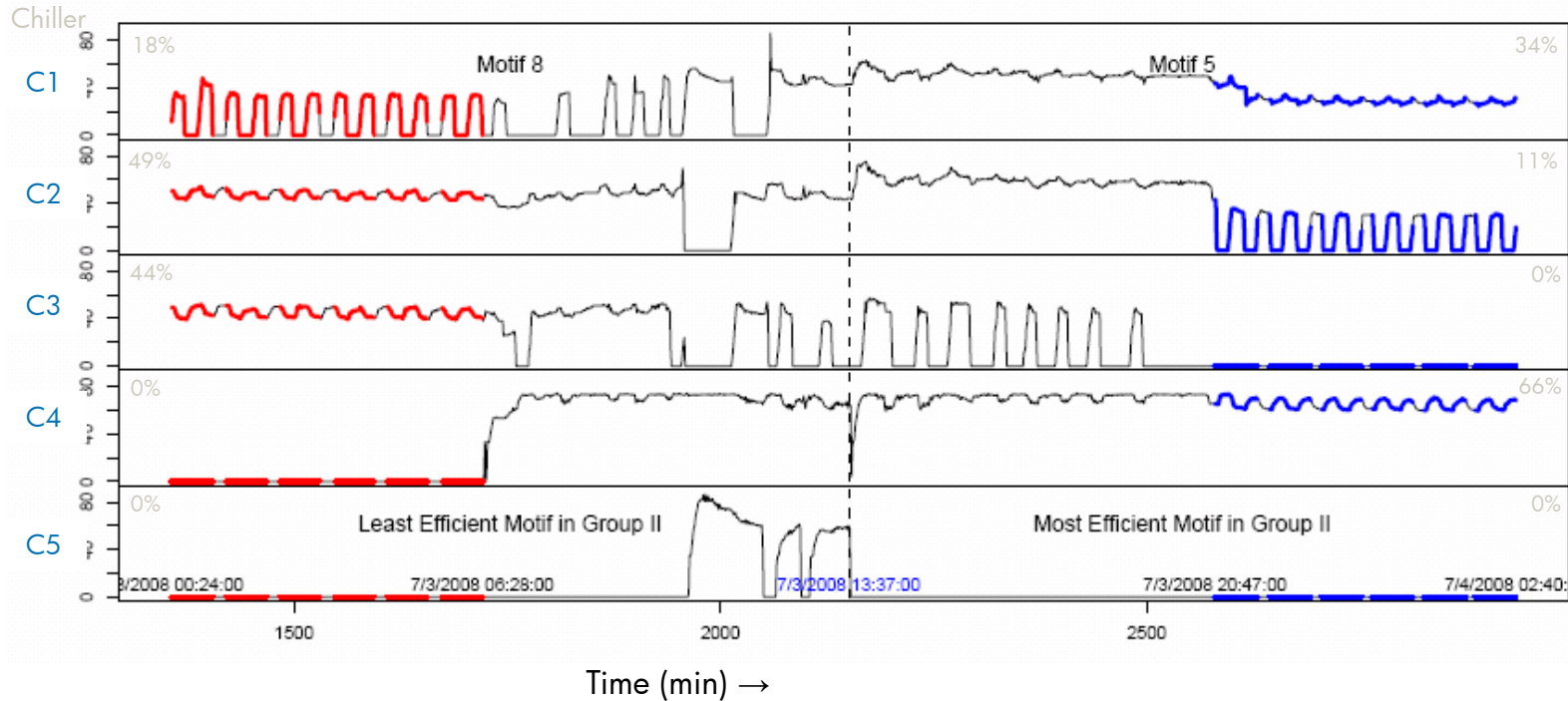


Experimental Results

- Data
 - From HP R&D data center in Bangalore
 - 70,000 sq ft
 - 2000 racks of IT equipments
 - Ensemble of five chiller units
 - 3 air cooled chillers
 - 2 water cooled chillers
 - 480 hours of data
 - July 2 – 7, Nov 27 – 30, Dec 16 – 26, 2008
- 22 motifs found in the data



Two Interesting Motifs



C1, C2, C3 → Air cooled

C4, C5 → Water cooled

— Motif 8
— Motif 5

	Motif 8	Motif 5
COP	4.87	5.40
Units operating	3 air-cooled	2 air-cooled, 1 water cooled



Potential Savings

	Load (KW)		Most Efficient Motif	Least Efficient Motif	Potential Power Savings	
	Ave.	Std			KW	%
Group II	2089	35	5	8	41	9.83%

- Annual saving from operating in Motif 5 instead of Motif 8
 - Cost savings = \$40,000 (~10%)
 - Carbon footprint savings = 287,328 kg of CO₂



Summary

- Data center chillers consume substantial power
 - Ensemble of chillers – part of data center cooling infrastructure – are challenging to operate energy efficiently
- Mine and characterize motifs
 - Symbolic representation
 - Event encoding
 - Motif mining
 - Sustainability characterization
- Demonstrated our approach on data from a real data center – indicates significant potential energy savings



