# Life @ the Edge: Challenges Accelerating and Specializing

Theophilus A. Benson

Brown University

# Azure Stack Edge Gets NVIDIA GPU To

**TECHNOLOGY NEWS**   MARCH 5, 2020 / 12:32 PM / 2 MONTHS AGO

# AT&T partners with Google Cloud for 5G edge computing

2 MIN READ
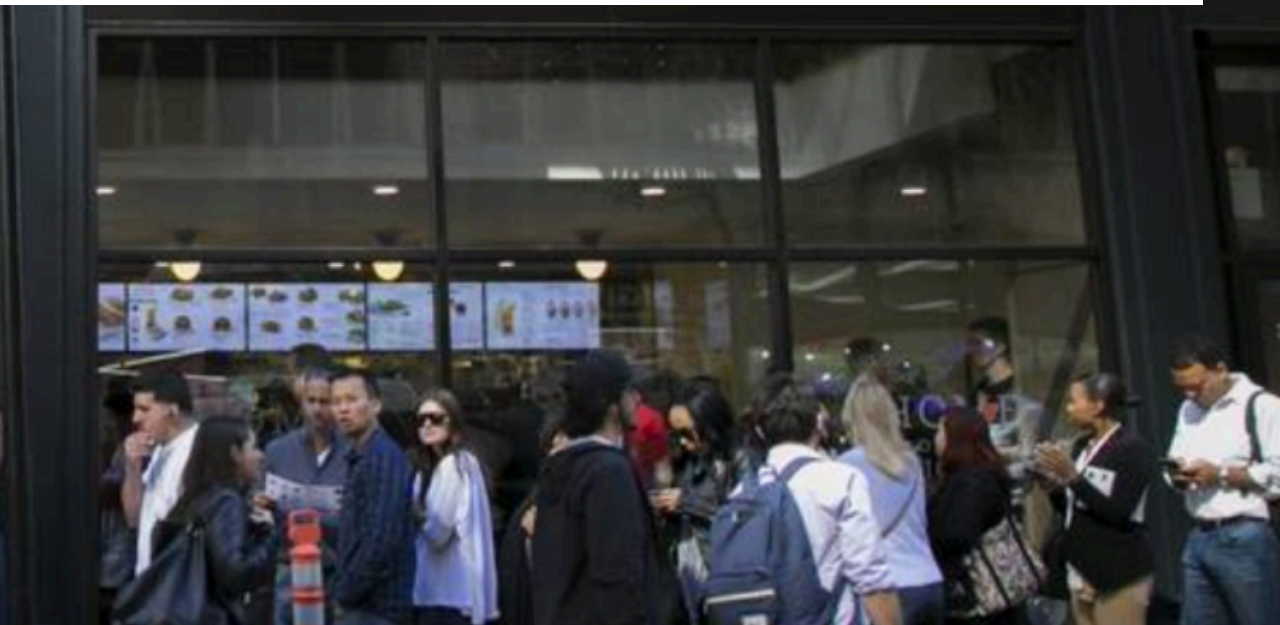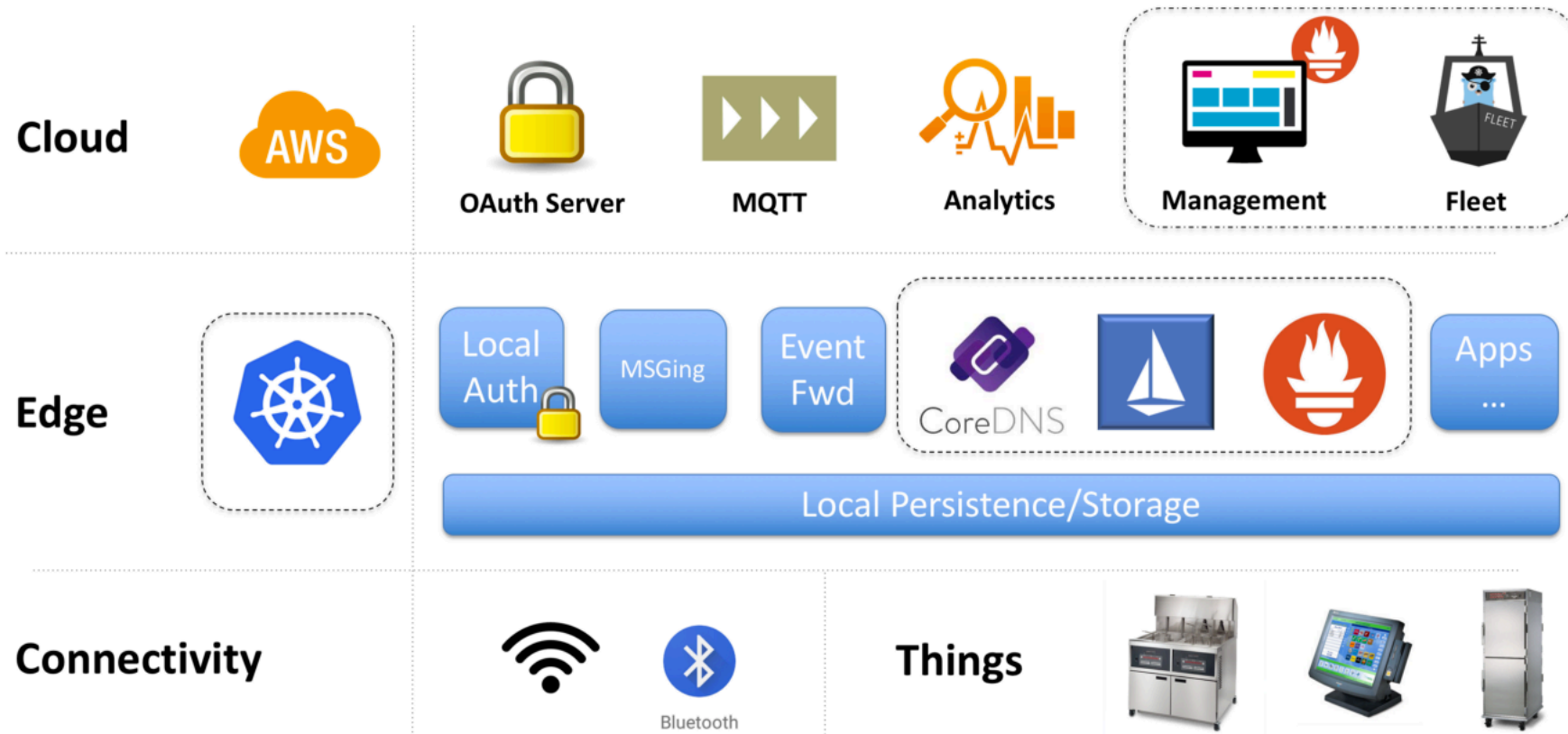
# 2000 edge deployments in the US?

Restaurant "Data Centers"
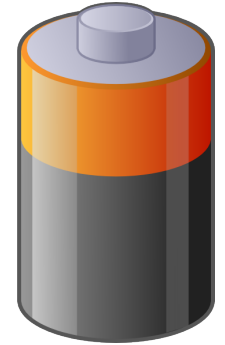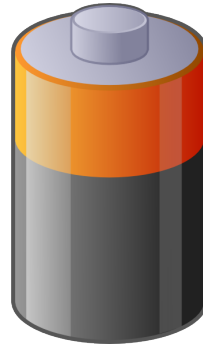
Intel: Quadcore processor, 8 GB RAM, SSD

Cloud

AWS

OAuth Server MQTT Analytics Management Fleet

Edge

Local Auth MSGing Event Fwd CoreDNS Apps ...

