

# High-Performance AI: A View from Systems and Frameworks

Deep Learning and  
Supercomputer Workshop  
of SC18

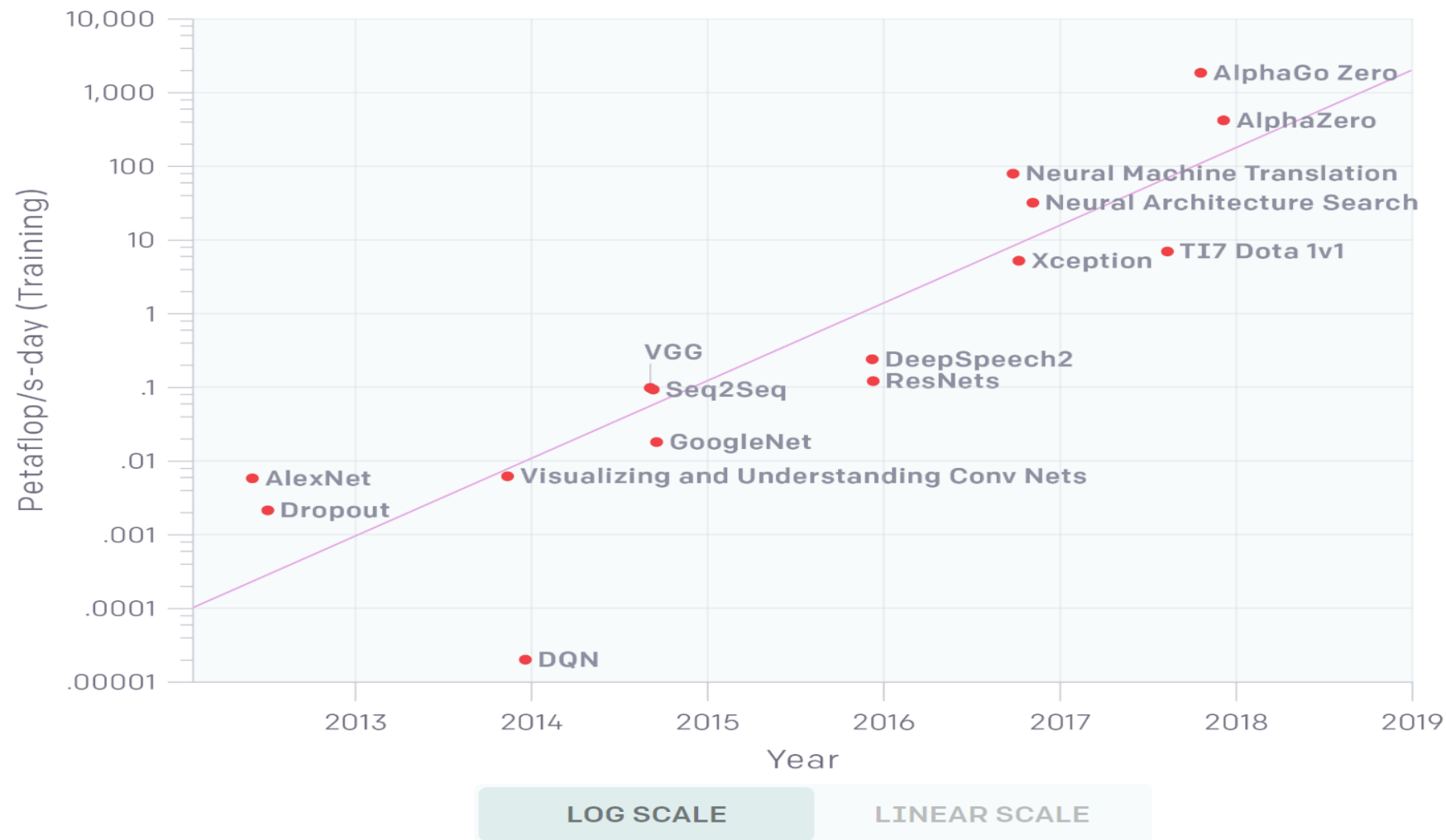
Digital Science Center

Judy Qiu, Langshi Cheng , Bo Peng , Chathura Widanage, Sahil Tyagi  
Intelligent Systems Engineering Department Indiana University  
Email: [xqiu@indiana.edu](mailto:xqiu@indiana.edu)

- **Challenges and Opportunities for AI and Systems**
- **Recent work**

**Outline**

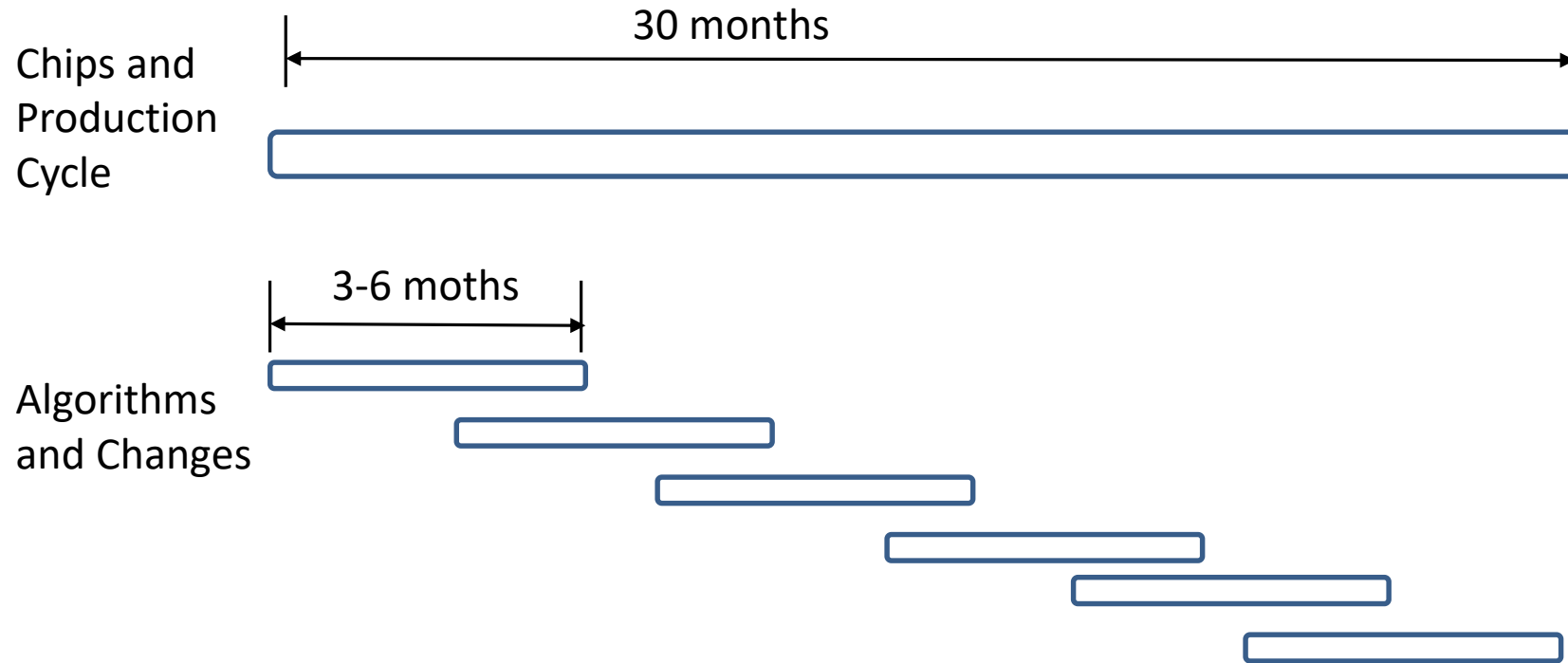
## AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



A petaflop/s-day (pfs-day) consists of performing  $10^{15}$  neural net operations per second for one day, or a total of about  $10^{20}$  operations (from OpenAI.com).