The Net-Zero Energy Data Center

Implementation in Palo Alto

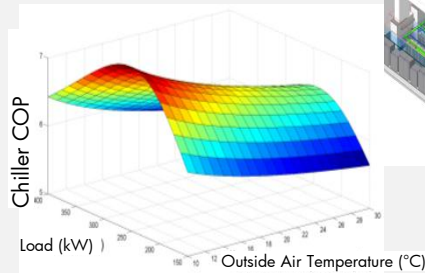
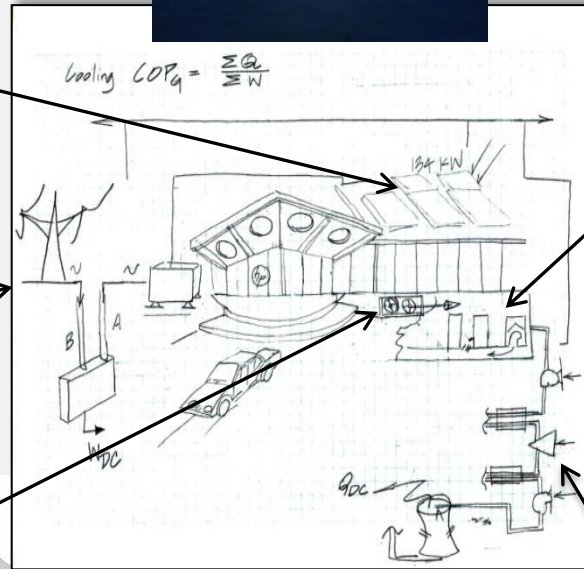
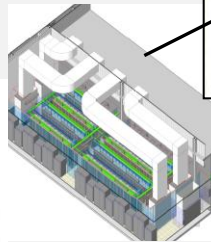
PV micro grid



Data center



Outside air



Cooling infrastructure power demand

Data center supply side

Data center demand side



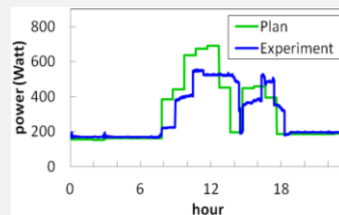
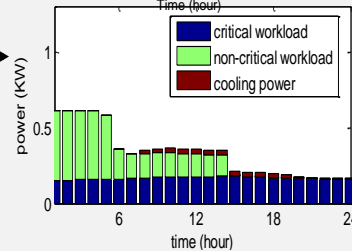
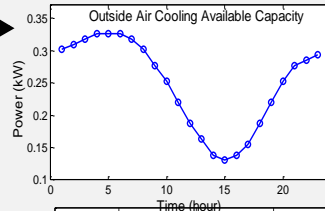
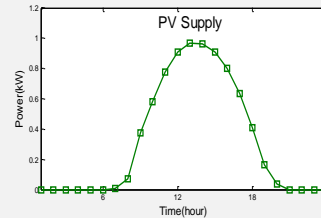
Net-Zero Energy Methodology and Integration

Prediction

Supply-side Prediction

- Renewable power prediction
- Cooling capacity prediction

IT Demand Prediction



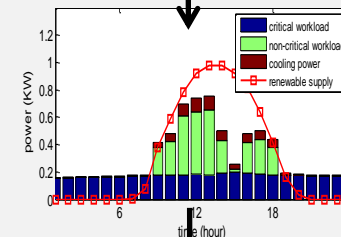
Planning

DC Operation Objectives:

- Net-zero energy operation
- Maximize use of renewable energy
- Minimize dependability on Grid

IT Workload Planning

- Integrate Supply and Demand Side
- IT Demand Shaping
- Power capping



Verification and Reporting

Measurement Verification

Execution

Dynamic IT Provisioning

Dynamic Cooling Provisioning

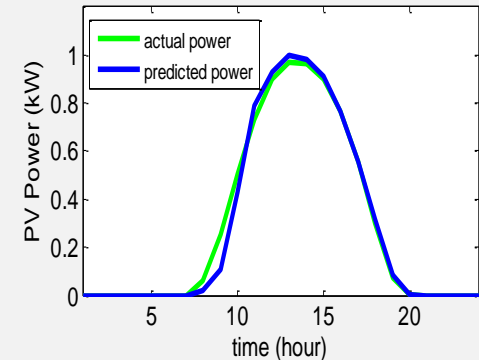


Prediction: Summary

PV Supply Prediction

- Search for most “similar” days in the recent past
- Hourly generation estimated from corresponding hours of “similar” days

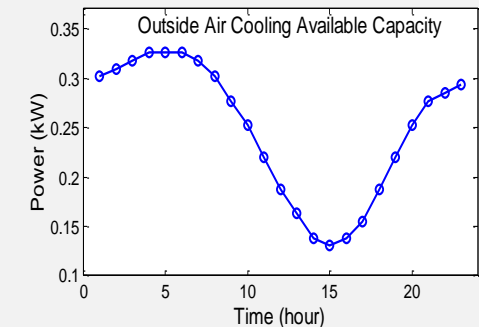
Ref: P. Chakraborty, M. Marwah, M. Arlitt, N. Ramakrishnan, Fine-grained Photovoltaic Output Prediction Using a Bayesian Ensemble, in Proceedings of the 26th Conference on Artificial Intelligence (AAAI'12), Toronto, Canada, July 2012



