# About



**Irfan Ahmad**
CachePhysics Cofounder
CloudPhysics Cofounder
VMware (Kernel, Resource Management),
Transmeta, 40+ Patents
University of Waterloo
@virtualirfan

# CachePhysics

**Data Path Monitoring and Modeling Software**
Real-time Predictive Modeling of Data Access Patterns
Increasing Performance & Cost Efficiency of Existing Caches
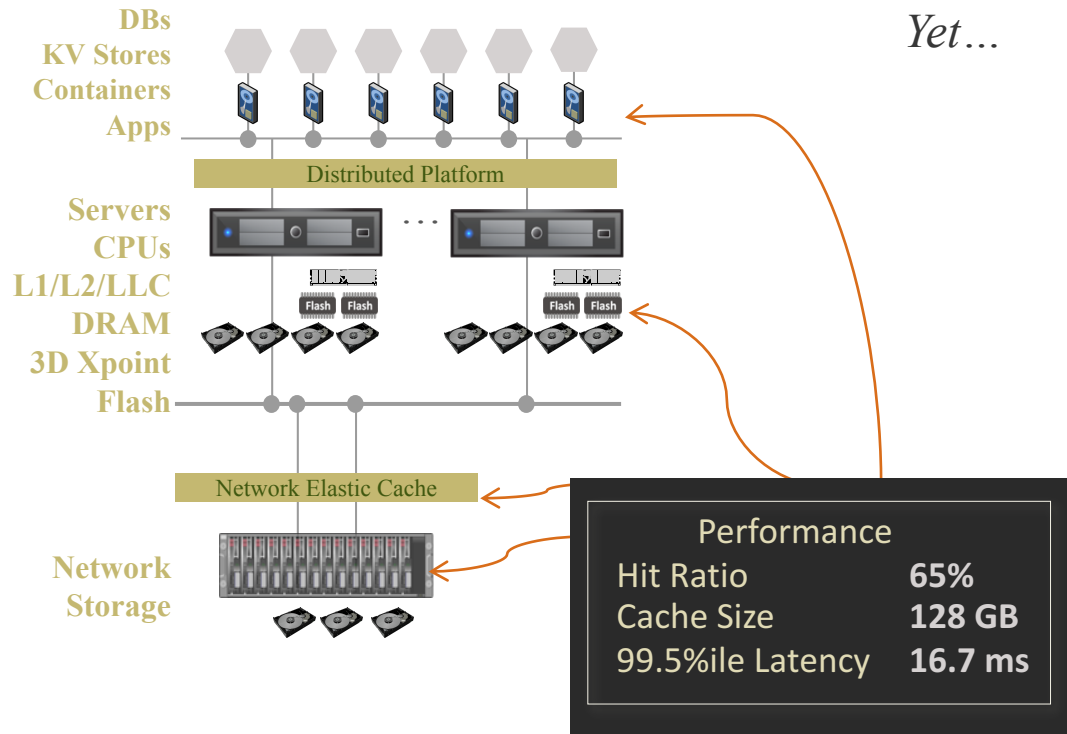Powering Next-Generation Self-Learning Caches

**2**

In-Memory Computing SUMMIT | NORTH AMERICA 2018

# In-Memory Computation is Hard



App Data needs changing daily. Providing QoS has become hard

In-Memory Computing SUMMIT NORTH AMERICA 2018

# Data Path Getting *More* Complex



The problems are getting much worse with increasing hardware complexity

# Data Path Performance is Critical



*Yet…*

Performance

| | |
|---|---|
| Hit Ratio | 65% |
| Cache Size | 128 GB |
| 99.5%ile Latency | 16.7 ms |