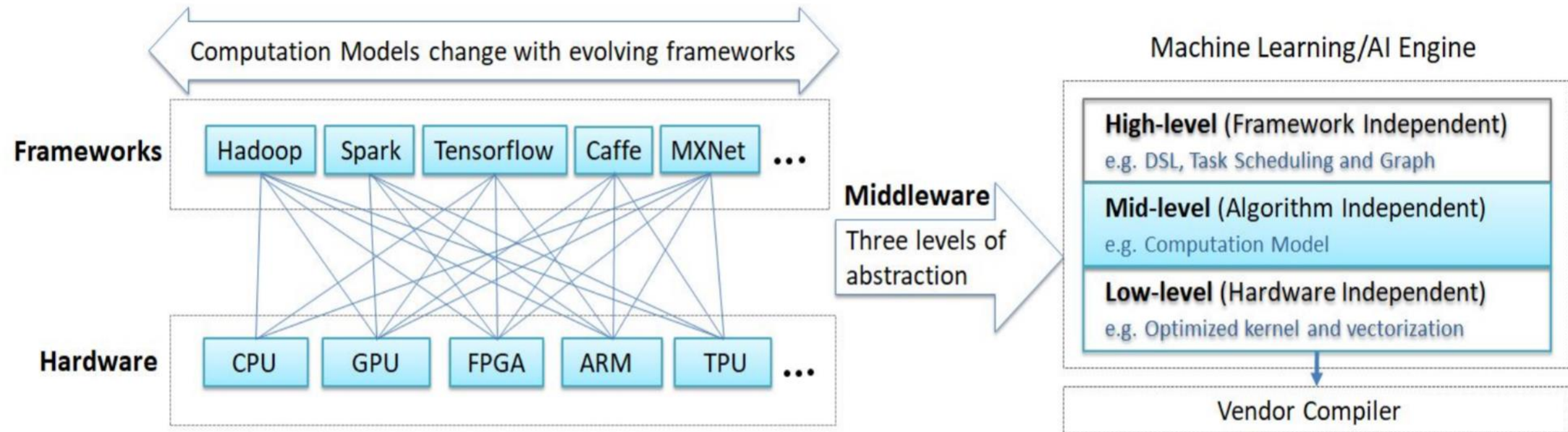
# As AI progresses, systems are far outside today's capabilities

