

# ADAPTIVE MODEL-DRIVEN FACILITY-WIDE MANAGEMENT OF ENERGY EFFICIENCY AND RELIABILITY

---

Ananta Tiwari

EP Analytics, Inc.

Performance Modeling and Characterization (PMaC) Lab, SDSC

Co-authors:

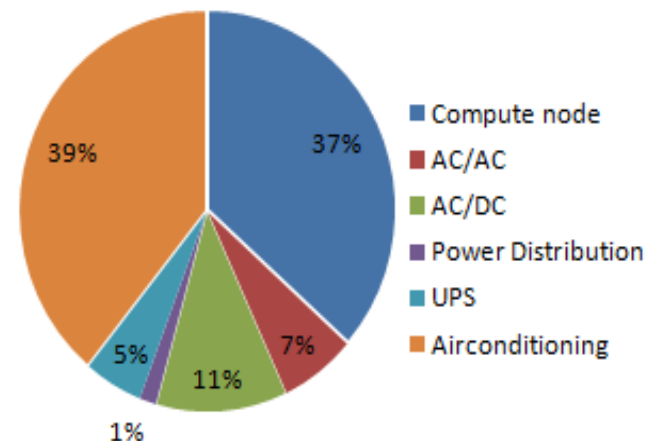
Michael Laurenzano, Adam Jundt (EP Analytics)

William A. Ward, Jr., Roy Campbell (DoD HPCMP)

Laura Carrington (EP Analytics, SDSC)

# Motivation

- Utility costs, constraints on power delivery, and system reliability issues are limiting the expansion of HPC systems
  - Electricity usage in datacenters increased from 7 GW to 10 GW between 2005 and 2010 and projected to reach 20 GW by 2015.
- Power consumption in datacenters comes not only from servers, but also:
  - Air conditioning
  - Power conversion
  - Other sources



# Motivation

- Nearly equaling power of servers is power required by the Computer Room Air Conditioning (CRAC) units
  - Previous work mostly just focus on node-level optimizations
- Current practice is to set the ambient temperature between 64F and 80F.
  - Some centers have suggested raising the set-point temperature higher than 80F to save additional energy.
- However, by raising the set-point temperature, server fan power and chip leakage power will increase
  - Chip leakage loss is projected to increase with shrinking feature sizes

# Motivation

- Set-point temperature is usually regulated to also reduce hardware error rates
- Studies have shown that HDDs and DIMMs are the 2 components that most commonly break down\*
  - Failure rate has been shown to depend on set-point and fluctuations in set-point

Multiple competing and offsetting events have to be considered simultaneously to manage the thermal, power and reliability aspects of large-scale computing installations.

\* N. El-Sayed, I. Stefanovici, G. Amvrosiadis, A. A. Hwang, and B. Schroeder. Temperature management in datacenters: Cranking up the thermostat without feeling the heat. *USENIX*, 38(1), February 2013.

# E2MP

- Our Solution: We propose to develop an Energy Efficient Management Platform (E2MP), a power-aware, green computing technology that can enact performance-neutral power and reliability management policies on high performance computing (HPC) centers
  - Policies can be both compute-node-level (e.g., selecting the right CPU speed) as well as facility-level (e.g., selecting the right set-point)
- We envision E2MP to have a holistic view of a datacenter's energy usage, from an encompassing room view down to application level models, that will enable finding energy and reliability optimal set-points at each level of operation

# Our position

*The operational settings that increase energy efficiency and reliability of a given HPC facility are not only governed by the capabilities and specifications of the computing and cooling hardware, but are also complex functions of the application software running in that facility.*