## Intelligent Management is Non-Existent

- Is this performance good?
- Can performance be improved?
- How much Cache for App A vs B vs …?
- What happens if I add / remove DRAM?
- How much DRAM versus Flash?
- How to achieve 99%ile latency of X μs?
- What if I add / remove workloads?
- Is there cache thrashing / pollution?
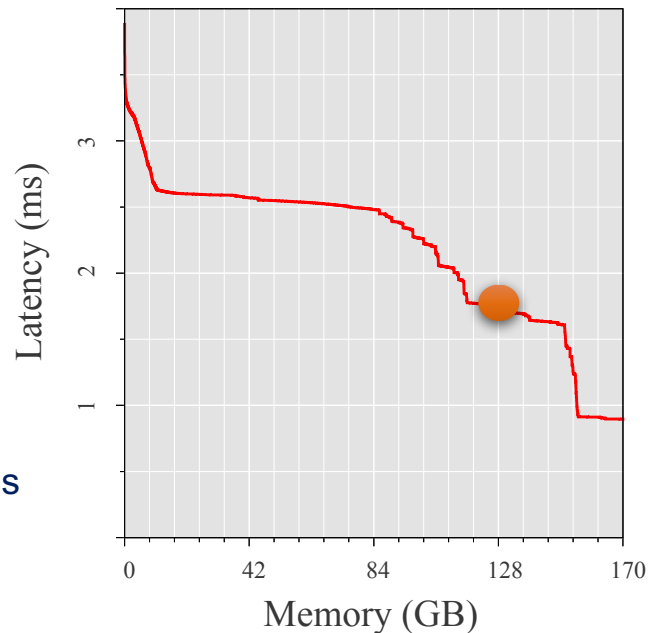- What if I change cache parameters?

5

In-Memory Computing SUMMIT | NORTH AMERICA 2018

# Algorithms to the Rescue

Lower is better

### Performance

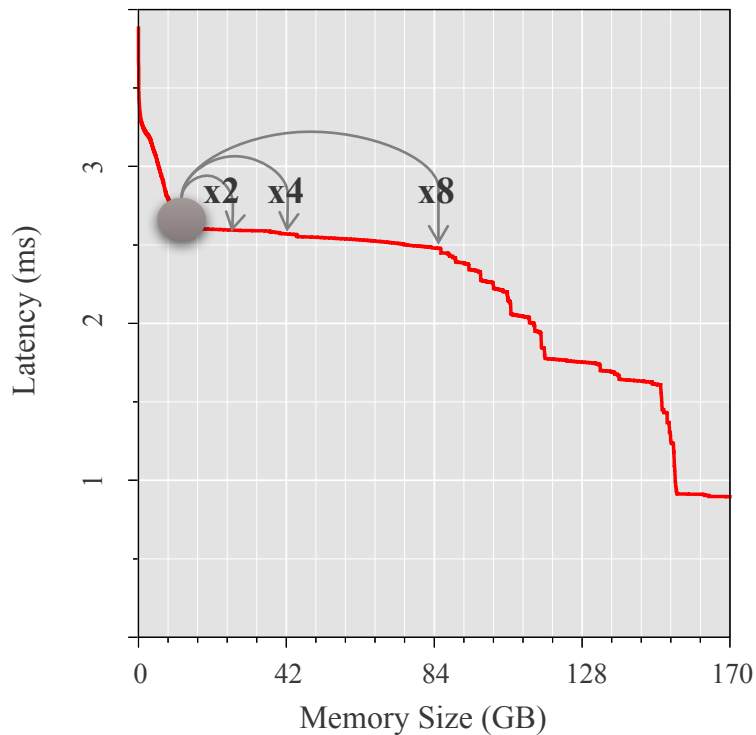| | |
|---|---|
| Hit Ratio | **65%** |
| Cache Size | **128 GB** |
| 99.5%ile Latency | **16.7 ms** |

Learn performance model of applications and cache

Predict the performance of workload as *f(memory size, params)*



Latency (ms) vs Memory (GB)

