# The Data-Driven Business Challenge
## From Reactive to Proactive

# Big and Slow or Small and Fast

## Batch Layer

ETL Tools → Change Log → Batch Processing → Data Lake (View 1, View 2) → Reports

**Data Sources**

## Real-time Layer

Stream → Stream Processing → Serving (In-Memory, NoSQL) → Real-time Dashboard

**Too slow**
- Big data but slow
- Not up to date
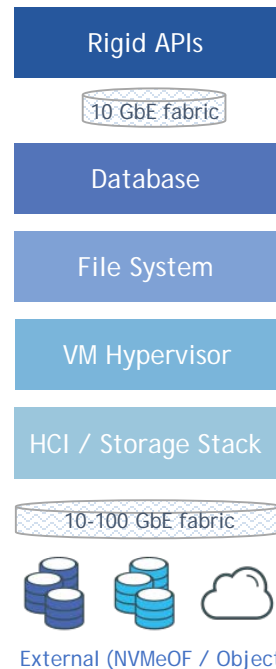- Complex

**OR**

**Limited context**
- Small amounts of data
- Expensive
- Lacks context

3

iguazio

# Traditional Approach, DB over File over Flash

**Traditional Layered Approach**

Rigid APIs

10 GbE fabric

Database

File System

VM Hypervisor

HCI / Storage Stack

10-100 GbE fabric

External (NVMeOF / Object)

- Slow
- Complex
- Expensive

## Ext3 classification illustrated

echo 'Hello, world!' >> foo; sync

- READ_10 (lba     231495 len     8 grp  9) <=4KB
- WRITE_10 (lba    231495 len     8 grp  9) <=4KB
- WRITE_10 (lba  16519223 len     8 grp  8) Journal
- WRITE_10 (lba  16519231 len     8 grp  8) Journal
- WRITE_10 (lba  16519239 len     8 grp  8) Journal
- WRITE_10 (lba  16519247 len     8 grp  8) Journal
- WRITE_10 (lba      8279 len     8 grp  5) Inode

7 I/Os (28KB) to write 13 bytes

– Metadata accounts for most of the overhead

Michael Mesnier, Jason Akers, Feng Chen, Tian Luo. Differentiated Storage Services.
23rd ACM Symposium on Operating Systems Principles (SOSP). October 2011.
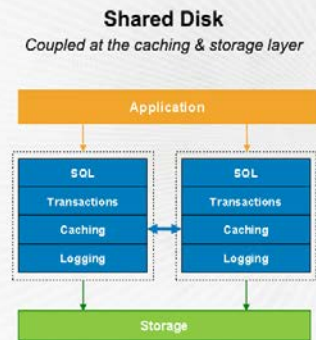
10/5/2016                    Intel Labs                    19                    intel

**For every file IOs conducted by the DB**
(Record, Redo/Undo, Metadata, ..)

4

iguazio

# New Cloud Databases Are Built to Scale Ops & Capacity

## Current DB architectures are monolithic

**Shared Disk**
*Coupled at the caching & storage layer*

Application

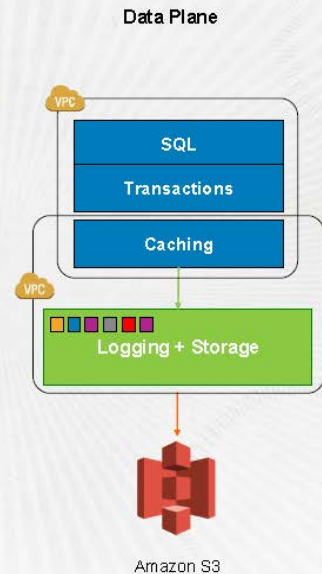| SQL | SQL |
| --- | --- |
| Transactions | Transactions |
| Caching | Caching |
| Logging | Logging |

Storage

Even when you scale it out, you're still replicating the same stack

## Amazon Aurora

Data Plane

- Service-oriented architecture applied to the database

- **Moved the logging and storage layer into a multi-tenant, scale-out database-optimized storage service**

- Integrated with other AWS services like Amazon EC2, Amazon VPC, Amazon DynamoDB, Amazon SWF, Amazon Route 53 for control plane operations

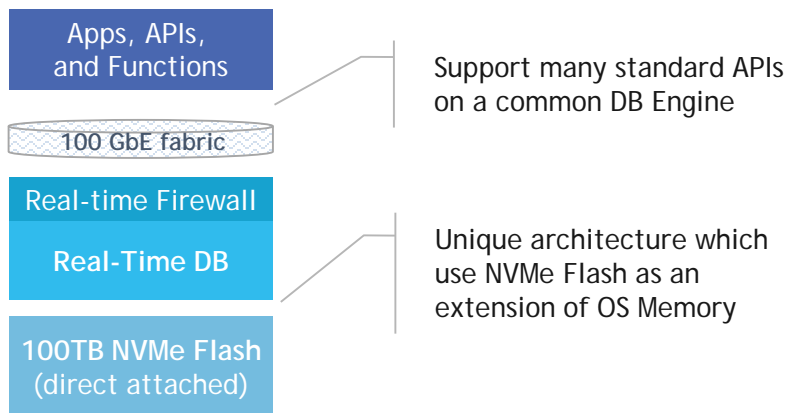- Integrated with Amazon S3 for continuous backup