Cooling Capacity Prediction

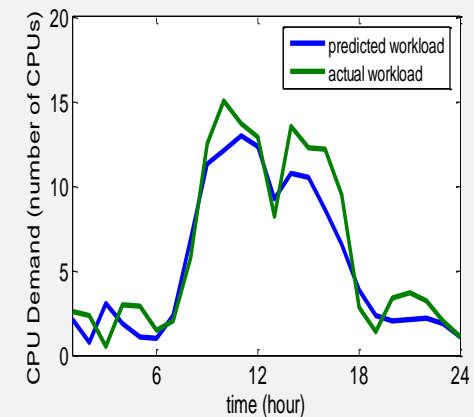
- End-to-End Energy Modeling

Ref: Breen, T.J. et. al. “From Chip to Cooling Tower Data Center Modeling: Validation of Multi-Scale Energy Management Model”, Proceedings of Itherm, June 2012



IT Workload Prediction

- Perform a periodicity analysis (e.g., Fast Fourier Transform)
- Use an auto-regressive model to predict workload from historical data



Fine grained PV Prediction using Bayesian Ensemble

- **Motivation**

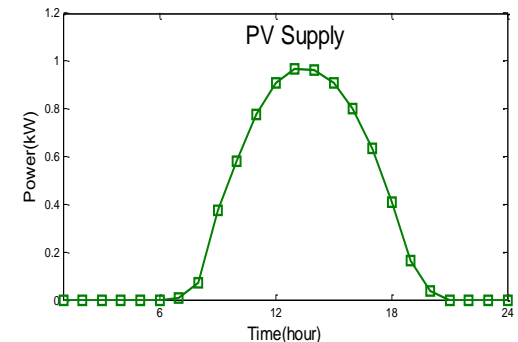
- Integration of renewable sources is an important goal of the smart grid effort
- PV output is variable and intermittent
- Knowledge of future PV output enables demand-side management and “shaping” in data centers

- **Problem addressed**

- Predict PV output for the next day

- **Data**

- Historical PV output data for about 9 months from the HPL Palo Alto site
- Weather data



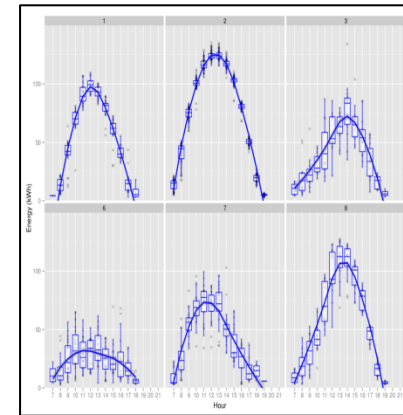
(AAAI 2012)



Fine grained PV Prediction using Bayesian Ensemble

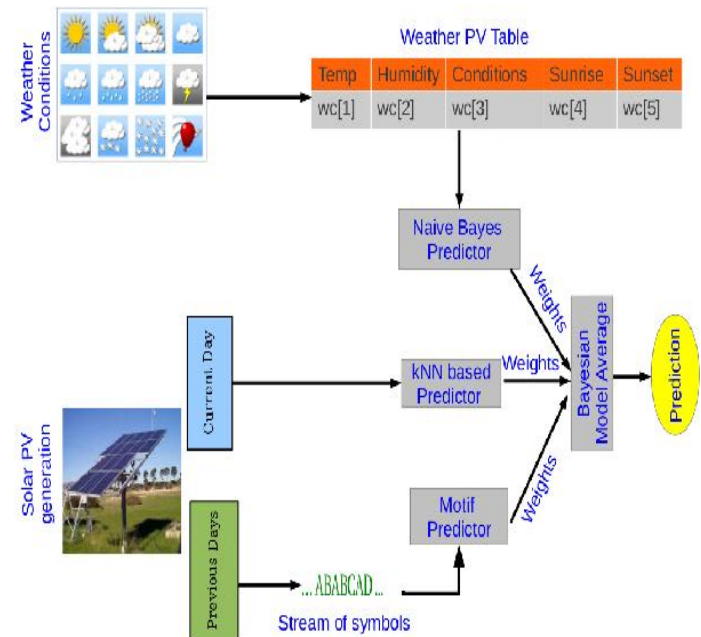
- **Approach**

- Extract daily profiles from training data
- Use ensemble of predictors
 - Naïve Bayes
 - K-NN
 - Motif based
- Perform Bayesian model averaging



- **Results**

Method	Testing Error		
	Per. Abs. Error	Per. RMS Error	Rel. Abs. Error
PreviousDay	20.54	20.65	20.81
ARWeather	18.54	18.31	19.73
Stagewise	12.77	12.68	15.66
Ensemble2	10.04	10.01	10.01
Ensemble3	8.13	8.21	8.34