- New Algorithms vs. Processors being Developed at Different Rate



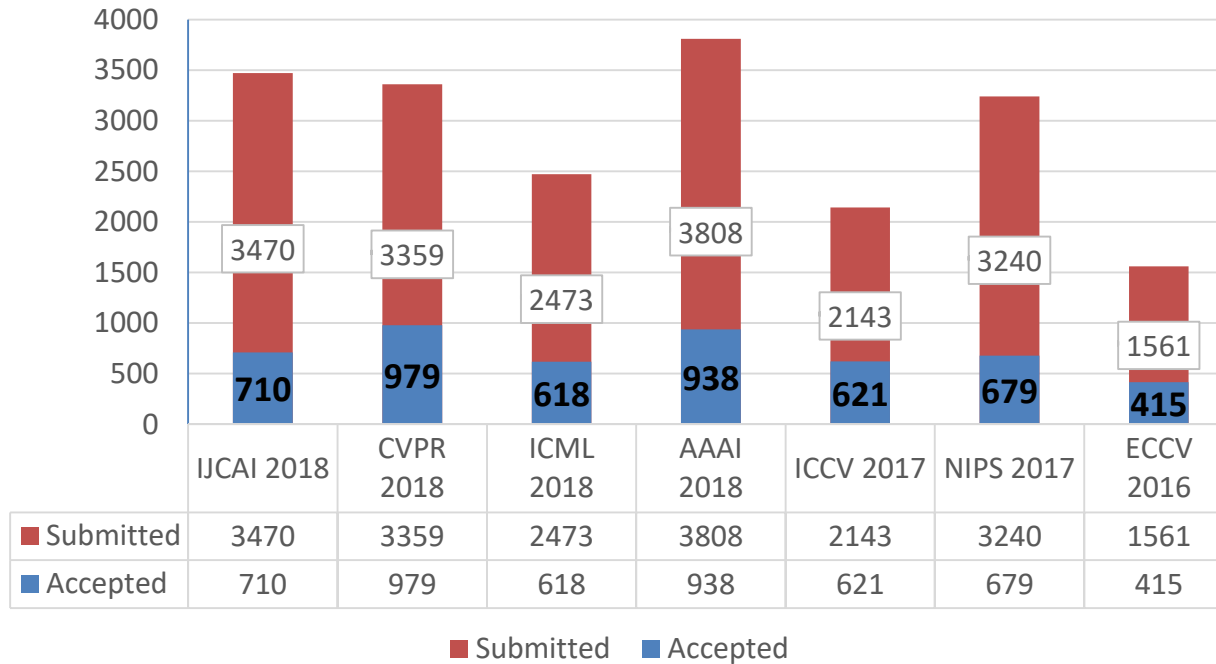
Depei Q. World AI Conference, 2018

# Challenges of Design in Computer Architecture

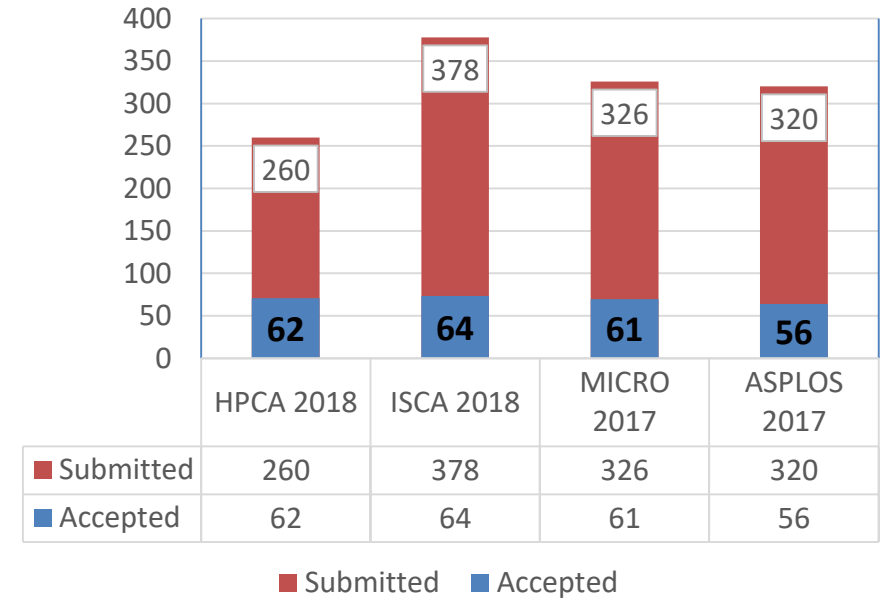


# AI Frameworks and Architecture of AI/Machine Learning

## AI Conference and Publications



## Systems Conference and Publications



Number of Publications 4960 vs 243

Depei Q. World AI Conference, 2018

# Gap between AI and Systems Research

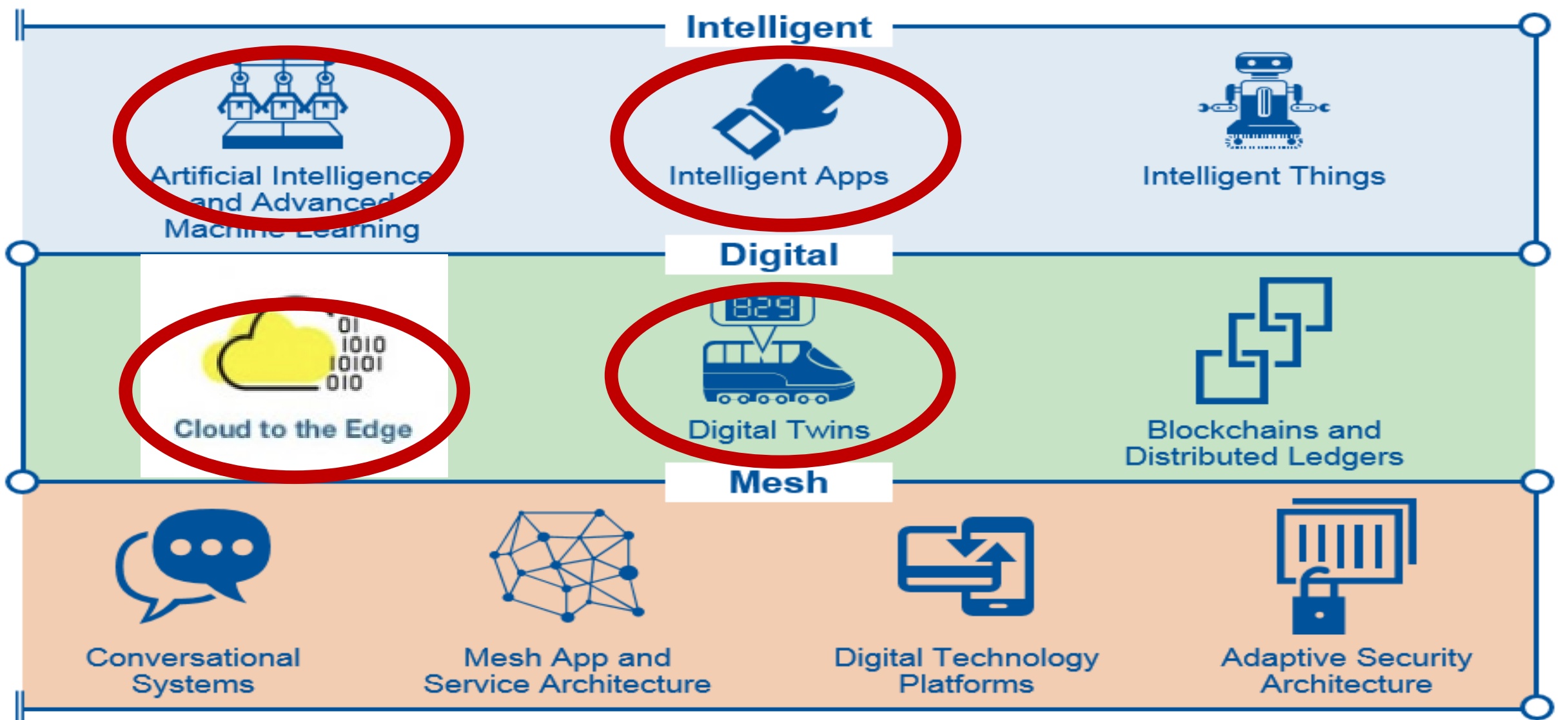
- AI: Challenges and Opportunities for AI and Systems
- **Our Recent Work**

**Outline**

# Real-Time Anomaly Detection from Edge to HPC-Cloud

Intro & Preliminary  
Results





© 2017 Gartner, Inc.

Gartner's report on Strategic Technology Trend for 2017-2018

- The [IndyCar Series](#) is a major open-wheel racing format in North America. The series' premier event is the [Indianapolis 500](#), held each May.
- Computing Systems and Data analytics is critical to the sport, both in improving the performance of the team to make it faster and in helping the race control to make it safer.