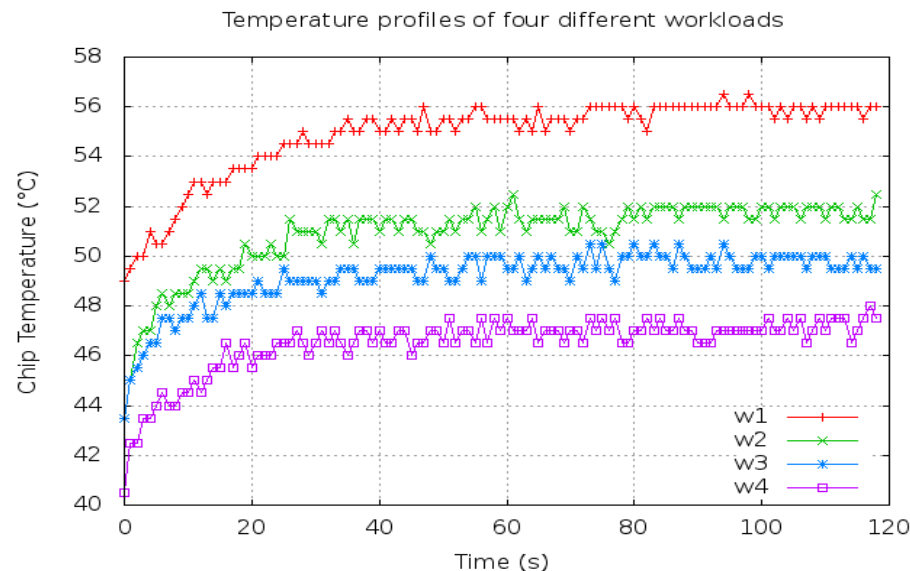
# Determining optimal set-point

- In order to determine energy optimal set-point, we need comprehensive understanding of:
  - Thermal footprint of applications (i.e., applications' effect on chip temperature)
    - Foot-print depends on computational characteristics and rotational speed of node fans
  - How ambient (or set-point) temperature affects the thermal footprint
  - How leakage power increases with the rise in chip temperature
  - Rotational speed of nodes' cooling fans

**E2MP provides a set of models that can inform the decision to select energy optimal set-point**

# Thermal footprint of applications

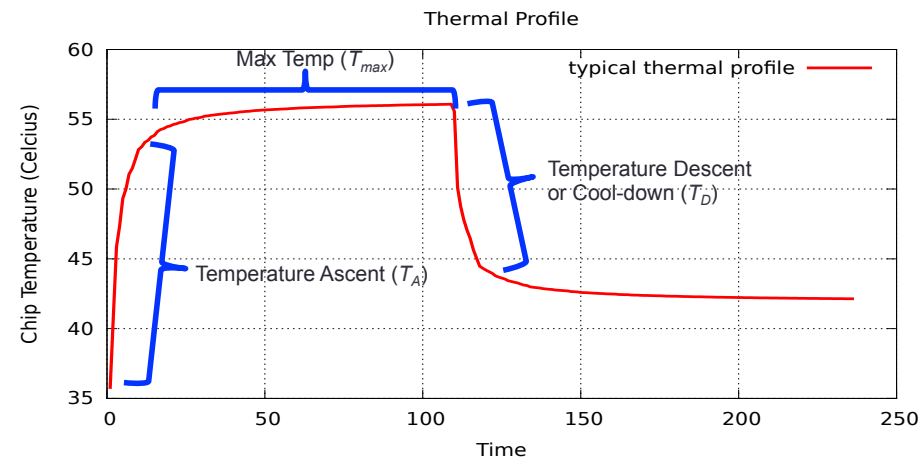
- Computational characteristics have a significant impact on the thermal state of a given system
- The figure shows the temperature for 4 different computations, the temperature difference between the “hottest” and “coldest” computation is 10C





# Typical thermal profile

- Figure shows a typical thermal profile of an application along with the cool-down phase
  - $T_{max}$  depends on computational characteristics of the application, instantaneous power draw and ambient temperature
    - Need models that can predict instantaneous power draw given some computational characteristics
  - $T_A$  and  $T_D$  can be modeled using non-linear regression methods