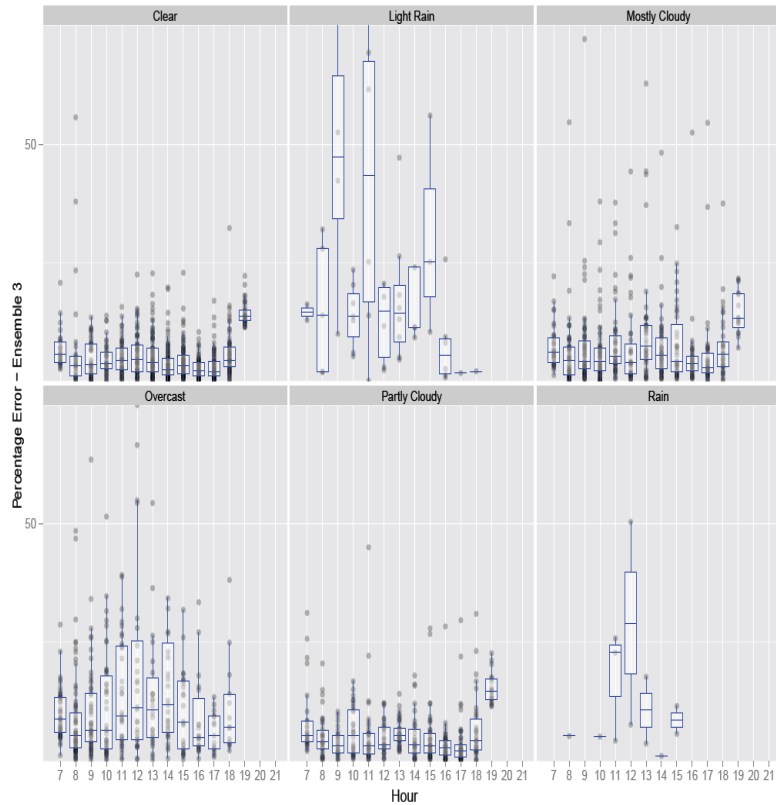


(AAAI 2012)