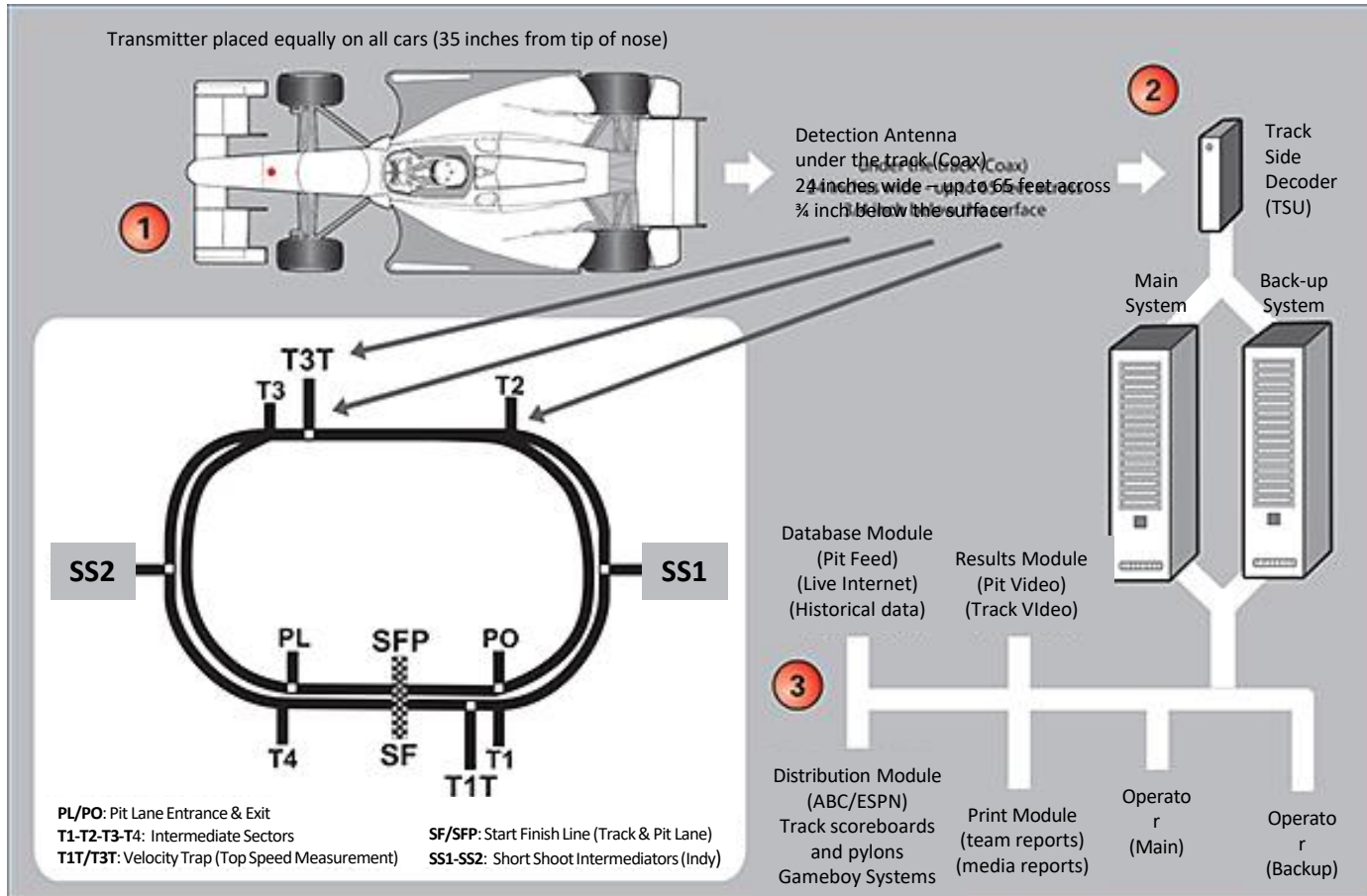


# IndyCar





Indianapolis 500 Car Racing, May 2017



- Sensors in the cars and under the track.
- Antenna and communication system.
- Telemetry data (including the many performance information of the cars like speed, gear, brake, throttle, etc) stream into the on-site computer system in a real-time fashion.

# Timing and Score Data



Command	Count	Protocol	Description	Frequency
A	2052	MLP	Announcement	Every 60 seconds
C	19432	MLP	Completed Lap Results	Upon Event (new and repeated)
D	2652	RP	Invalidated Lap Information	Every 30 seconds
E	7737	MLP	Entry Information	Every 60 seconds
F	725	MLP	Flag Information	Upon Event (new and repeated)
G	7892	RP	Car Display Pit Stop Timer Information	Every 120 seconds
H	17260	MLP	Heart beat	Every Second
I	53	MLP	Invalidated Lap Information	Upon Event (new and repeated)
L	79884	MLP	Line Crossing Information	Upon Event
M	1738	eRP	Messages	Upon Event
N	3861	MLP	New Leader Information	Upon Event
O	33263	RP	Overall Results	Upon Event
P	3693653	eRP	Telemetry Data	
R	701	MLP	Run Information	Every 20 seconds
S	102272	MLP	Completed Section Results	Upon Event (new and repeated)
T	233	MLP	Track Information	Every 60 seconds
U	235	RP	Track Information	Every 30 seconds
V	117	MLP	Version Stream Information	Every 120 seconds
W	287	RP	Weather Data	Every 60 seconds
X	12124	RP	Heart beat	Every Second

- The INDYCAR Timing system supports retrieving timing data from the primary timing system – serial or sequential data feed for live data and report querying for historical or archived data.
- The **Results Protocol** is designed to deliver more detailed results information through the use of a single record command.
- Example: one Indianapolis 500 car race on the 28th of May 2017 which contained 750 MB of data and total of 3986170 records.

# Dataset

\$P – Telemetry (does not include full Record Header and not in the XML file)

Fieldname	Data description	Comments
No	Characters	Car number – 4 characters max
Time Of Day	Integer	TOD in ms
Lap Distance	Float	Metres since start of lap (12345.67)
Vehicle Speed	Float	MPH ie. 123.456
Engine Speed	Integer	RPM ie. 12345
Gear	Integer	0 = Neutral, 1..6 = Gear 1 through 6
Brake	Float	% brake
Throttle	Float	% throttle
Steering	Float	-1.00 .. 0.00 .. 1.00
Long. Accel.	Float	G's – might not be valid for all cars
Lat. Accel.	Float	G's – might not be valid for all cars
Vert. Accel.	Float	G's – might not be valid for all cars
Boost Pressure	Float	Turbo boost pressure
Tire Type	Character	P = Primary A = Alternate W = Wets U = Unknown
OT_Event	Integer	Seconds remaining in current instance, may be zero if not valid or available
OT_Remain	Integer	Number of seconds or pushes remaining, depending on current rules.
OT_Status	Boolean	0 = false, not activated 1 = true, P2P/OT active

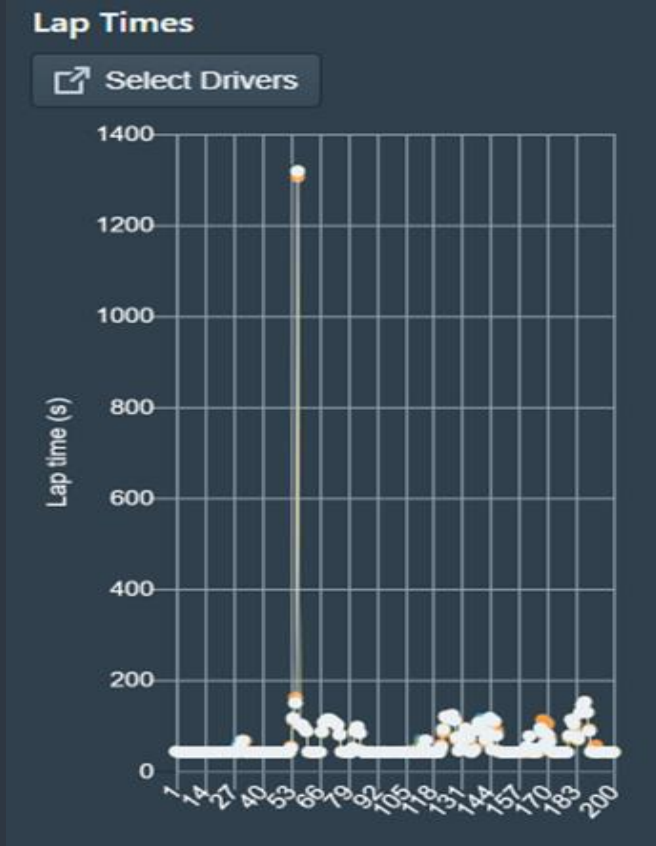
- Telemetry is a radio device that relays information such as engine, tire, steering and throttle performance to team engineers in the pit box. The team can monitor car and driver activity to ensure the car is performing properly.
- During a race, IndyCar Series teams use telemetry to gather data live from their cars as they formulate their race strategy. In the IndyCar Series, teams receive their data from their own on-board systems as well as from the league's Timing & Scoring operation.