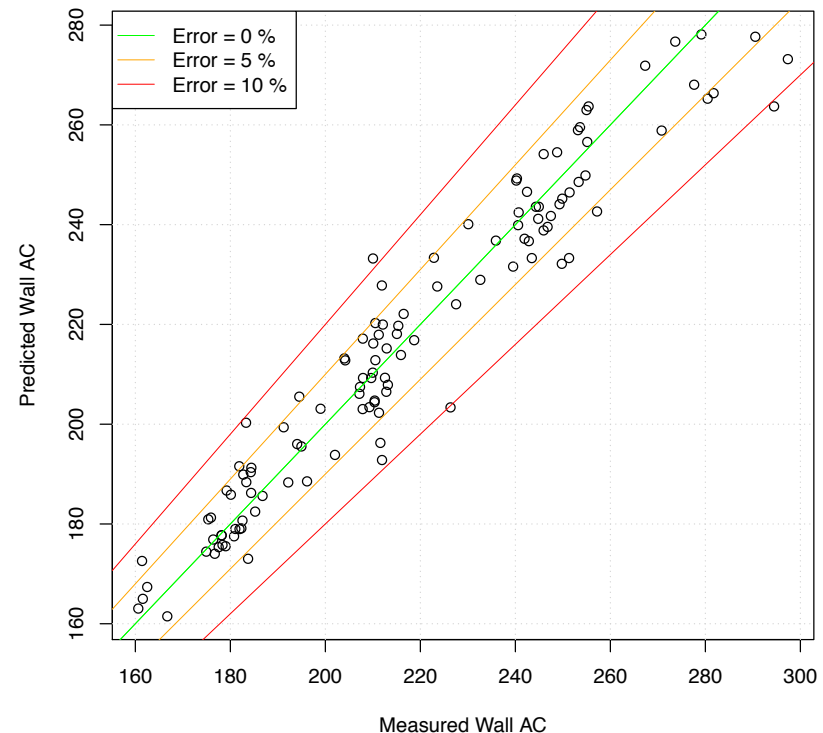
*Accurate models for  $T_{max}$ ,  $T_A$  and  $T_D$  will provide full knowledge of the thermal behavior of a given application, which can then be used to determine the optimal set-point.*

# Modeling $T_{max}$

- $T_{max}$  can be expressed as:  
$$T_{max} = f(\text{computational characteristics, power draw, set-point})$$
- Computational characteristics are extracted using our binary analysis tools developed on top of PEBIL
  - Static analysis tools provide operation counts (memory and floating point), operation parallelism and program structure (e.g., function and loop boundaries)
  - Dynamic analysis tools provide cache hit rates, execution counts, loop length, etc.
- Power draw can either be measured or modeled
  - We take the modeling route to derive power draw at fine granularities

# Modeling Power Draw

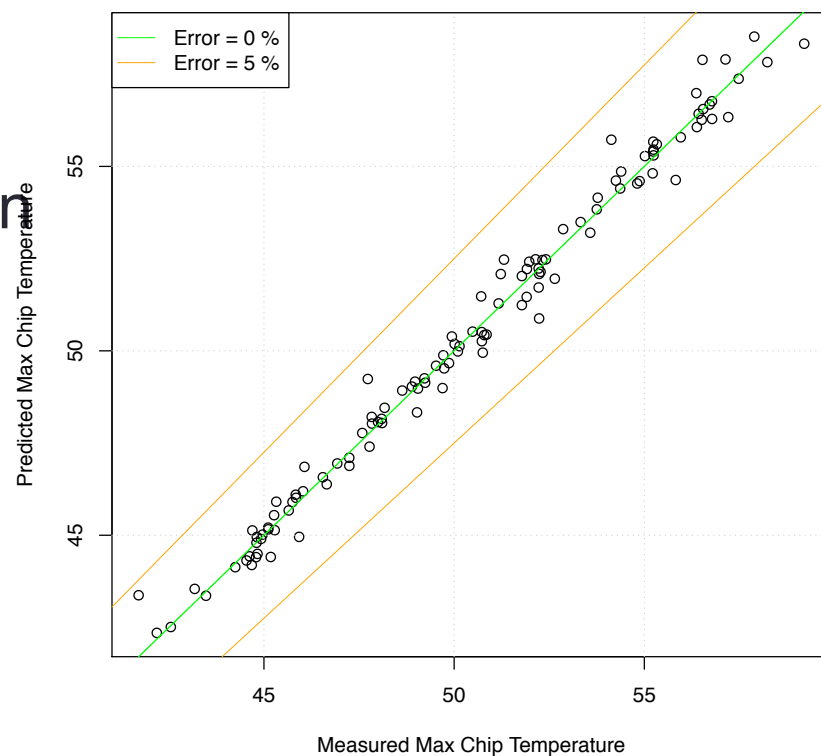
- Power draw also depends on computational characteristics and CPU speed (clock frequency)
- The graph shows the measured vs. modeled wall AC
  - Our model is highly accurate (absolute mean error: 3%)