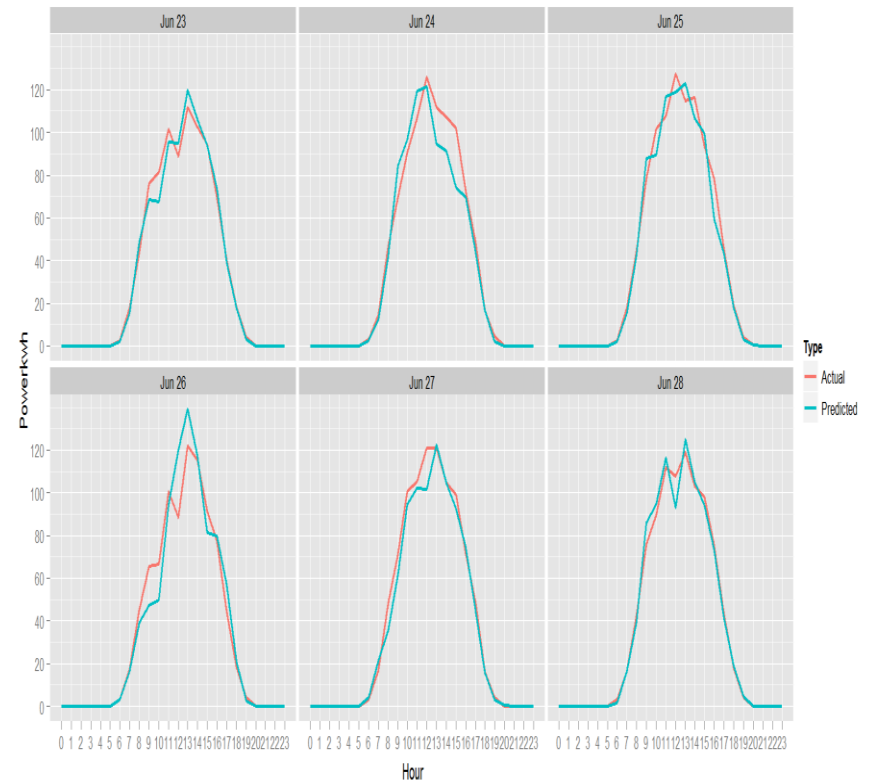


Fine grained PV Prediction using Bayesian Ensemble

- Results**



Error by weather condition



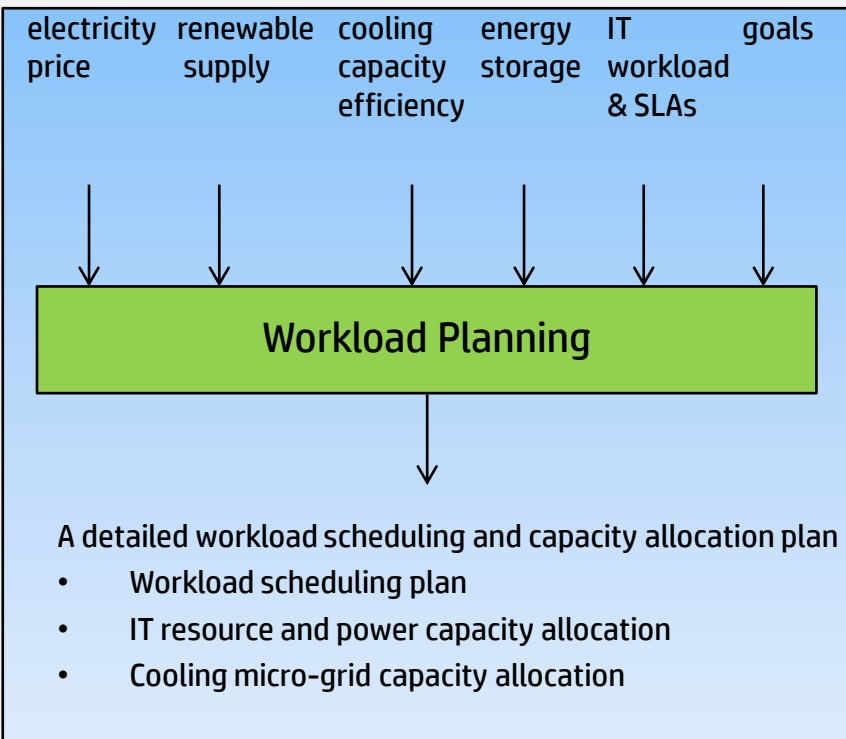
Actual versus predicted

(AAAI 2012)

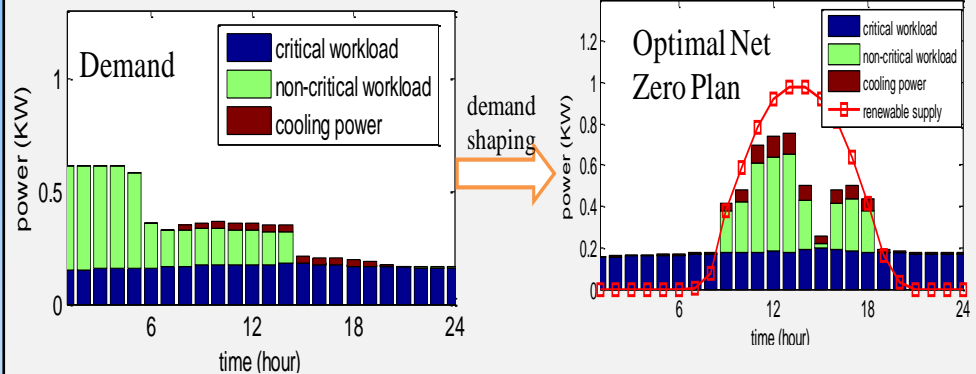


Planning: Supply-Side Aware IT Workload Planning

Planning Flow



Demand Shaping

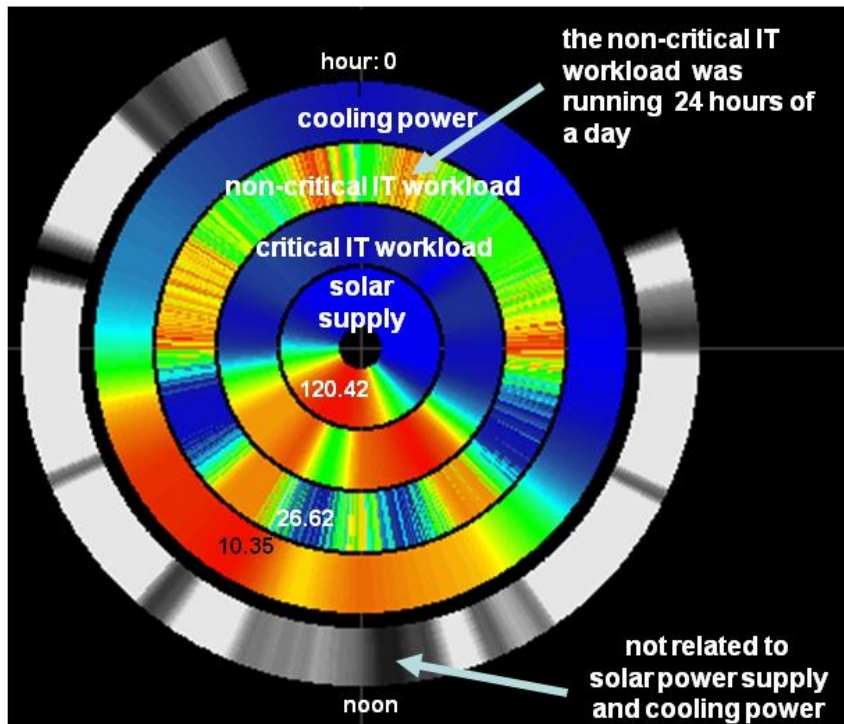


Overall demand is “shaped” according to input constraints and operation objectives