# Telemetry Data

# SIMULATION



<https://goo.gl/WRet7Q>



20

# Ed Carpenter

Indianapolis IN  
Ed Carpenter Racing

Chevy

 Veteran


2C0

**#1**

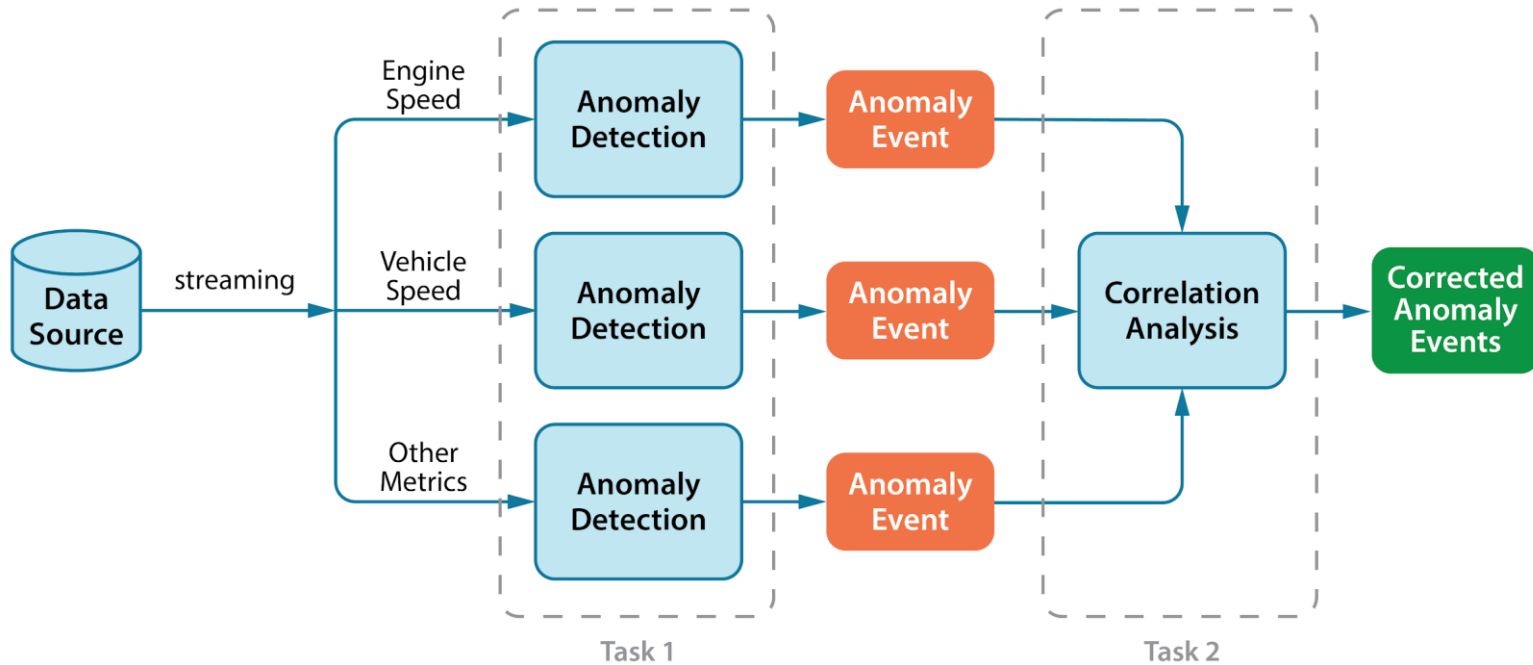
# Real-time data analysis is in need

"We want to know if it's an expected event or a minor deviation that we need to be worried about. It helps race control people. "

" And we want to know the data corresponding to the **anomaly**; when car got into problems what kind of event it is and what is **causing** it."

## Problem





- Anomaly Detection: Learning algorithms themselves can only find the abnormal pattern in the data with best efforts under predefined assumption of what is “normal”.
- Correlation Analysis: Learning algorithms can find out the “events” from data and mine correlation relationship among the events.

# Tasks

- Dataset

- One of the Indianapolis 500 car races on the 28th of May 2017, which consists of 750 MB of data and a total of 3986170 records in the log file.

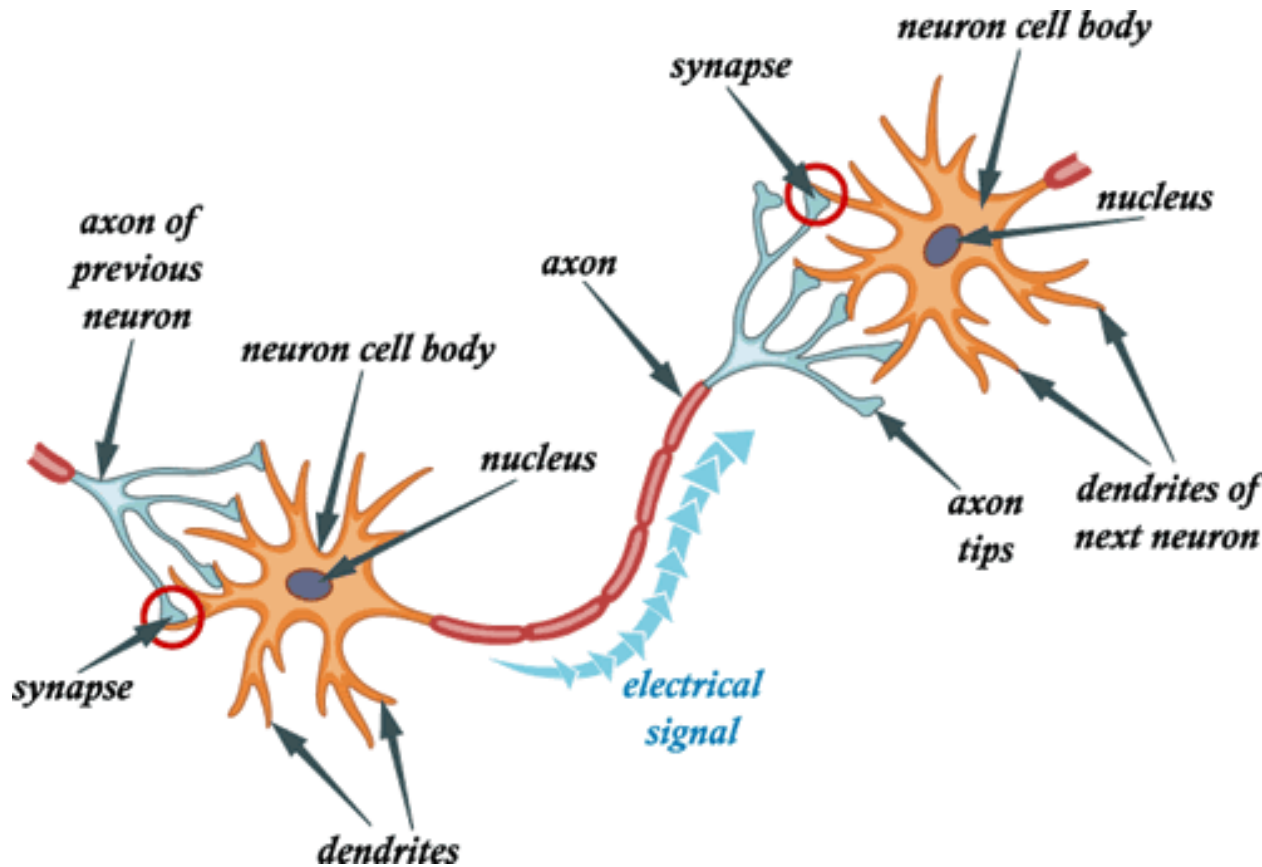
- Approach

- Hierarchical Temporal Memory (HTM) algorithm, a state-of-the-art online algorithm to detect anomaly in real-time data source.