# Modeling $T_{max}$

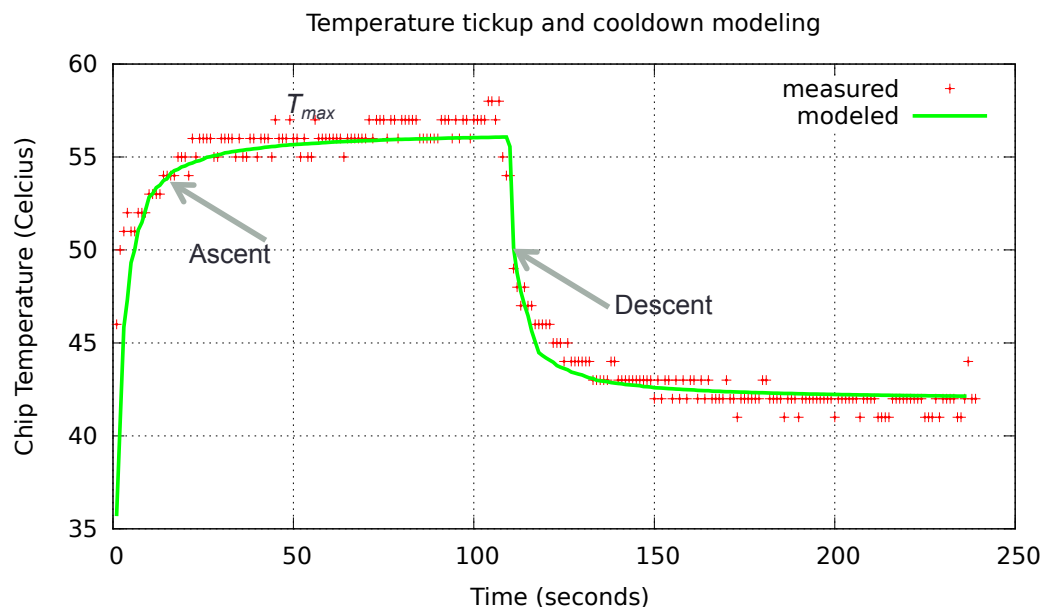
$T_{max} = f(\text{computational characteristics, modeled power draw, set-point})$

- The graph shows the model accuracy for a set of HPC kernels
- Our model predicts  $T_{max}$  with mean absolute error in prediction  $< 2\%$ .



# Modeling $T_A$ and $T_D$

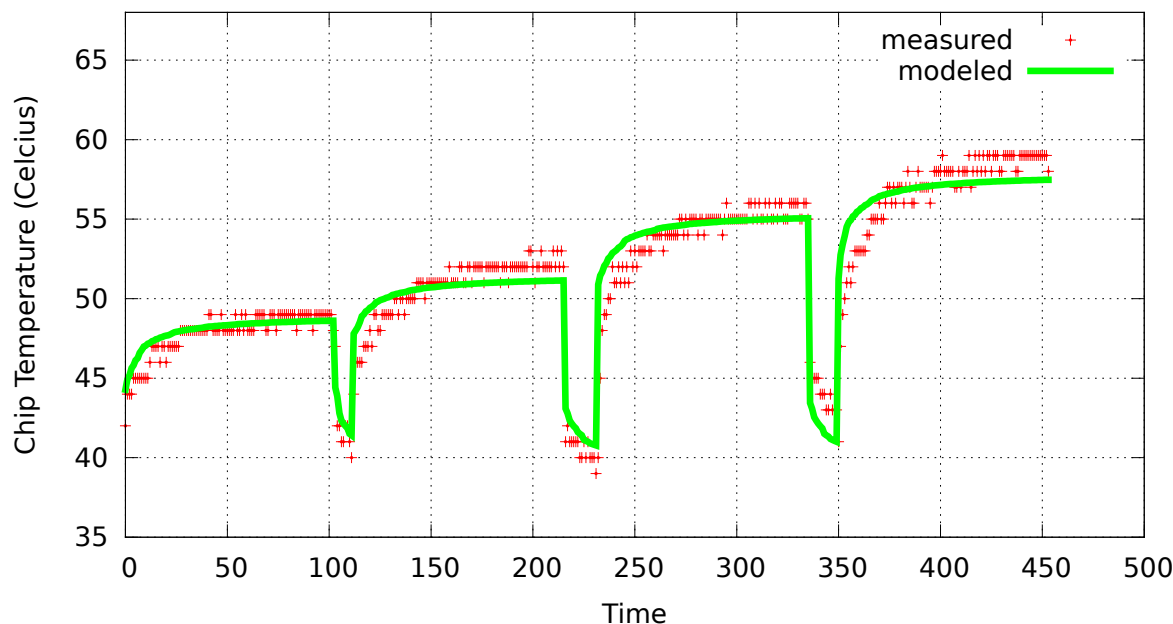
- $T_A$  and  $T_D$  are modeled using non-linear regression technique
- The figure shows the capability of our models to predict the temperature ascent and descent