Local Persistence/Storage

Connectivity

Bluetooth

Things

Holistic Management of
Edge Accelerator for
maximum efficiency

Dynamic Specialization of
Edge Network Stack for Varying Objectives

Applications of
Edge Computing

AI   5G   IoT

Platforms for
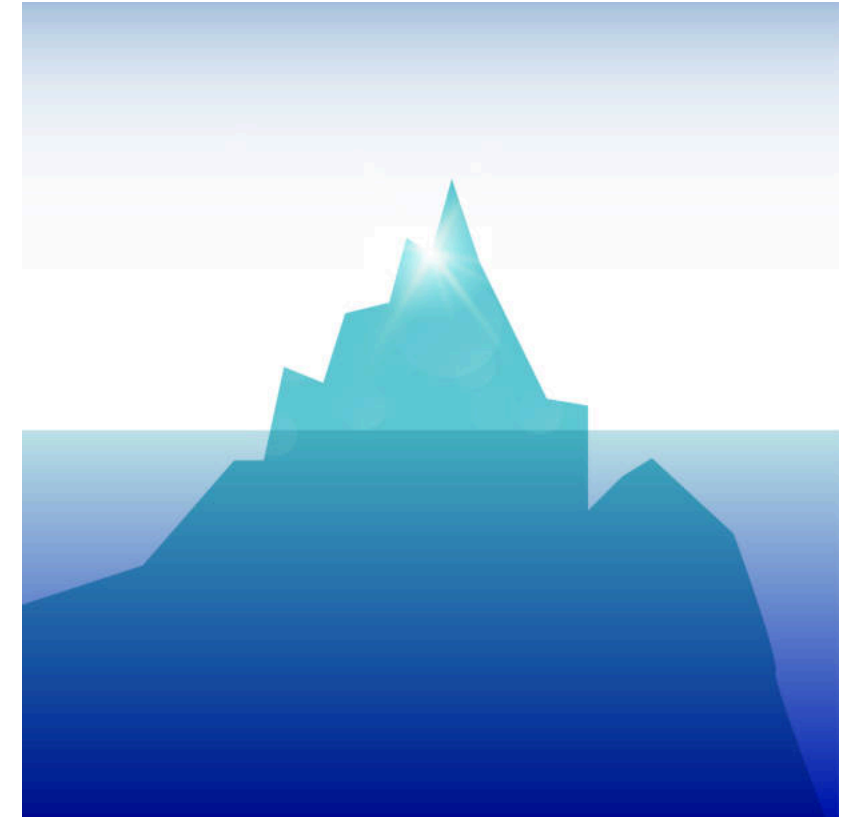Edge Computing

**Networking Stack**
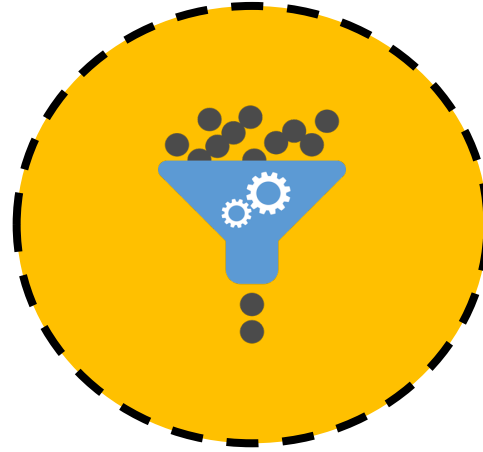
**Operating System**
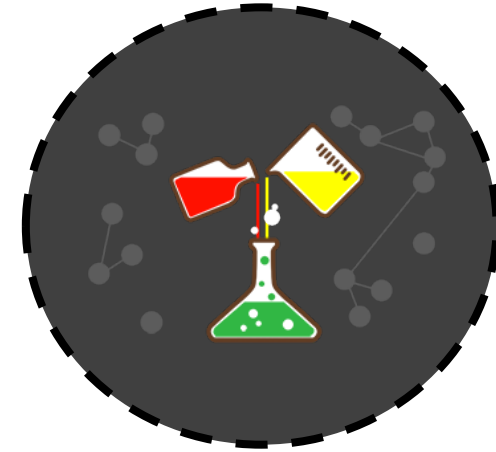
**CPU**   **GPU**   **FPGA**
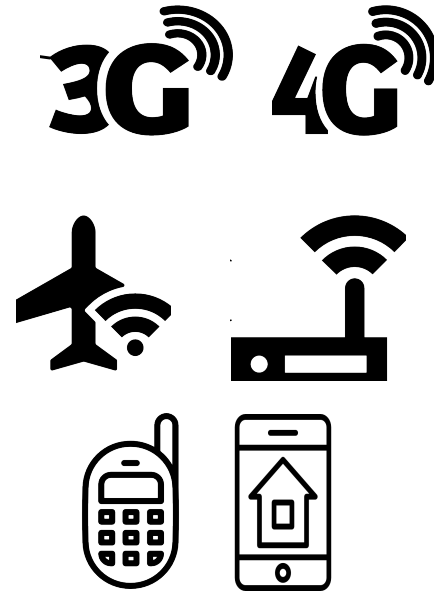
# RoadMap