## Experiments and Preliminary Results

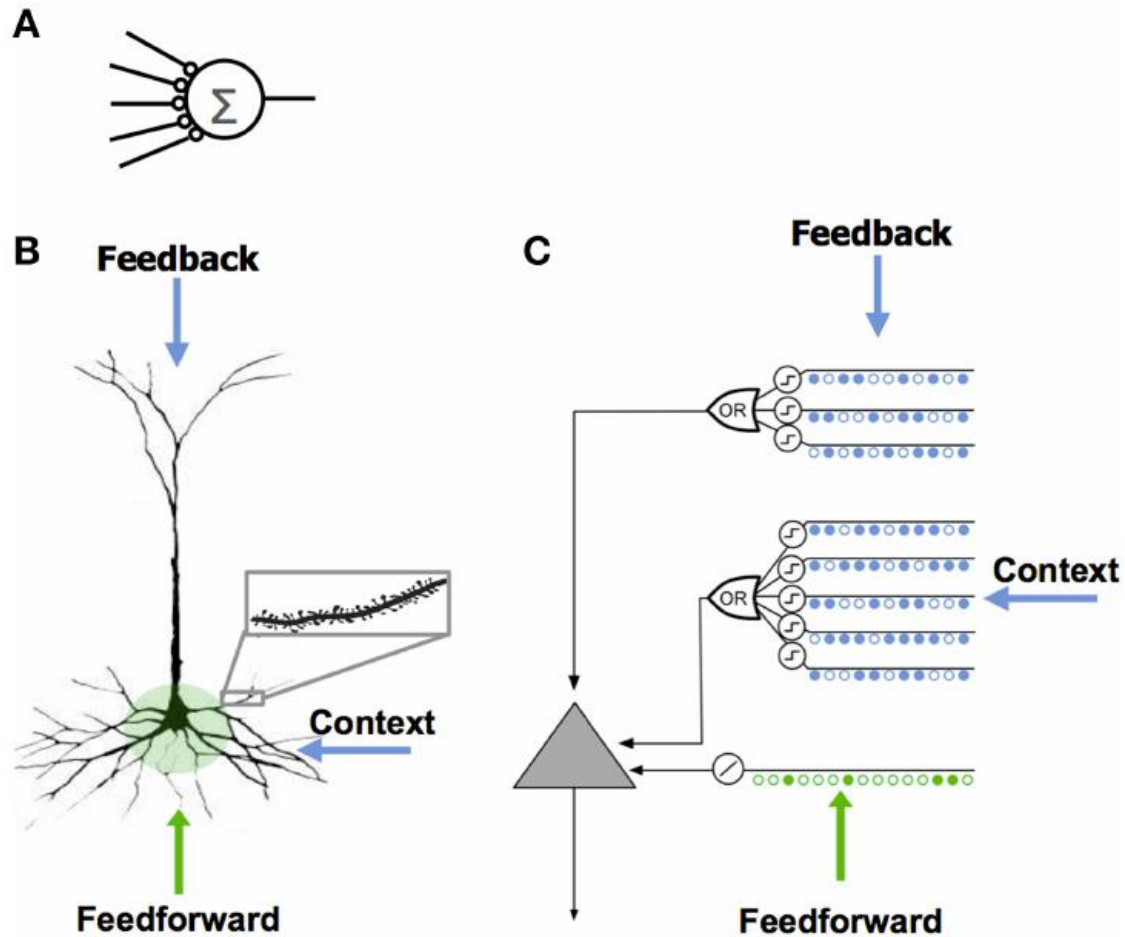
# HTM neuron model

Streaming Data Analysis



- Neurons have specialized projections called **dendrites** and **axons**.
- Dendrites bring information to the cell body and axons take information away from the cell body.
- Information from one neuron flows to another neuron across a **synapse**. The synapse contains a small gap separating neurons.

# Axon, Dendrite and Synapse



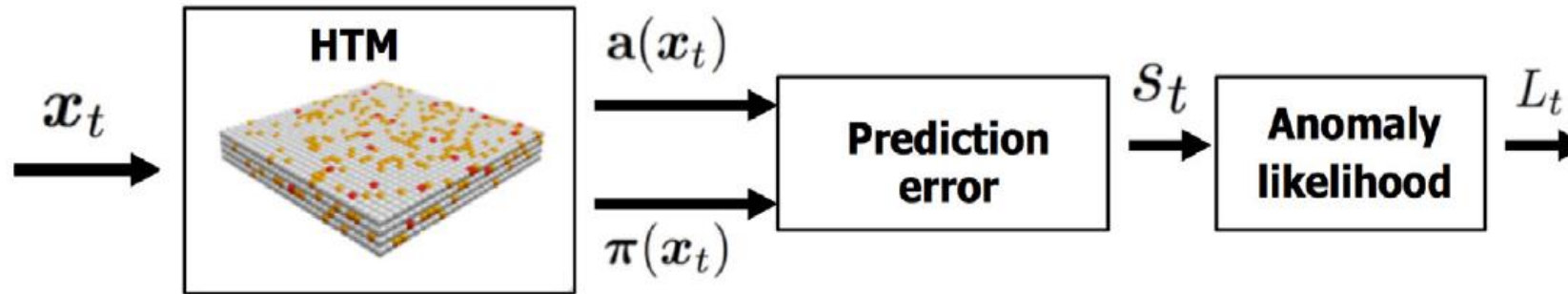
- A. The neuron model used in most artificial neural networks has a few synapses and no dendrites
- B. Different source of input: feedforward, feedback and context
- C. HTM sequence memory models dendrites with an array of coincident detectors each with a set of synapses

# Hierarchical Temporal Memory(HTM)



- Hierarchical Temporal Memory (HTM) is a machine learning technology that aims to capture the structural and algorithmic properties of the neocortex.
- HTM models **neurons**, which are arranged in columns, in layers, in regions, and in a hierarchy. In this regard HTMs are a new form of neural network.
  - HTM consists of a **layer** of HTM neurons organized into a set of **columns**.
  - It models high-order sequences (sequences with long-term dependencies) using a composition of two separate **sparse representations** (spatial and temporal).
  - When receiving the next input, the network uses the difference between predicted input and the actual input to update its **synaptic connections**.
- HTMs are online learning and prediction machines that can be applied to many types of problems.

# HTM



- The input time series  $x_t$  are fed to the HTM component. It models temporal patterns in  $a(x_t)$  and output a prediction in  $\pi(x_t)$ .
- Then by building a statistical model on the prediction error,  $\pi(x_t) - a(x_{t-1})$ , anomaly likelihood score can be calculated on  $x_t$ .