# Putting everything together

- We designed a synthetic workload that consists of four different loops with increasing thermal footprint.
  - Sleep intervals are inserted in between the loops to show the cool-down or descent modeling

Modeling Thermal Profile of Multiple Applications

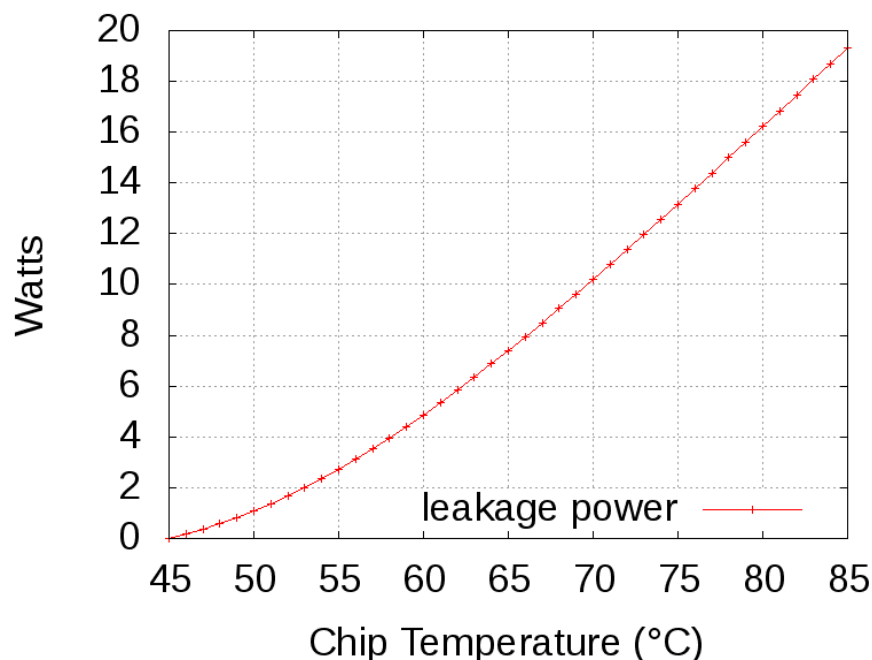


# Determining optimal set-point

- In order to determine energy optimal set-point, we need comprehensive understanding of:
  - Thermal footprint of applications (i.e., applications' effect on chip temperature)
    - Foot-print depends on computational characteristics and rotational speed of node fans
  - How ambient (or set-point) temperature affects the thermal footprint
  - How leakage power increases with the rise in chip temperature
  - Rotational speed of nodes' cooling fans

# Chip leakage power characterization

- As chip temperature increases, leakage power increases exponentially
- We have developed models for SandyBridge processors that predict the static leakage power as chip temperature increases, caused by ambient temperature or computational properties
- We can combine the thermal application models and the chip leakage models to see how much leakage power an application will incur





# Fan power

- Another area for energy savings in the datacenter is from fan power
- On a typical server, 14% of power is dedicated to running fans (see figure)
- Fans are typically run at their highest speed to ensure processors stay below their thermal limits
  - Informed selection of fan speed can reduce fan power significantly