VPC

| SQL |
| --- |
| Transactions |

Caching

VPC

Logging + Storage

Amazon S3

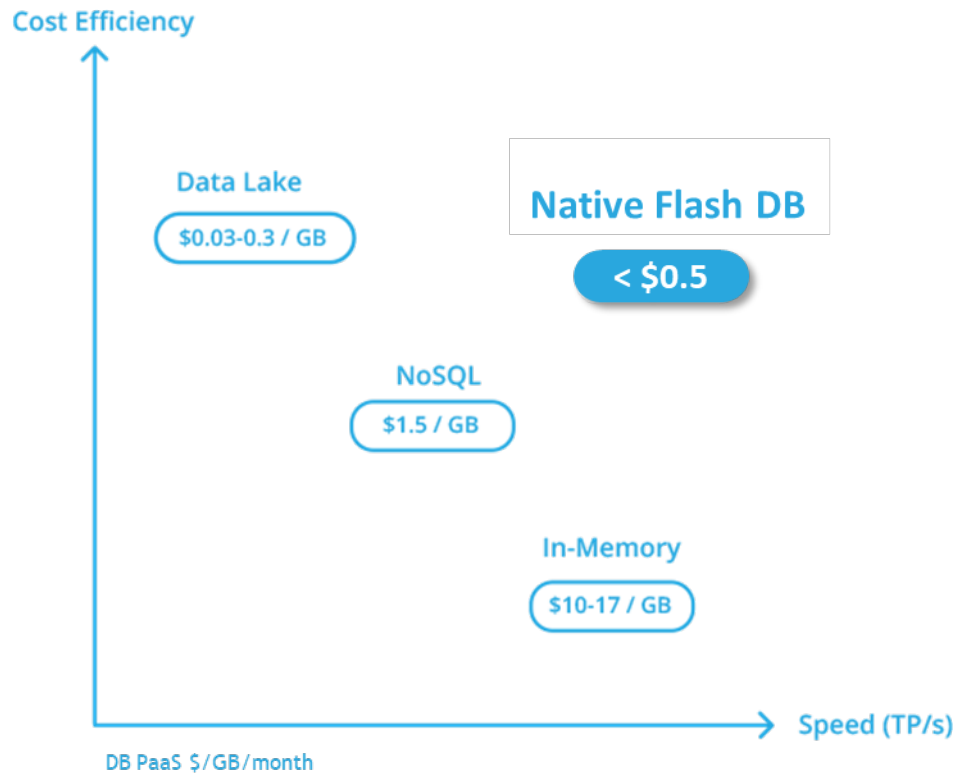**API & Transaction**

**Distributed Processing & Cache**

**Capacity (Object)**

**Decouple access, processing, and capacity and eliminate storage serialization**
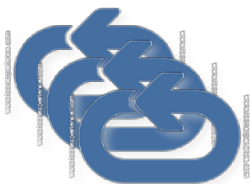
iguazio

# Breaking The Volume and Velocity Barrier

Apps, APIs, and Functions

100 GbE fabric

Real-time Firewall

Real-Time DB

100TB NVMe Flash (direct attached)

Support many standard APIs on a common DB Engine

Unique architecture which use NVMe Flash as an extension of OS Memory

Re-engineer the stack to deliver memory speed with Flash density

**Cost Efficiency**

**Data Lake**
$0.03-0.3 / GB

**Native Flash DB**
< $0.5

**NoSQL**
$1.5 / GB

**In-Memory**
$10-17 / GB

Speed (TP/s)

DB PaaS $/GB/month

iguazio

# Breaking Performance Barriers – Design Principles

*Never blocking, never locking, 100% parallelism*
*Latency optimized, QoS aware, data scheduler*
*Lockless, preempt less memory management*
*True scale out through parallelism*

*Zero processing wastes*
*CPU cache optimization and prediction*
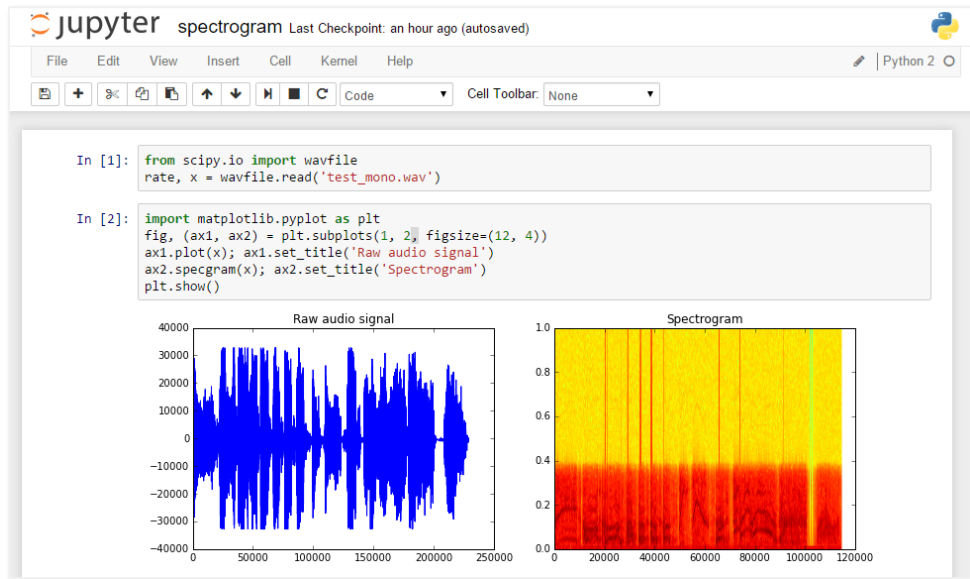*E2E zero buffer data flow (NIC to Disk, accelio)*
*Complete OS bypass*