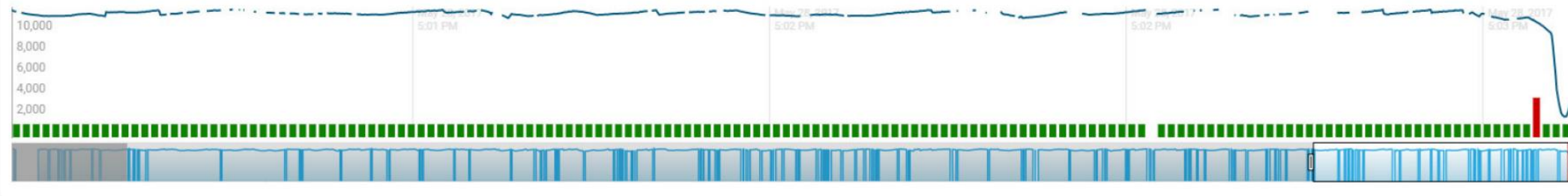
# Anomaly Detection Based on HTM

EngineSpeed(RPM)



(a). Car#9 in middle of the race

EngineSpeed(RPM)

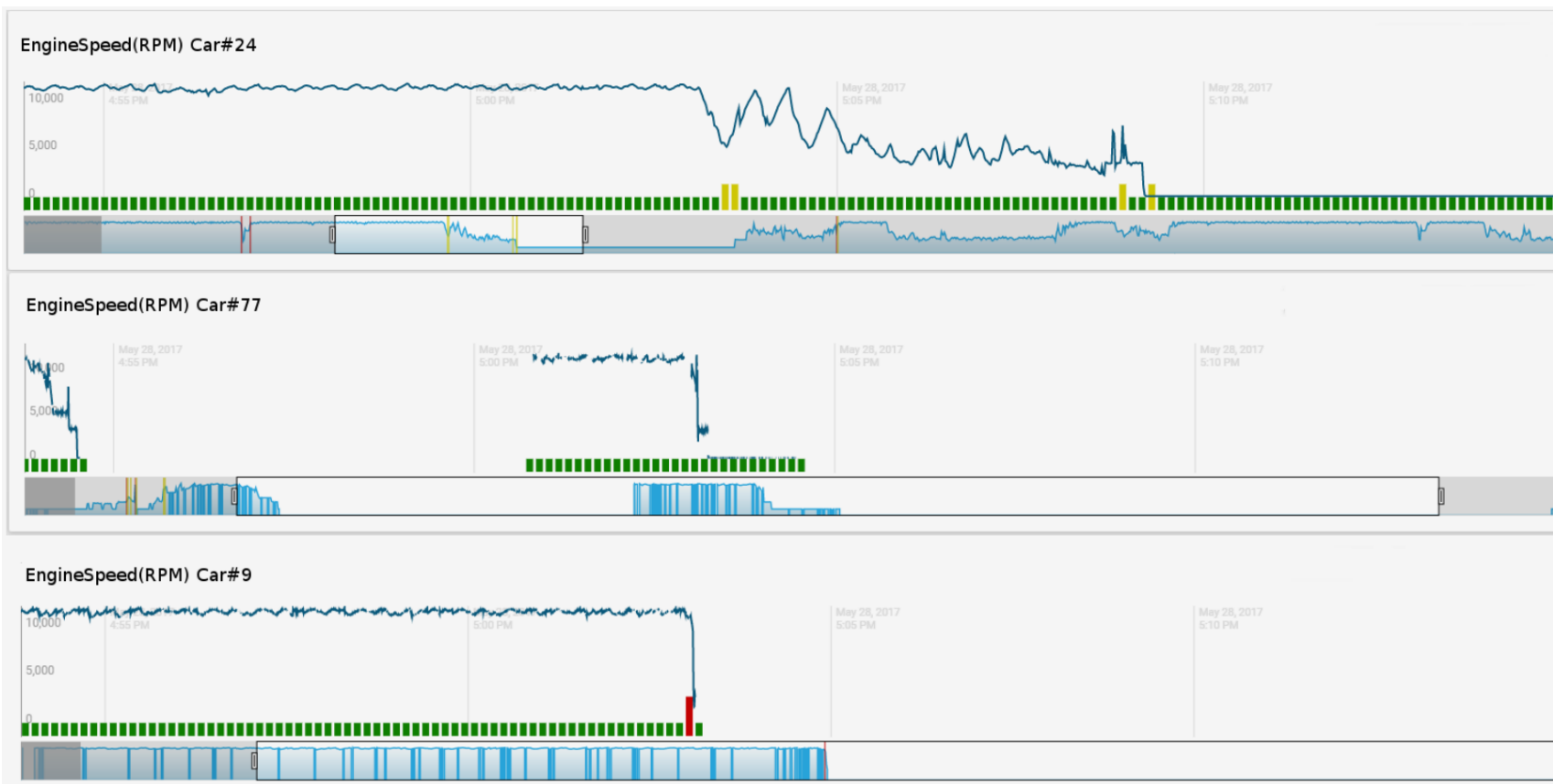


(b). Car#9 at end of the race

- Detection algorithm used has the capability to detect certain type of anomaly in few seconds ahead of the time.

# Anomaly Detection



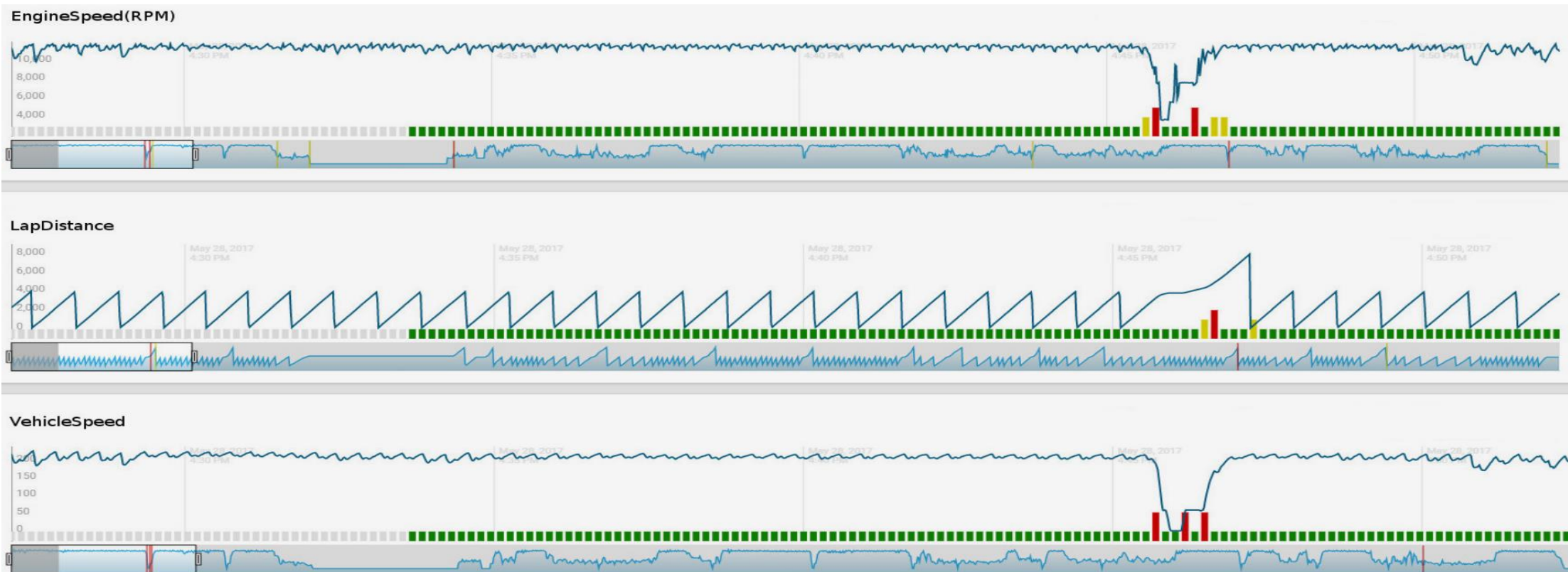


Correlation of “events” among different cars: #9 #77 and #24

- The anomaly occurs around 15:03 pm, where the RPM of car #9 totally disappeared. In fact, car #9 got totaled due to a collision with car #77. The others cars, including car #24, all slowed down after the crash.