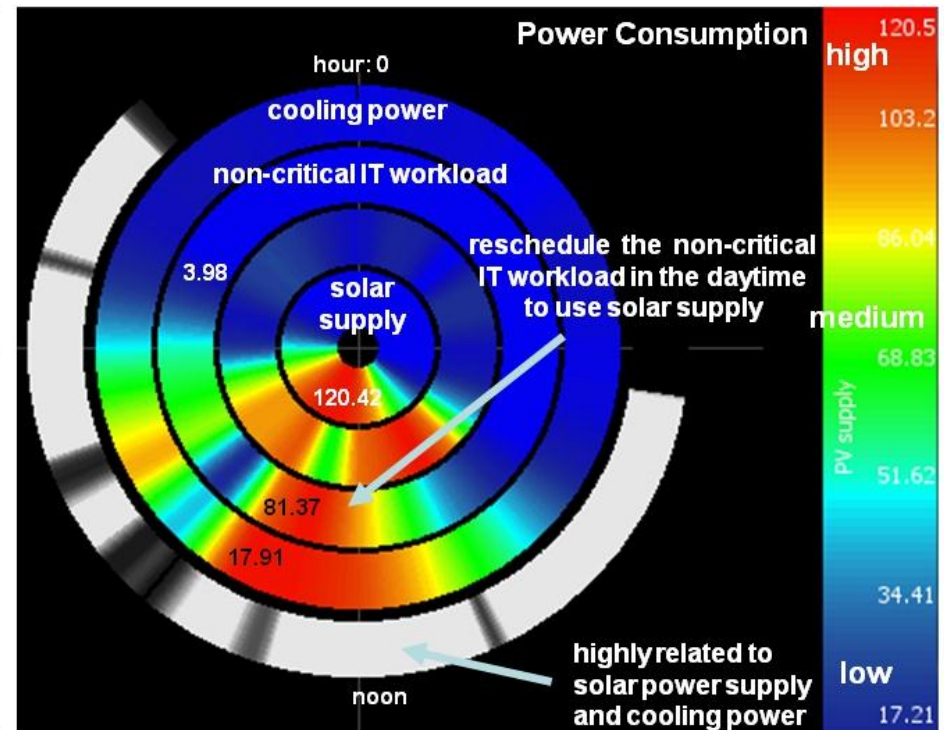
Satisfy critical workload resource requirements

Ref: Z. Liu, Y. Chen, C. Bash, A. Wierman, D. Gmach, Z. Wang, M. Marwah, C. Hyser, "Renewable and Cooling Aware Workload Management for Sustainable Data Centers", ACM SIGMETRICS/Performance, June 11-15 2012, London, UK.

Power and Workload Visualization



Before optimization



After optimization

Some other projects

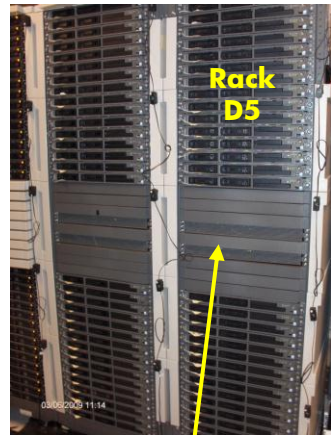
- Anomaly detection (SensorKDD 2010)
- Energy Disaggregation (SDM 2011, AAAI 2013)
- Automating Life Cycle Assessment (IEEE Computer 2011)
- Building Energy Management (BuildSys 2011)



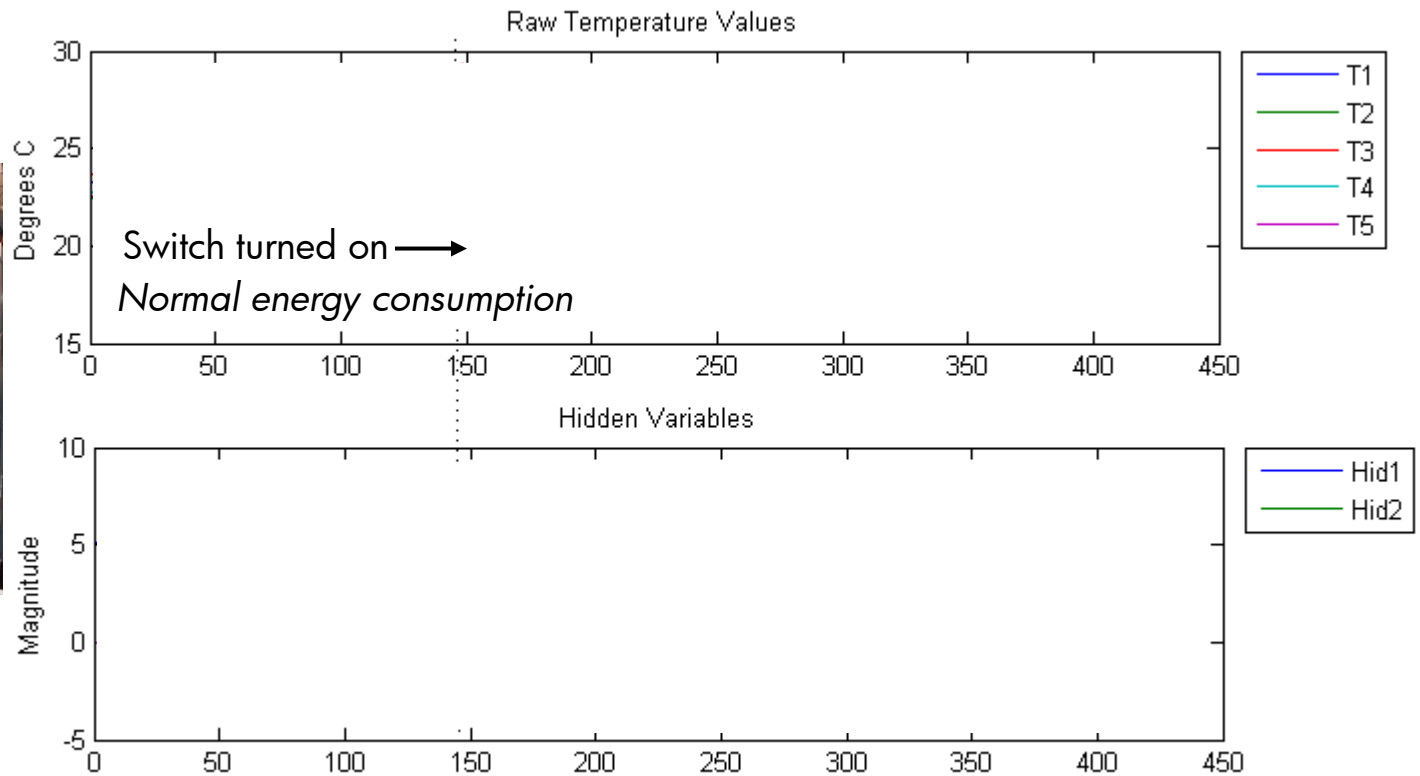
Anomalous Thermal Behavior Detection using PCA