**Motivations**
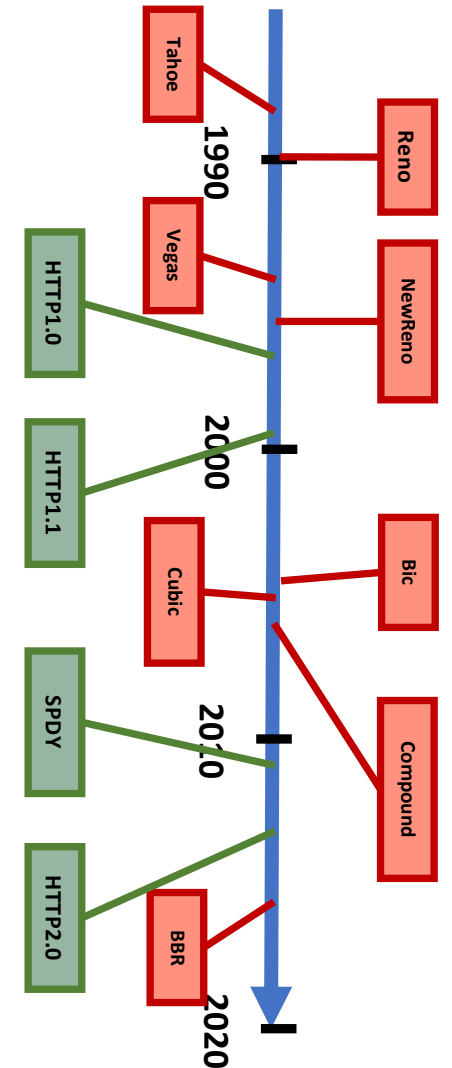
**Data-Driven
Specialization**

**Accelerator
Management**

- Networks evolve organically
  - Adaptation of new technology

- Protocols are redesigned
  - Encapsulate domain-specific insights

- Protocols are statically deployed
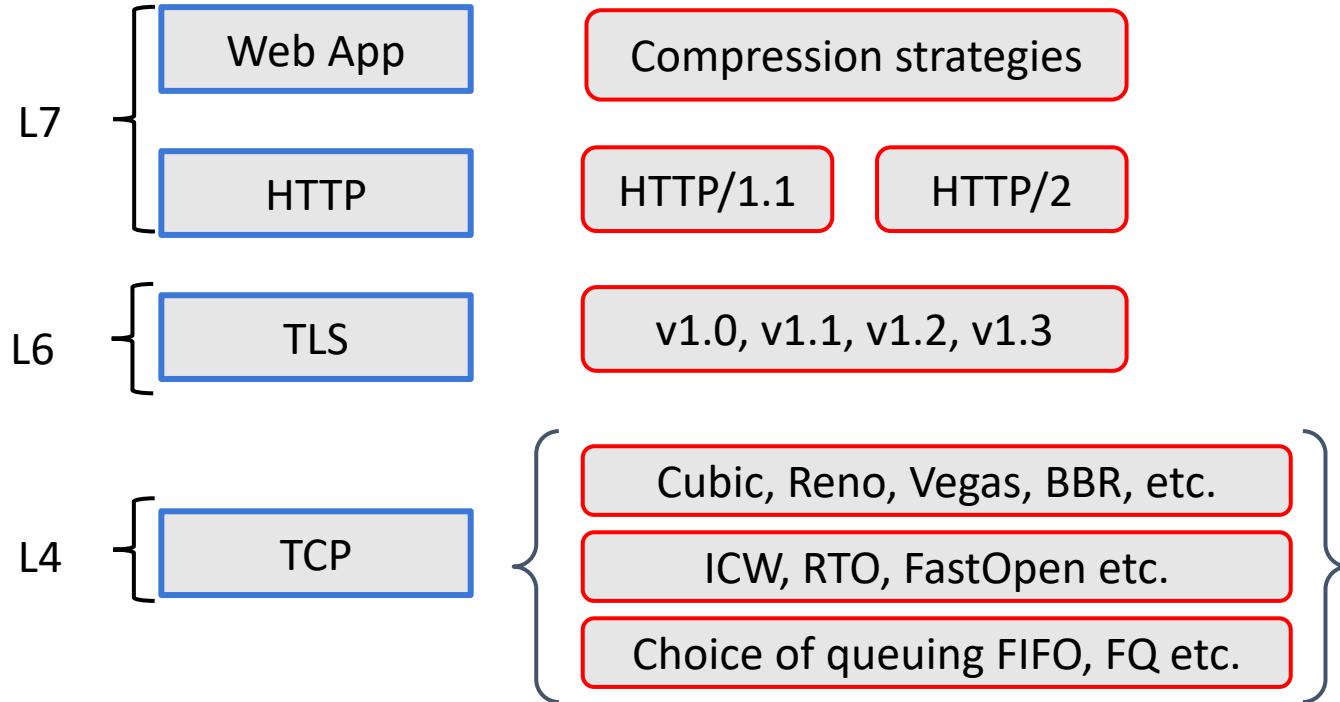  - Protocol use is agnostic of conditions

**What we can have**

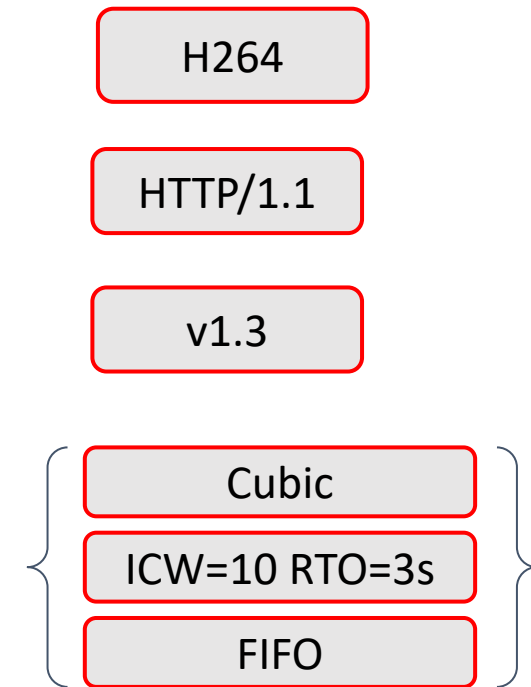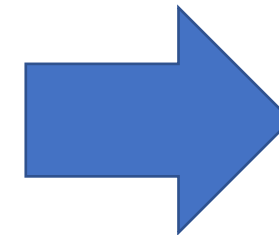**What is used**

L7

Web App

Compression strategies

H264

HTTP

HTTP/1.1  HTTP/2

HTTP/1.1

L6

TLS

v1.0, v1.1, v1.2, v1.3

v1.3

L4

TCP

Cubic, Reno, Vegas, BBR, etc.

ICW, RTO, FastOpen etc.
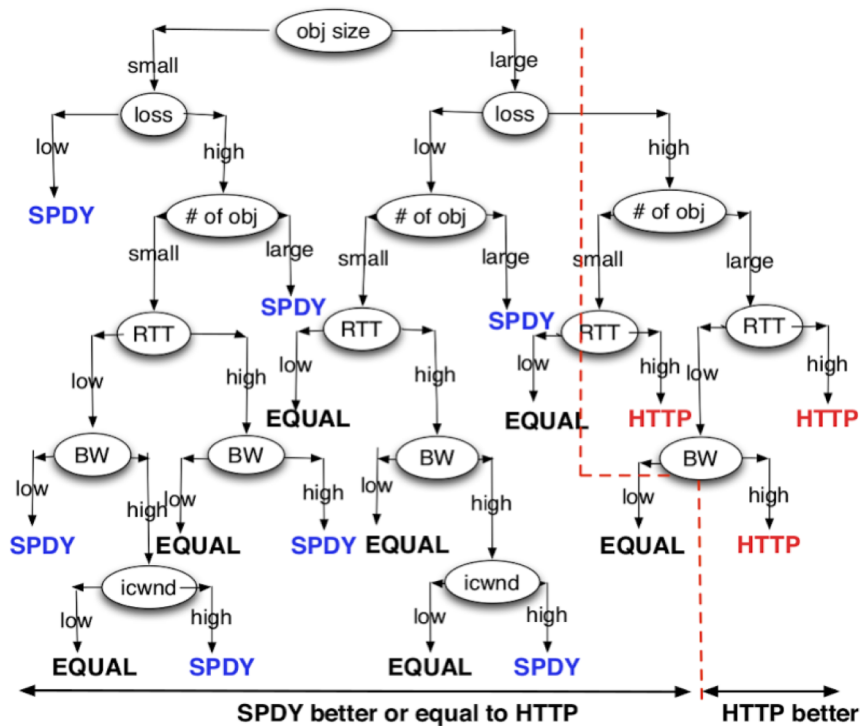
Choice of queuing FIFO, FQ etc.