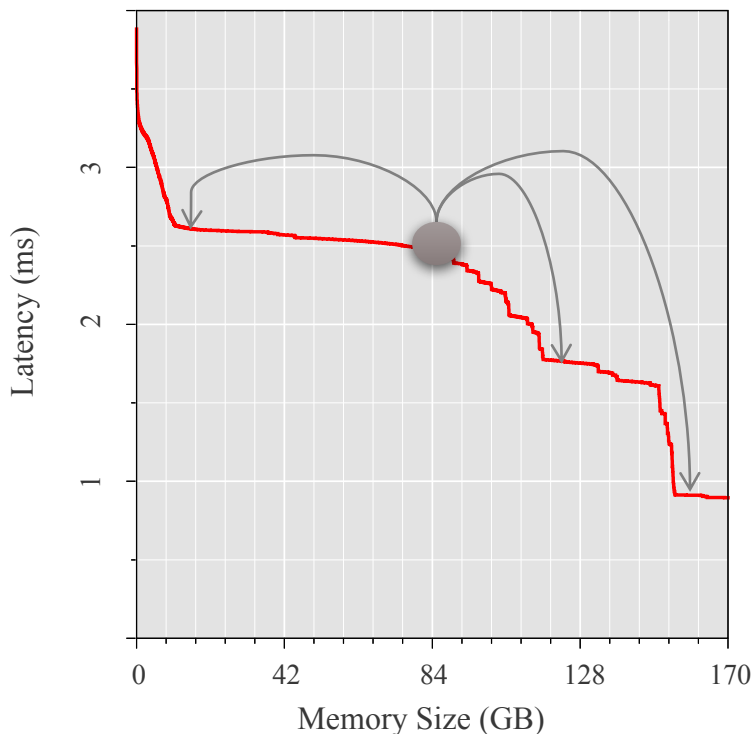# Understanding Performance Models



Lower is better

Latency (ms) vs Memory Size (GB)

x2   x4   x8

Decide useful
increments of change.

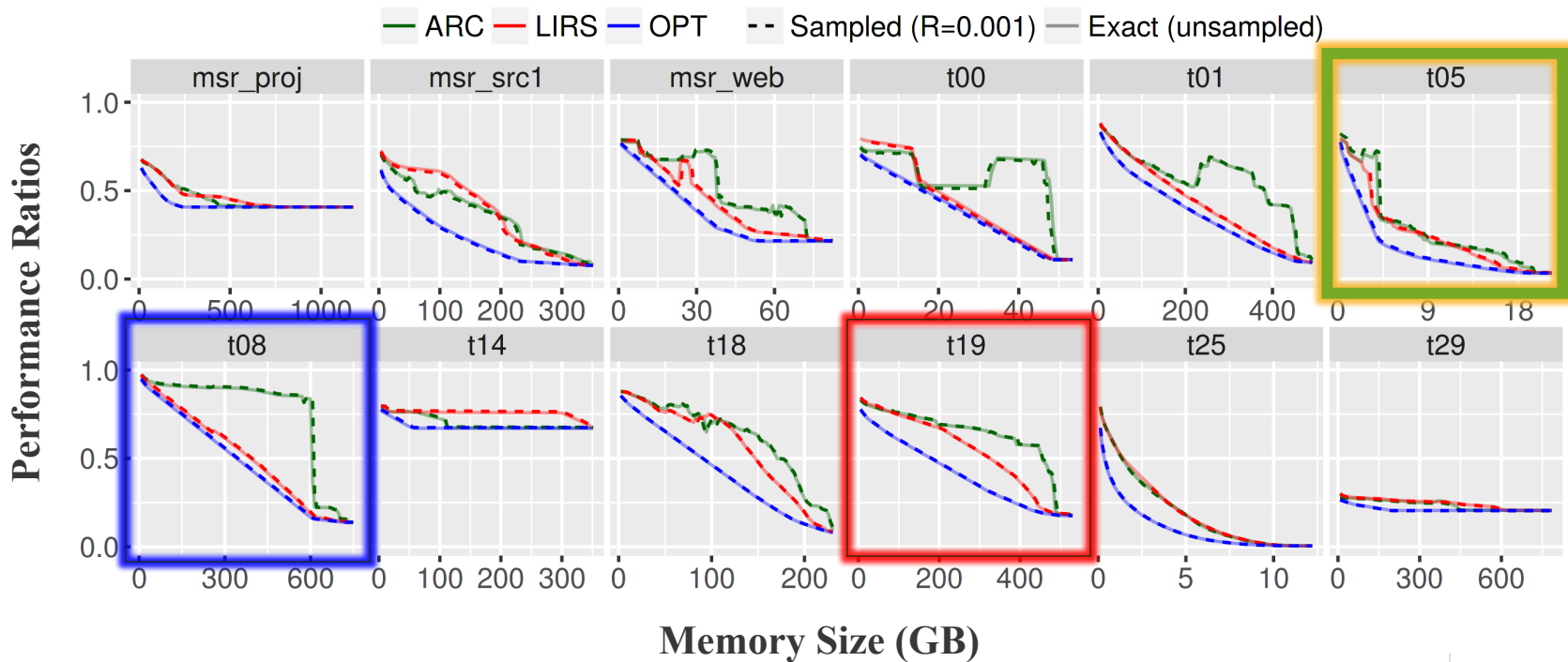In-Memory Computing SUMMIT | NORTH AMERICA 2018

# Understanding Performance Models (2)

Lower is better



The only 3 efficient operating points for this workload.