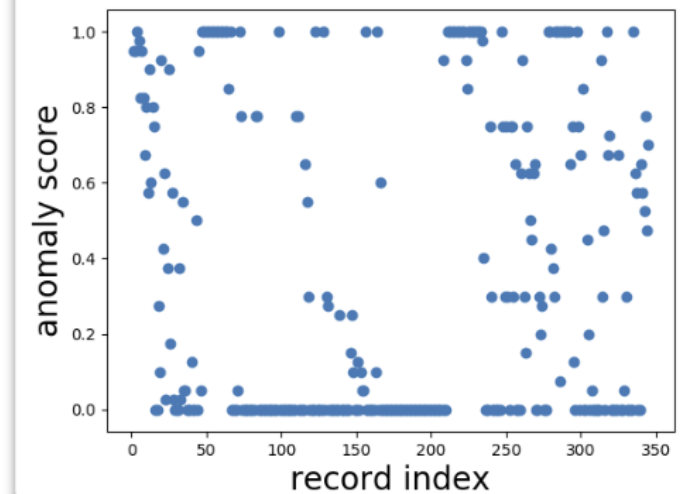
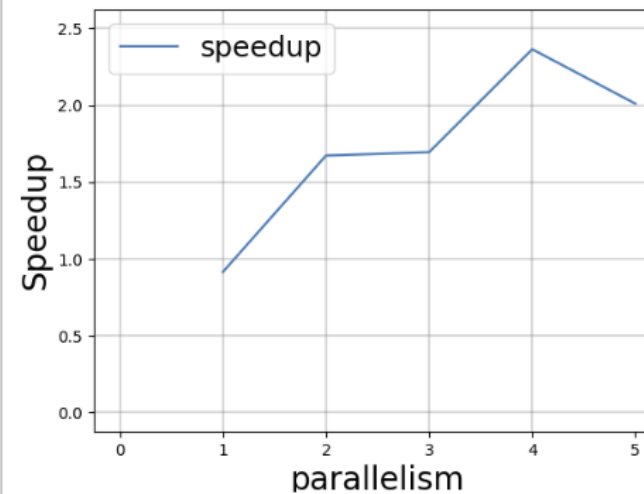
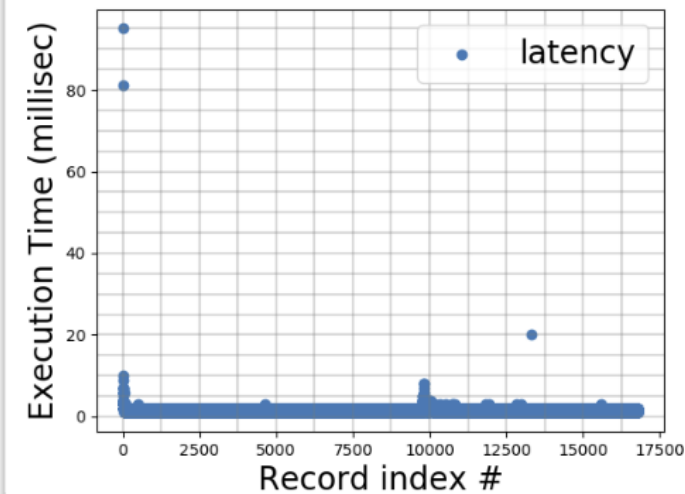
# Correlation Analysis





Correlation of “metrics” for one car #11.

# Correlation Analysis

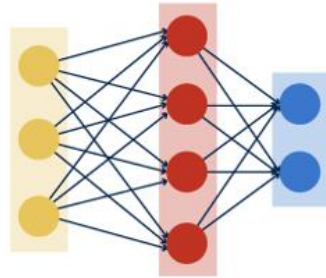


- The preliminary tests were done on the eRP data for car #9 (Scott Dixon), with ~17,300 records.
- The first plot shows the execution latency (time taken to predict anomaly on a single input record) to process each record in car #9 data on Apache Storm framework with parallelism of 1. Average execution time to predict an anomaly ~ 1.43 milliseconds.
- Plot 2 shows time speedup gain on increasing the parallelism in Storm to process car #9 telemetry. Speedup is the ratio of time to predict anomalies on a serial process to the time taken by running Storm at various parallelism levels on same data.
- To establish a sense of ground truth on labels and validate our approach, we deliberately inject anomalies at 2% data fraction (~ 345 data points) at known indexes. We inject speeds of 0.00 mph (absolute vehicle halt) at various indexes. The anomaly scores of each of the 345 data points is shown in plot 3. Anomaly score of 0.0 implies a normal event.

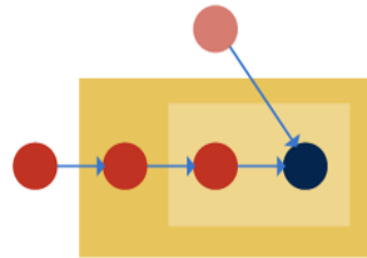
# Validation



K-Means



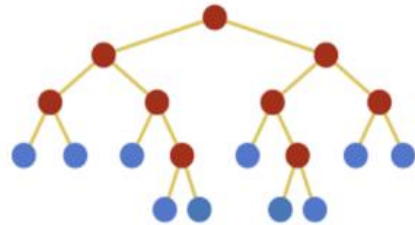
Neural Networks



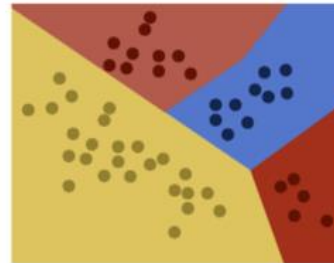
Latent Dirichlet Allocation



Support Vector Machine



RandomForest

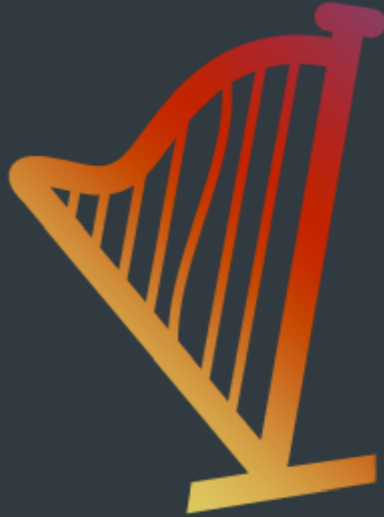


Multi-Class Logistic  
Regression

Harp/Harp-DAAL is an HPC-Cloud convergence framework that aims to automate ML as a service for both ease of use and scalability. Harp is designed to cover a full range of data-intensive computation from pleasingly parallel to machine learning and simulation.

# Data Analytics and Machine Learning

Open Source Github Website (<https://dsc-spidal.github.io/harp>)



# Harp-DAAAL

in collaboration with



is a high performance framework with the fastest machine learning algorithms on Intel's Xeon and Xeon Phi architectures.

See how it works

Performance

Explore algorithms

Hands on

Slide deck