Cubic

ICW=10 RTO=3s

FIFO

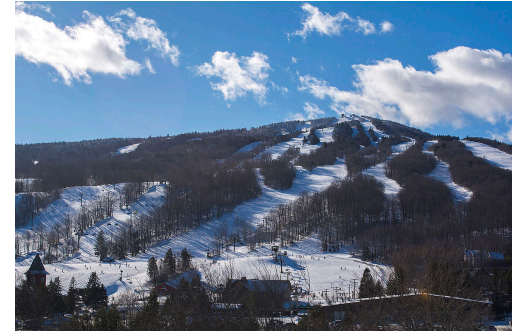**one size fits all to address diverse networks and end user devices!!!**

13

# Web Performance



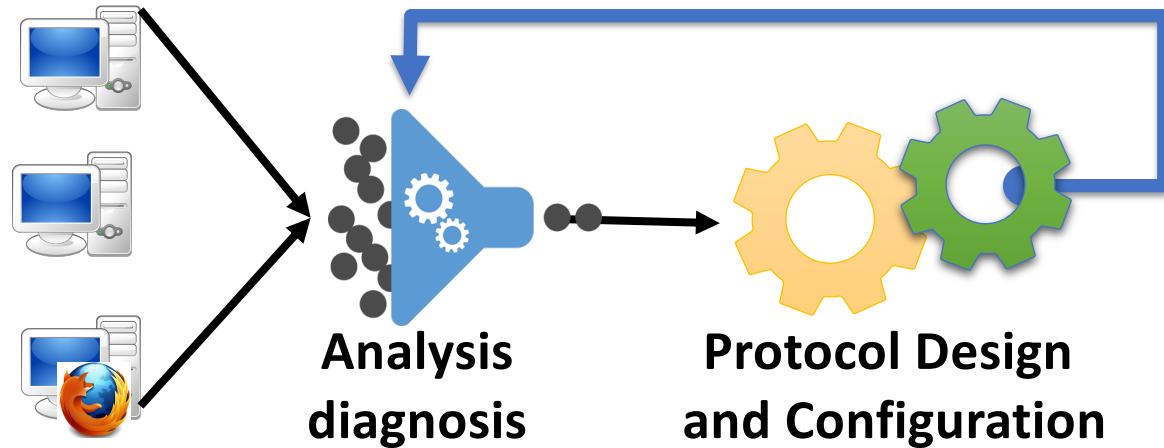How speedy is SPDY? [NSDI '14]

- How Speedy is SDPY? [NSDI '15]
  -HTTP versions

- Verus [NSDI '15]/Pantheon [USENIX ATC'18]
  -TCP versions

- How quick is QUIC? [IEEE ICC '16]
  -QUIC vs HTTP2 vs HTTP1.1

- Overclocking the Yahoo! [IMC '11]
  -HTTP pipelining, TCP ICW, ABC

How do you pick which car?



**Garage of Cars**

Infrastructure consolidation:
- Providers own data center, edge, client
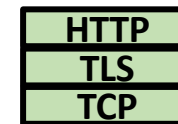
Well defined web metrics
- Web pages: SpeedIndex
- Video: Netflix metric

Open Challenges:
- Domain specific learning algorithm
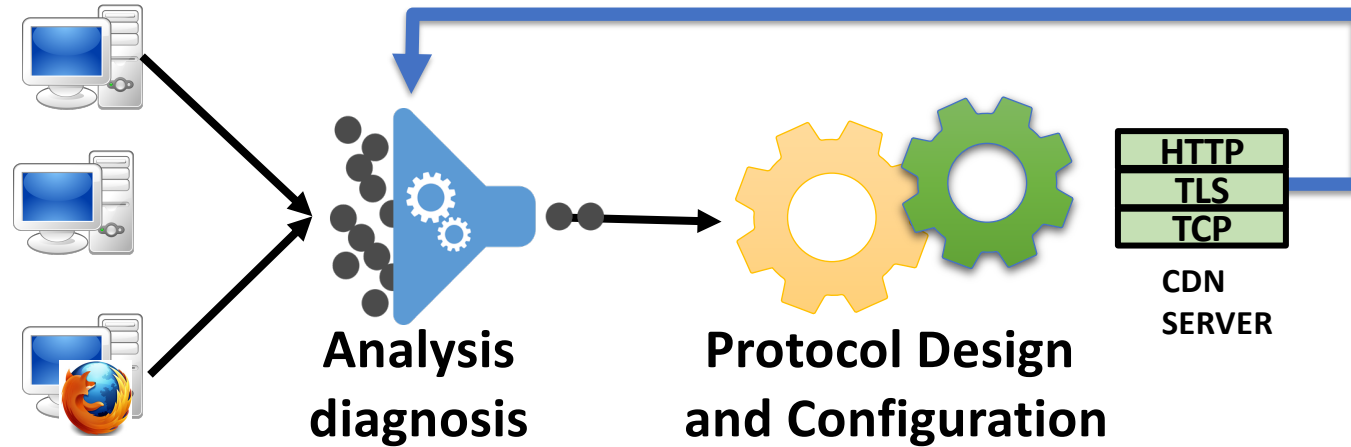- Fine-grained, low-overhead knobs

**Analysis diagnosis**

**Protocol Design and Configuration**

**HTML** **CSS** **JS**

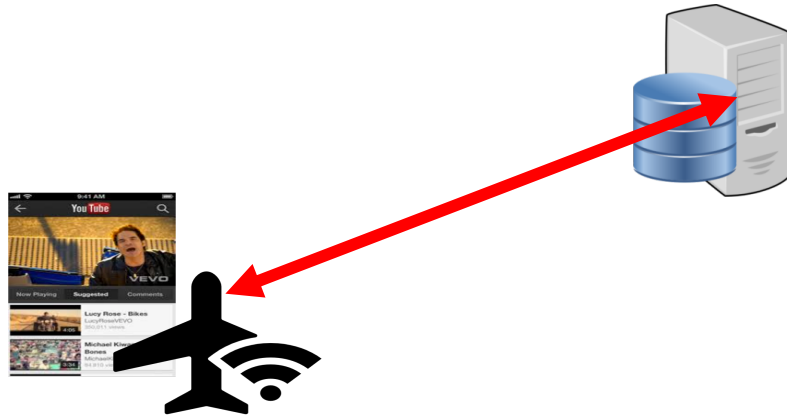**CDN Application**

| HTTP |
| TLS |
| TCP |

**CDN SERVER**

# Barriers to Fine-grained and Dynamic Tuning



- Algorithm learning challenges
  - Need effective and accurate predictive algorithms

- System design challenges
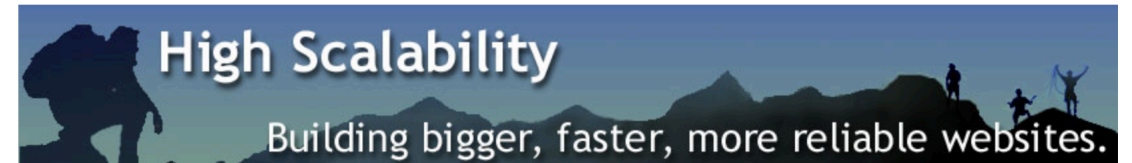  - Server must support low overhead fine grained tuning

# Strawman Approach

- For each client
  - Greedily test each configuration
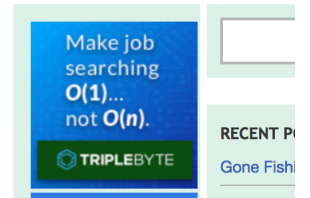  - Pick optimal configuration

**Challenges**
- **Network is dynamic**
- **Data collection is costly**
- **Data can be noisy (non-gaussian noise)**



High Scalability
Building bigger, faster, more reliable websites.