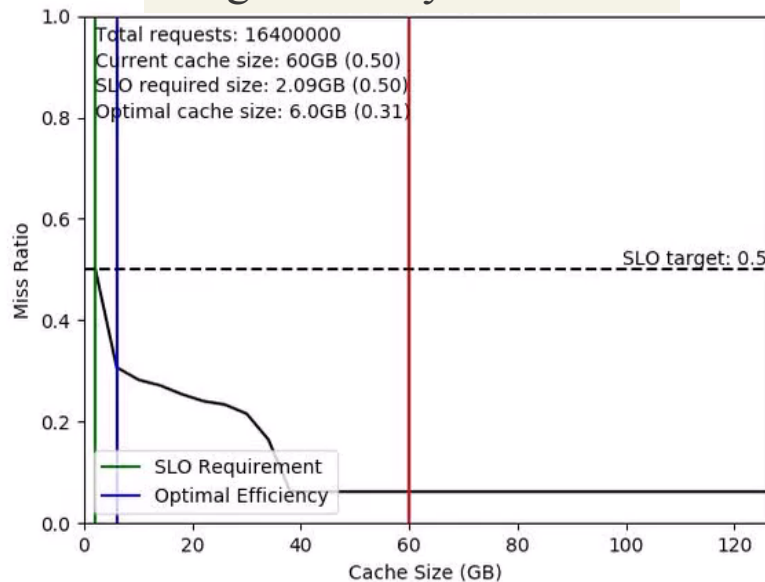Note: most operating points are highly inefficient.

8

In-Memory Computing SUMMIT | NORTH AMERICA 2018
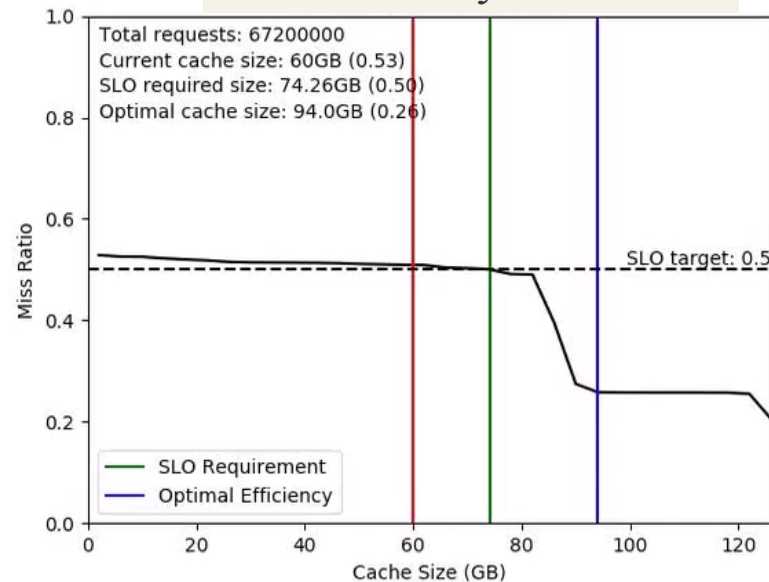
# Production Workload Performance

9

# Workload Behavior: Highly Dynamic



High Velocity Workload

Low Velocity Workload

# Use Case: Achieving Latency Targets



Latency Target (7 ms)

Cache Allocation (>16 GB)

Client target 95th %ile latency is 7 ms

Autoset cache partitions size to 16GB to guarantee avg latency SLOs

* Throughput targets can be implemented similarly

**11**

In-Memory Computing SUMMIT NORTH AMERICA 2018

# Use Case: Memory Thrash Remediation

# Use Case: Multi-Tier Sizing



Tier 0 (DRAM) allocation

Tier 1 (3D Xpoint)