– Example: Event (Anomaly) Detection

Period of increased energy consumption (17 % increase)



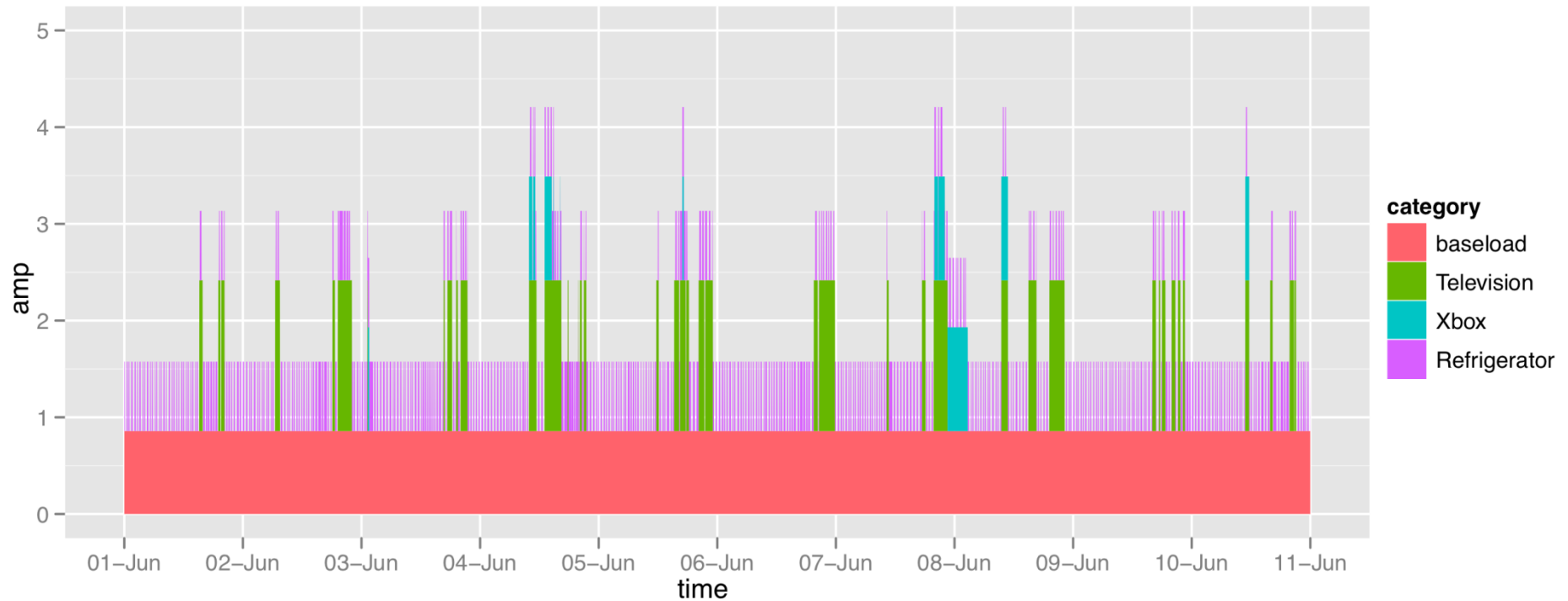
Network Switch



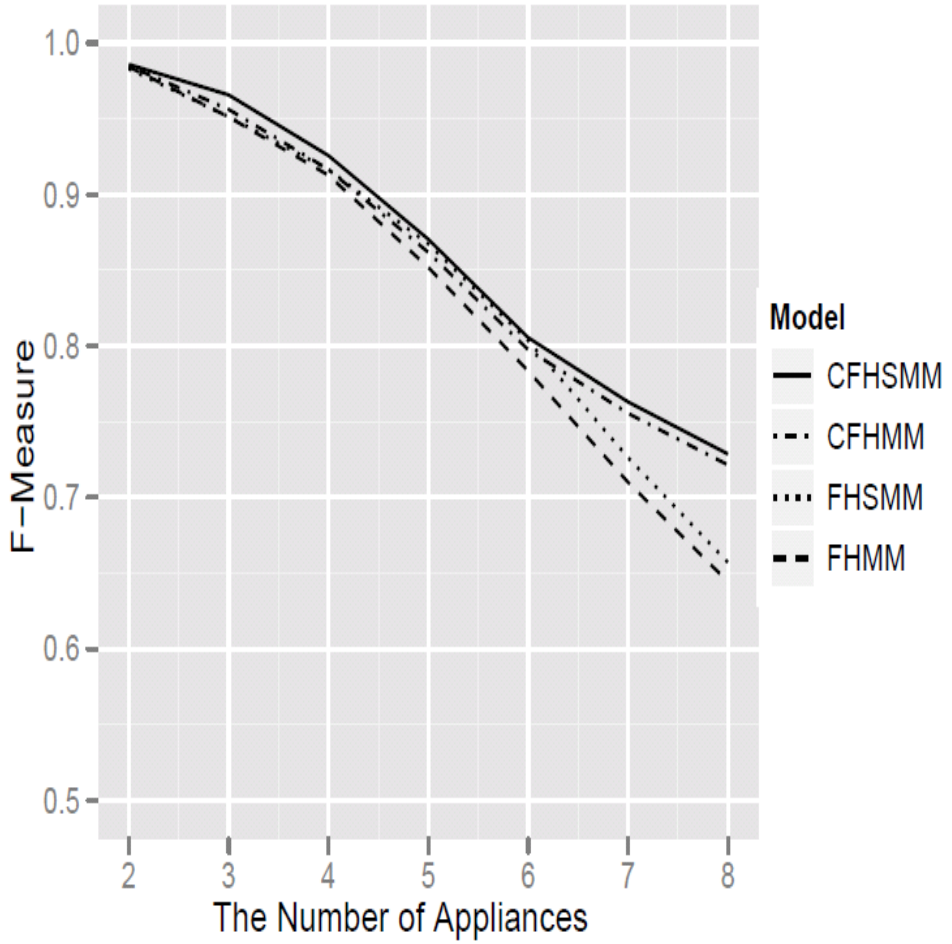
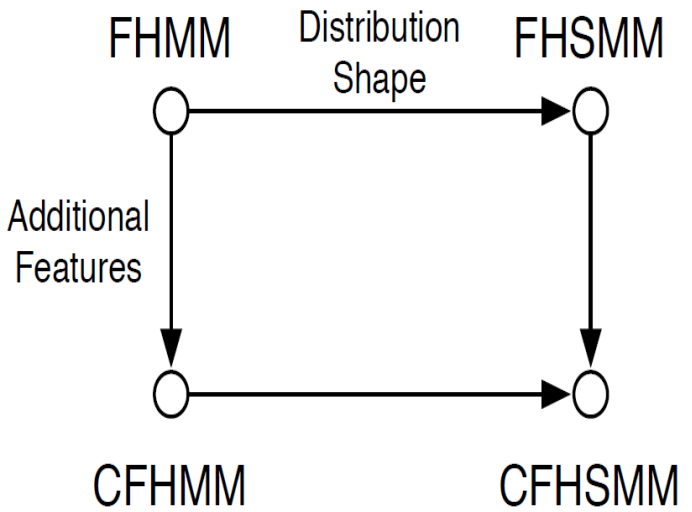
Start: 2009-09-28 16:44:34

End: 2009-09-28 23:58:34

Energy Disaggregation



Proposed Variant of Factorial HMM's (SDM 2011)



References

- P. Chakraborty, M. Marwah, M. Arlitt, and N. Ramakrishnan. Fine-grained Photovoltaic Output Prediction using a Bayesian Ensemble, in *Proceedings of the 26th Conference on Artificial Intelligence (AAAI'12)*, Toronto, Canada, 7 pages, July 2012, To appear.
- Z. Liu, Y. Chen, C. Bash, A. Wierman, D. Gmach, Z. Wang, M. Marwah, C. Hyser, "Renewable and Cooling Aware Workload Management for Sustainable Data Centers", ACM SIGMETRICS/Performance, June 11-15 2012, London, UK.
- Manish Marwah, Amip Shah, Cullen Bash, Chandrakant Patel, Naren Ramakrishnan, "Using Data Mining to Help Design Sustainable Products," IEEE Computer, August 2011
- Hyungsul Kim, Manish Marwah, Martin Arlitt, Geoff Lyon and Jiawei Han, "Unsupervised Disaggregation of Low Frequency Power Measurements", SIAM International Conference on Data Mining (SDM 11), Mesa, Arizona, April 28-30, 2011.
- Gowtham Bellala, Manish Marwah, Martin Arlitt, Geoff Lyon, Cullen Bash, "Towards an understanding of campus-scale power consumption." In ACM BuildSys, November 1, 2011, Seattle, WA.
- Manish Marwah, Ratnesh Sharma, Wilfredo Lugo, Lola Bautista, "Anomalous Thermal Behavior Detection in Data Centers using Hierarchical PCA," in SensorKDD in conjunction with KDD 2010.
- D. Patnaik, M. Marwah, Sharma, Ramakrishna, "Sustainable Operation and Management of Data Center Chillers using Temporal Data Mining," In ACM KDD, June 27 - July 1, 2009, Paris, France.
- Amip Shah, Tom Christian, Chandrakant D. Patel, Cullen Bash, Ratnesh K. Sharma: Assessing ICT's Environmental Impact. IEEE Computer 42(7): 91-93, July 2009.