Home    Explain The Cloud Like I'm 10    All Time Favorites    Real Life Architectures    Strategies    Support In Patreon    Advertising    Start Here    Contact    RSS    Twitter    Facebook    G+    All Posts    Amazon Store

« Handle 700 Percent More Requests Using Squid and APC Cache | Main | Paper: Parallelizing the Web Browser »

**Latency Is Everywhere And It Costs You Sales - How To Crush It**

Make job searching $O(1)$... not $O(n)$.
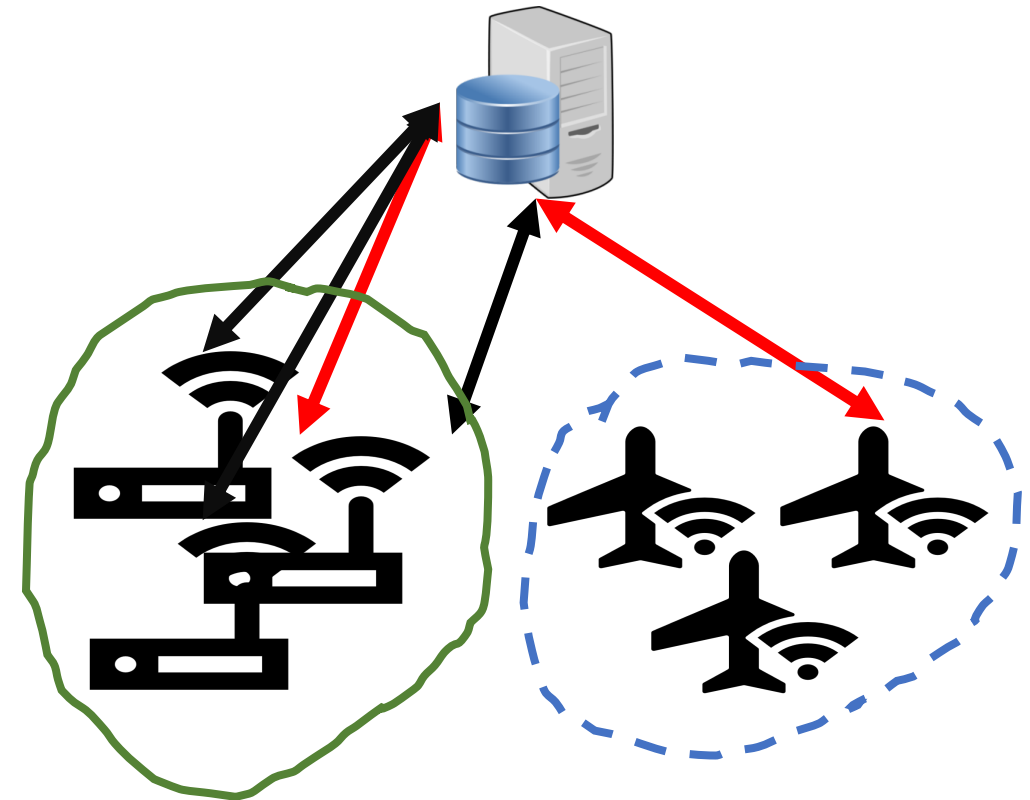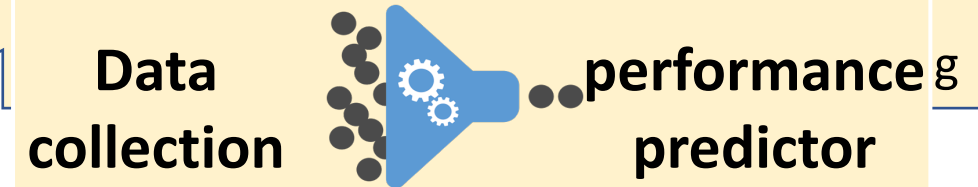
RECENT P
Gone Fishi

18

# Challenges

- Costly data collection
  - Distribute cost over similar users
  - Generalize observations

- Network dynamics
  - Online exploration of network

- Noisy data (non-gaussian noise)
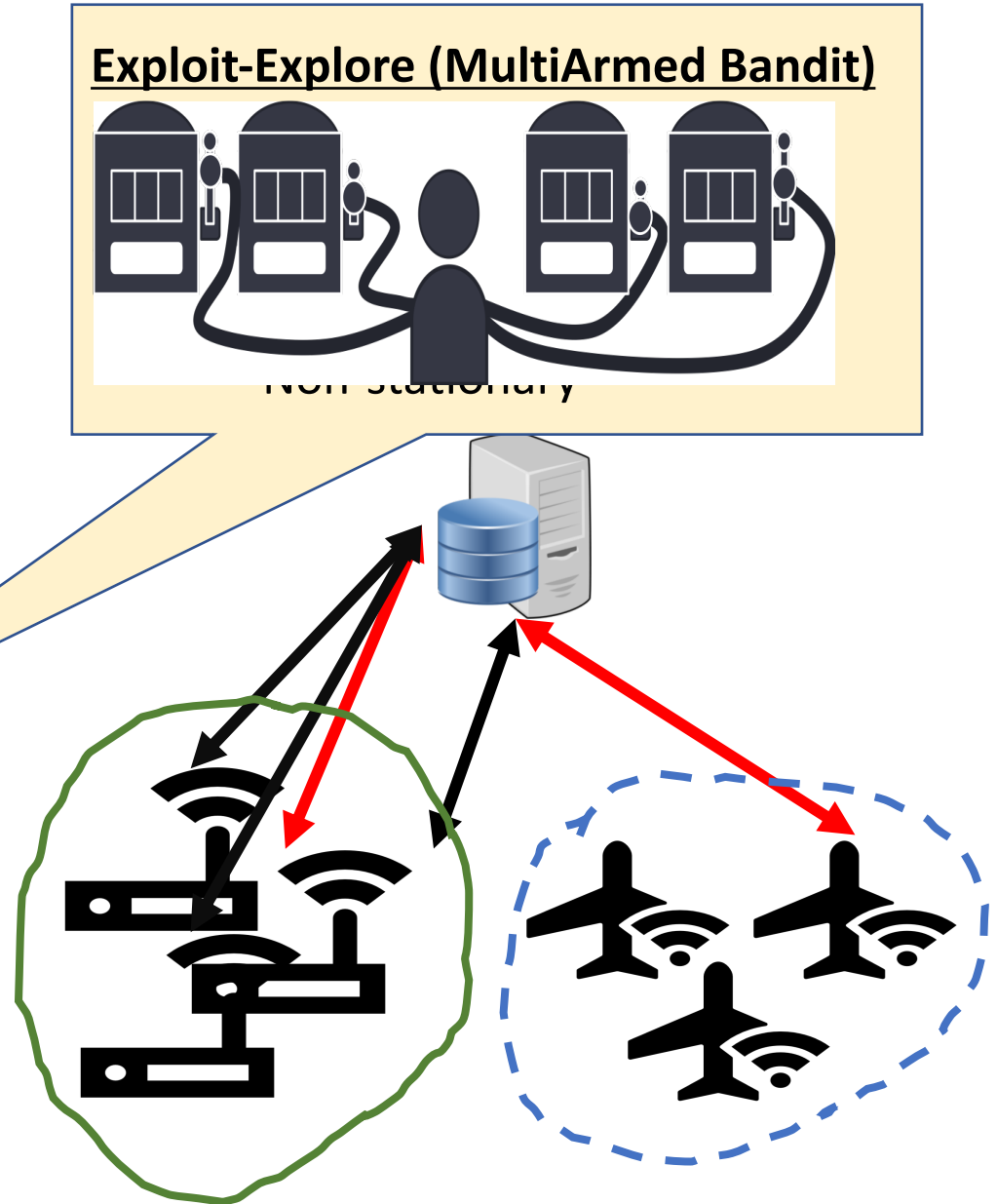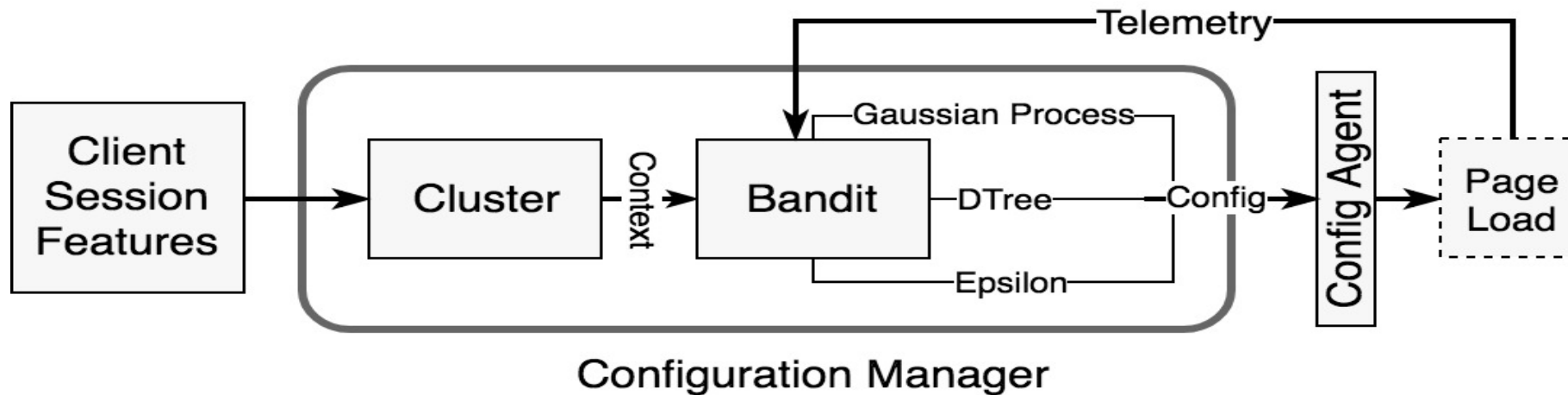  - Resample to get clean data