Tier 2 (Local Flash)

Tier 3 (Remote Flash)

Network Misses

\* Can model network bandwidth as a function of cache misses from each tier

In-Memory Computing SUMMIT | NORTH AMERICA 2018

# Use Case: Per-App Memory Pools

- Improve aggregate cache performance

- Allocate memory based on benefit

- Prevent inefficient utilization / thrashing
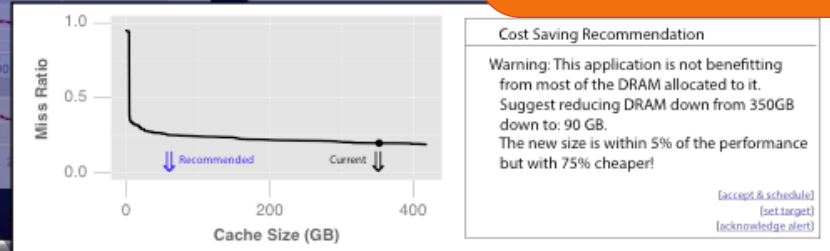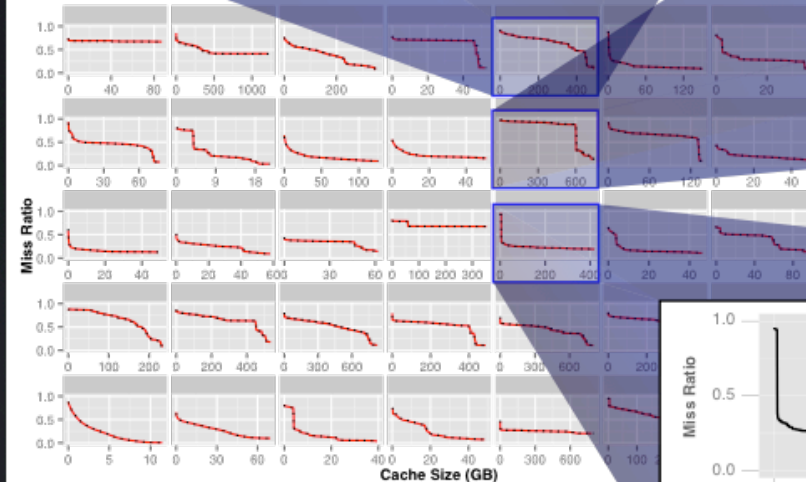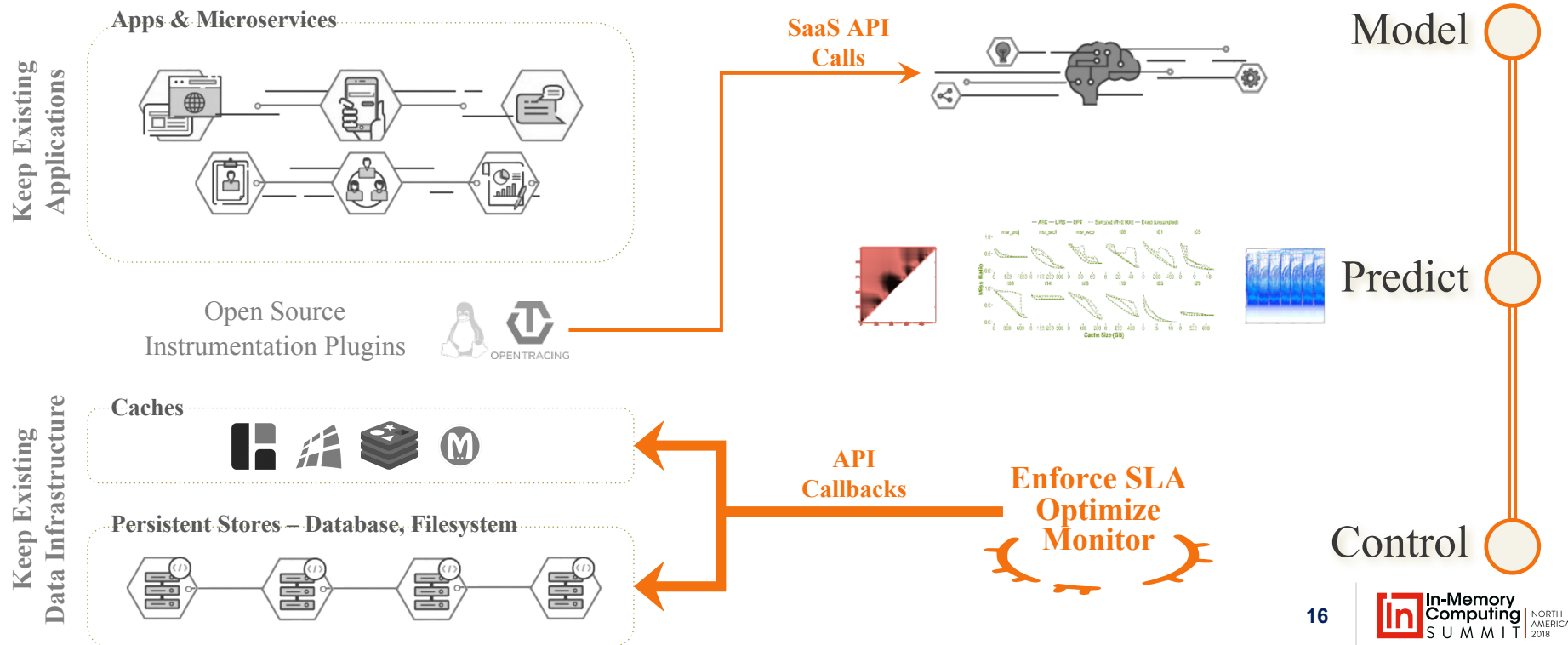
- Adapt to changing workload behavior