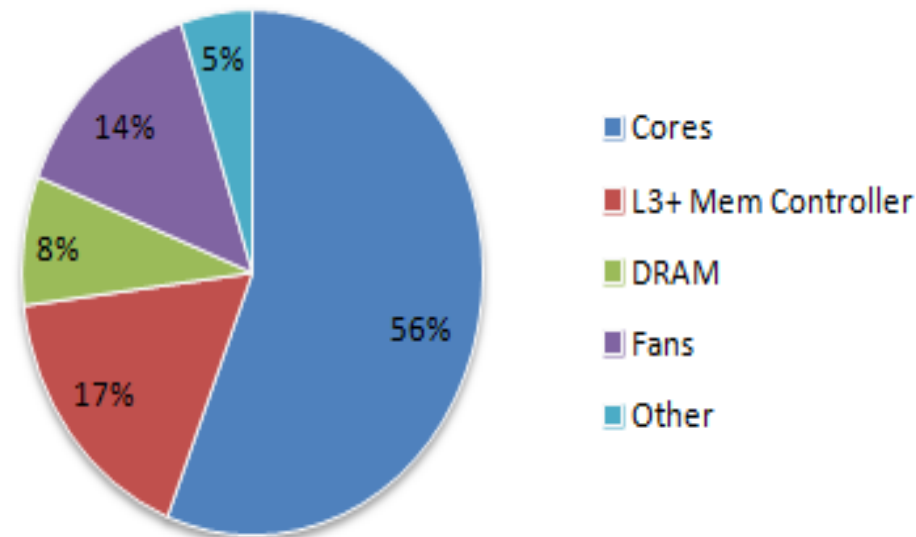
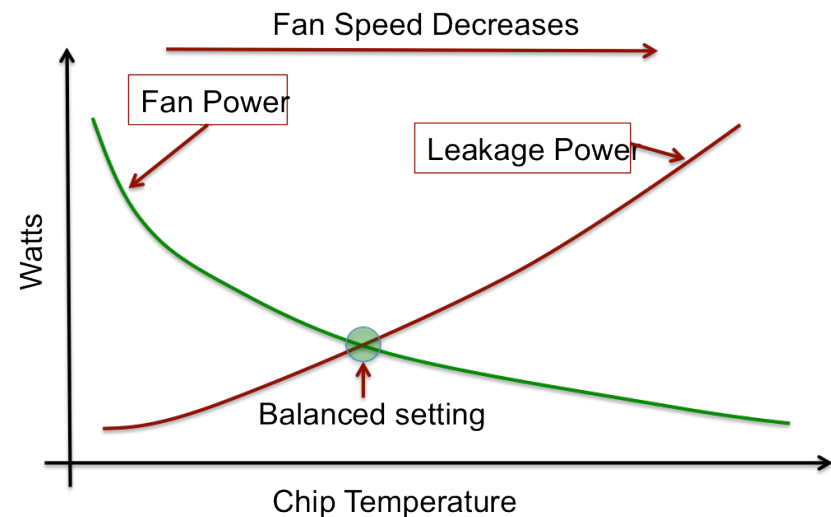


Figure : Typical power breakdown for a datacenter server

# Use-cases of our models

- Using the  $T_{max}$ ,  $T_A$ ,  $T_D$  and leakage power models, datacenter operators can select the optimal server fan speeds
  - Tradeoff scenarios like the one depicted in the figure can also be explored to find the right balance in fan power and leakage power
- Since our models can inform how “hot” an application will make a node apriori, datacenter operators can take steps to reduce energy usage using a number of methods:
  - Adjusting fan power
  - Moving “hot” applications to the nodes that are closer to CRAC vents
  - Lowering CPU frequency



# Summary

- Modeling an entire datacenter is a difficult problem
- By breaking down each component's power usage in the center, from the CRAC to the applications run on nodes, E2MP is able to allow for energy savings without sacrificing performance.

# Committee Questions

1. *What is the major contribution of your research?*
  - We provide model-based integrated power, thermal, performance and reliability view that will enable HPC centers to develop a set of best practices to increase energy efficiency by enacting performance and reliability neutral strategies.
2. *What are the gaps you identify in the research coverage in your area?*
  - Combined view of inter-related power, reliability, thermal and energy efficiency related issues.
3. *What major opportunities do you see for cross-pollination between your project and others?*
  - Per-device failure probability models based on various metrics – temperature, age of device
4. *What would you like to most see solved/addressed other than what you are working on?*
  - Resiliency characterization and modeling that correlates faults and failure probability with system and environmental events
    - Hardware event counters -- performance counters, power state, thermal state, errors corrected
    - Logs and state descriptors (e.g. age of devices, usage, disabled units, thermal alarms triggered, prior fault history, etc.)
    - External variables (room temperature, system power draw, etc.).

# QUESTIONS / THANK YOU

---

Contact:

[ananta.tiwari@epanalytics.com](mailto:ananta.tiwari@epanalytics.com)

[www.epanalytics.com](http://www.epanalytics.com)