**Predictors**
- Machine learning
  - Deep learning (CNN, RNN)

**Data collection** → **performance predictor**

# Challenges

- Costly data collection
  - Distribute cost over similar users
  - Generalize observations

- Network dynamics
  - Online exploration of network

- Noisy data (non-gaussian noise)
  - Resample to get clean data



**Exploit-Explore (MultiArmed Bandit)**
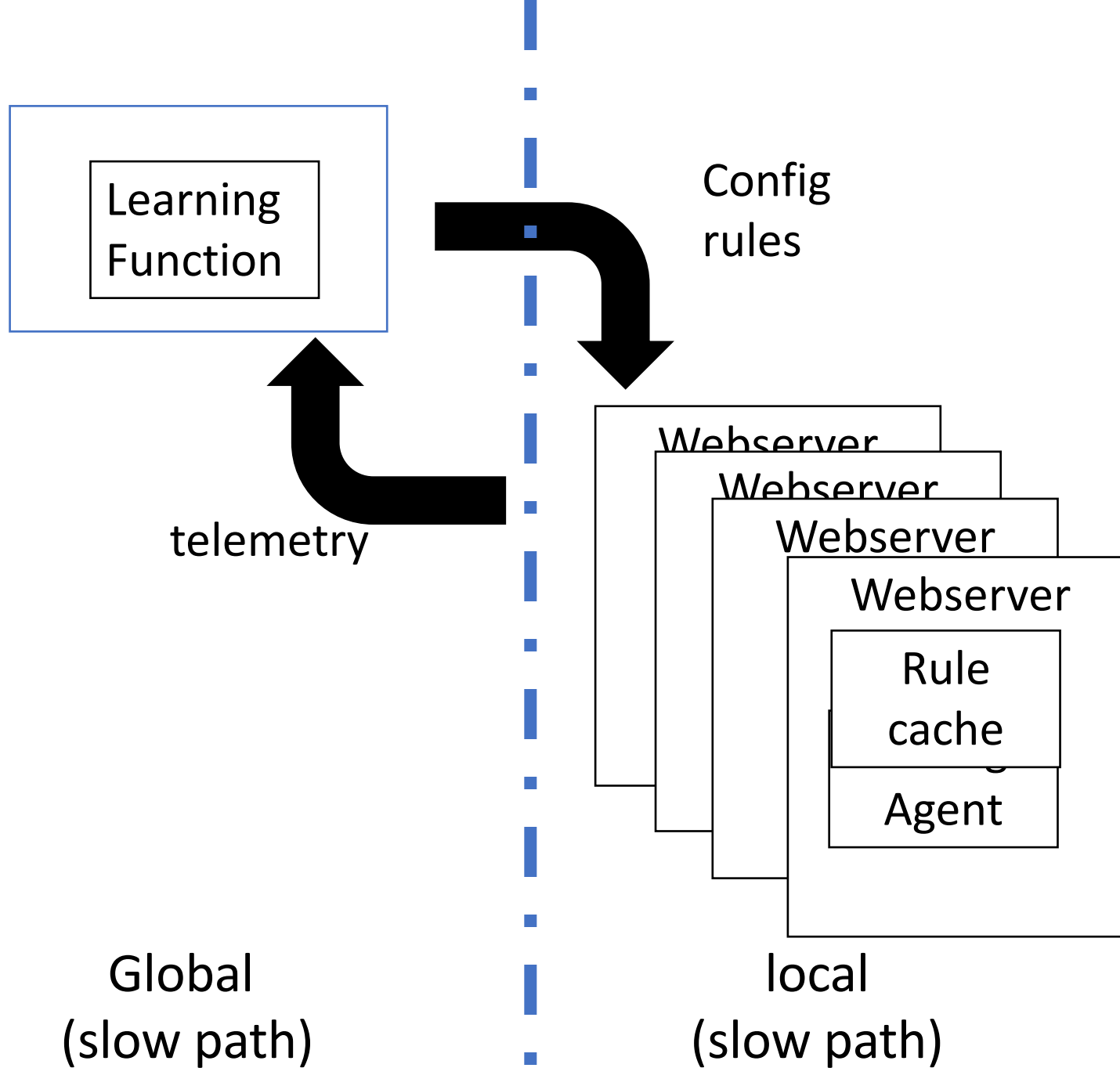