Client Allocation



Client 0                    Client 1

SLA Violation

Recommendation: increase DRAM Memory
This application has a performance SLA
than can no longer be achieved due to a
change in the workload.
The current size of 325 GB should be
increased to: 650 GB

[accept & schedule]
[change target]
[acknowledge alert]

SLA Risk: Application Workload Shifting

Warning: The working set profile of this
application has changed by more than 10%.
This application is still within it's SLA but it is
50% closer to a violation than before.
The Miss Ratio Curve started to change
on Nov 7, 2017. Contact the developers.

[file a bug report]
[change target]
[acknowledge alert]

**Features**
- Self-learn predictions for each client
- Alert, recommendations
- Recommendation/SLA API
- Capacity planning, what-ifs

Cost Saving Recommendation

Warning: This application is not benefitting
from most of the DRAM allocated to it.
Suggest reducing DRAM down from 350GB
down to: 90 GB.
The new size is within 5% of the performance
but with 75% cheaper!

[accept & schedule]
[set target]
[acknowledge alert]

MacBook

In-Memory
Computing
SUMMIT | NORTH AMERICA 2018

# Architecture

# Finally

## Takeaways
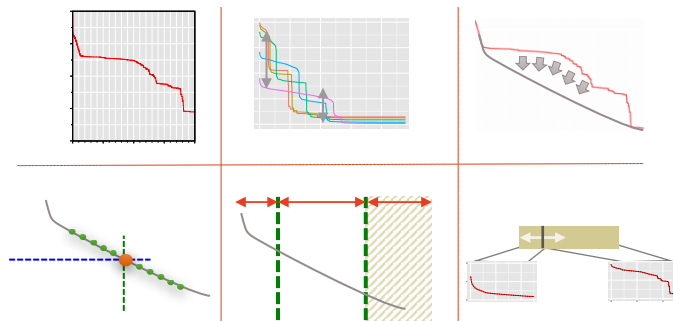
Tech breakthrough for In-memory Computing

- Optimal cost, same performance
- World 1st/only latency SLOs
- Self-tuning data path
- Auto-scaling data paths

Award-winning technology

## Asks

New Customer Projects

- 50% Efficiency Gain Guarantee!
- Latency SLO guarantee
- In-memory compute
- Database, Key-Value store, Filesystem, Disk system Optimization

# CachePhysics

irfan@cachephysics.com       650-417-8559    @virtualirfan

In-Memory Computing SUMMIT | NORTH AMERICA 2018