Non-stationary

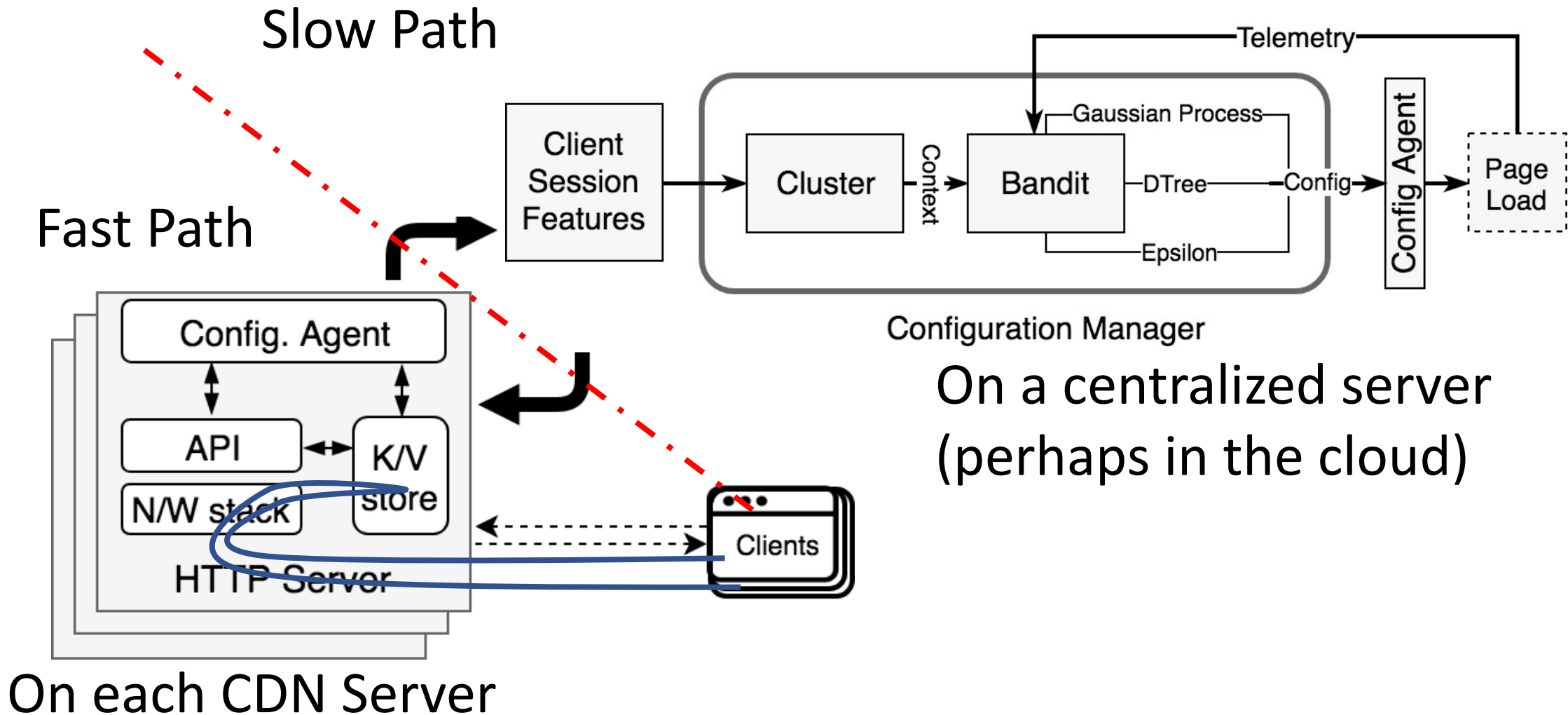# Domain-specific learning ensemble



- Contextual Multi-armed bandit: explore-exploit!
  - Contextual-exploration: cluster provides context
    - Non-Guassian: random exploration
    - Efficient learning: Gaussian exploration
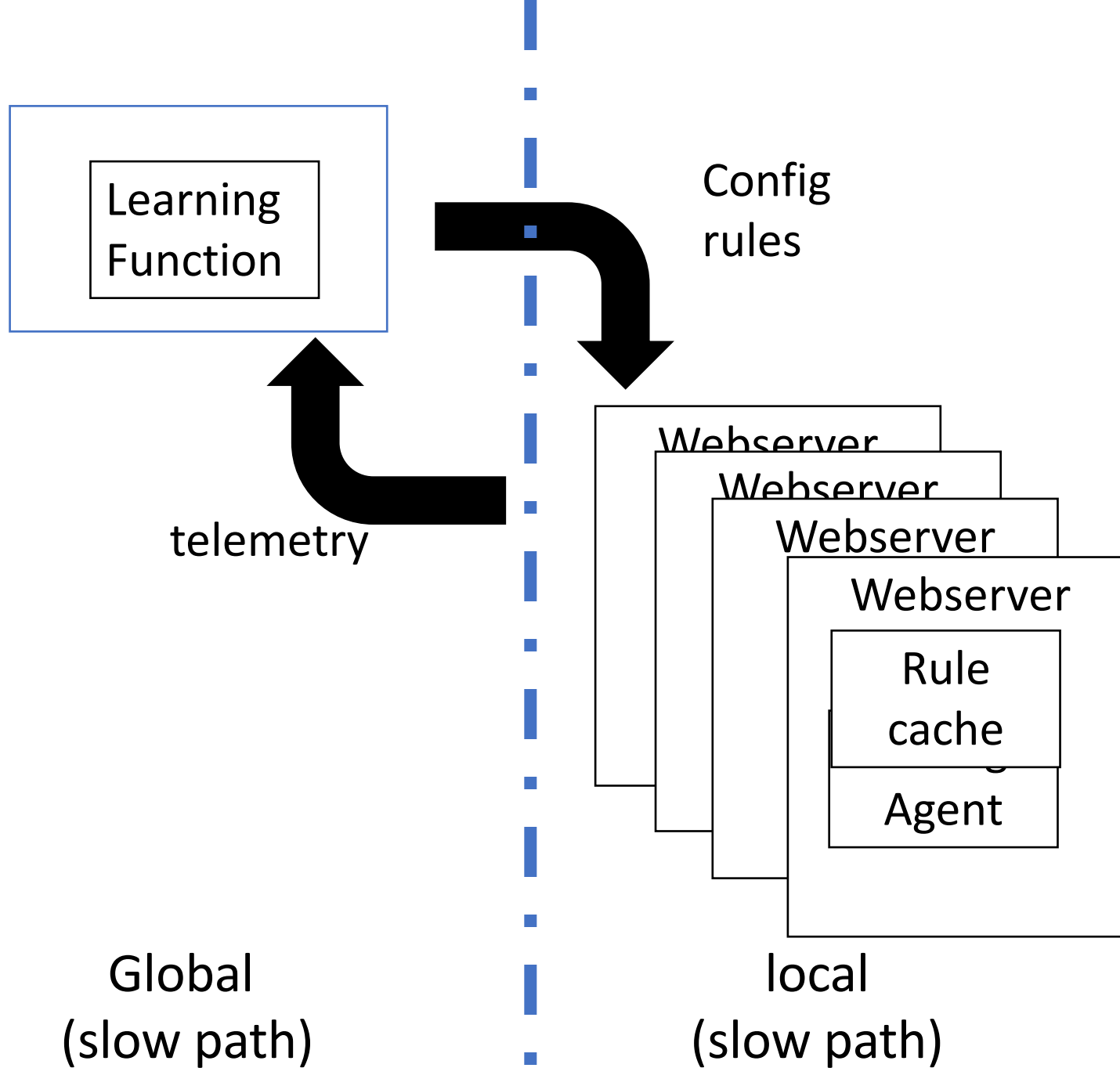  - Exploitation Arm:
    - D-Tree: Generalize observation

# Systems Support for Reconfiguration

| | Flexibility | Overheads |
|---|---|---|
| One VM per Configuration | 👍 | 👎 |
| One Container per Configuration | 👎 | 👍 |
| User Space Tuning + Kernel Module | 👍 | 👎 |
| Kernel Space Tuner+ kernel Module | 👍 | 👍 |

Learning Function

Config rules

telemetry

Webserver

Webserver

Webserver

Webserver

Rule cache

Agent

Global (slow path)

local (slow path)

23

# ConfigTron



Slow Path

Fast Path

On a centralized server (perhaps in the cloud)

On each CDN Server

24

Learning Function

Config rules

telemetry

Webserver
Webserver
Webserver
Webserver

Rule cache

Agent

Global
(slow path)

local
(slow path)

25

# Systems Approach to Configuration Tuning

- Manual protocol tuning
  - Coarse-grained and slow
    - Labor intensive
    - Slow to adapt to changes

- Bayesian-based tuning
  - DB: iTuned [VLDB'9], StarFish [CIDR'11]
  - Apps: WISP [SoCC'17], Ernest [NSDI'16]
  - Cloud: CherryPick [NSDI'17]
  - Fine-grained but static
    - Quickly finds "good" configuration
    - Doesn't adapt to network dynamics
    - Doesn't account for non-gaussian noise

- ConfigTron is fine-grained and dynamic

# Objective Functions



ClientDevice Specific

Provider Specific

Application Specific

# Configuration Knobs



**CDN Application**

**CDN Caching**

A
B

**CDN SERVER**

HTTP
TLS
TCP

**CDN Network**

BGP PEER

BGP PEER

| Stability | Fairness | Security |
|-----------|----------|----------|

http://tradingeveryday.com/wp-content/uploads/2016/08/fastdotcom.png

# RoadMap



**Motivations**
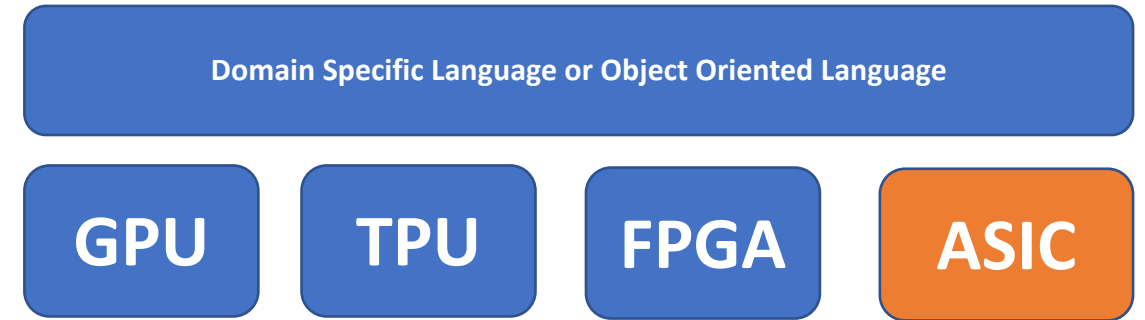
**Data-Driven Specialization**

**Accelerator Management**

GPU Virtualization
- CAVA[ASPLOS'20]
- Time-division multiplexing

FPGA Virtualization
- Virtal [ASPLOS'20]
- Space-division multiplexing

- Rich body of work on hardware accelerators
  - Class one: Siloed approach → limits global orchestration and management

- Rich body of work on hardware accelerators
  - Class one: Siloed approach → limits global orchestration and management
  - Class two: Unified approach → disconnected from lower layer optimizations/advances

- Rich body of work on hardware accelerators
  - Class one: Siloed approach → limits global orchestration and management
  - Class two: Unified approach → disconnected from lower layer optimizations/advances

- Neither allows efficient management of a fabric of accelerators

# Ideal Cloud Accelerators Fabric

- ``Accelerators'' and not developers are the first-class citizens
  - Simplify life-cycle management

- Provide ``big fabric'' abstraction
  - Abstract device level details

- Ensure efficient utilization of accelerators
  - Capture data and execution pipelining
  - Tackle various level of parallelism available

# Unifying Instruction Set Architecture

- Unify at the instruction set level
  - Decouple fabric properties. From unification

- Enable innovation above and below ISA
  - Above (OS): global properties, e.g., security and efficiency
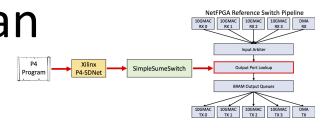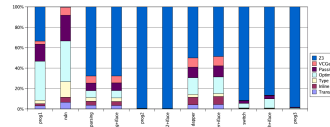  - Below (Device drivers): local properties, e.g., security and efficiency

# Preliminary Evidence of ISA for Accelerators

- Rich body of work on P4
  - Porting application to use P4
  - Developing P4 compilers for Accelerators

| KeyValue Store | Streaming App. | Machine Learning |
|---|---|---|
| P4 Compiler | P4 Compiler | P4 Compiler |
| GPU | FPGA | ASIC |

# Primer on Programmable Data Planes

P4
Code

- Composed of multiple components:
  - Parser (FSM)
  - Control flow graph (weighted DAG)
  - Match/Action Tables

**Language Limitations**
- **No loops, pointers**



```
table ipv4_lpm {
    reads {
        ipv4.dstAddr : lpm;
    }
    actions {
        set_next_hop; drop;
    }
}
```

# Primer on Programmable Data Planes

**Language Limitations**
- **No loops, pointers**

**Easy to Analyze**

**Hardware Limitations**
- **Limited resources**
- **Limited processing**

**Predictable Performance**

P4 Code

- Composed of multiple components:
  - Parser (FSM)
  - Control flow graph (weighted DAG)
  - Match/Action Tables



Programmable Parser

M    A

Switch Config

# Preliminary Evidence of ISA for Accelerators



- Rich body of work on P4
  - Porting application to use P4
  - Developing P4 compilers for Accelerators

- Initial benefits
  - Practical verification → limited language
  - Demonstrated acceleration → manual partitioning
  - Simply development → More intuitive than Verilog

- Many open challenges
  - Security: device level and fabric level
  - Global management
  - Virtualization
  - Efficient utilization of hardware

# The Road towards AcceleratorFabric

- Classic but Open Management Challenges
  - Security: device level and fabric level
  - Global resource management
  - Device virtualization
  - Efficient utilization of hardware
  - Program diagnosis and debugging



Consolidated Accelerator Deployments

Disaggregated Accelerator Deployments



Harmony [HotOS'19]



P4Visor [CoNEXT'18]

# Life @ the Edge

- The edge is real!! (and it gets you burgers)

- Heterogeneity dominates the edge
  - Let's lean into this heterogeneity
  - New systems techniques and principles

- Think holistically about management
  - Performance, efficiency AND correctness, security

# Questions?