*HW awareness*
*RDMA, NVMe (3DXP)*
*Vector processing operations*
*IRQ balancing and throttling*

Ok, any other challenges on the way to real-time AI ?
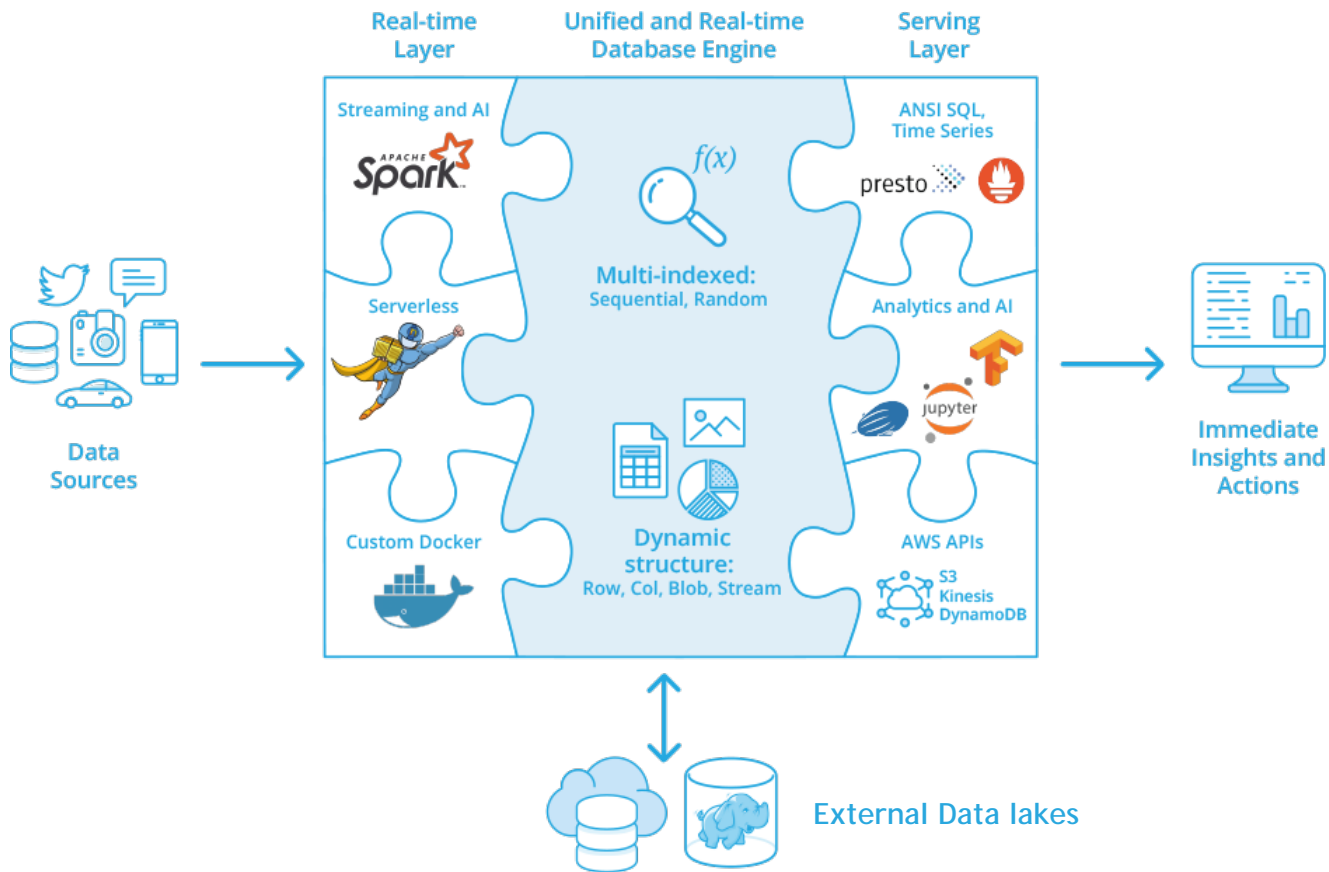
# 90% of AI Today



Build feature vectors using batch and CSVs

Inspect, Improve

**How do we form complex feature vectors in real-time?**
**How do we visualize or act on the results in real-time?**
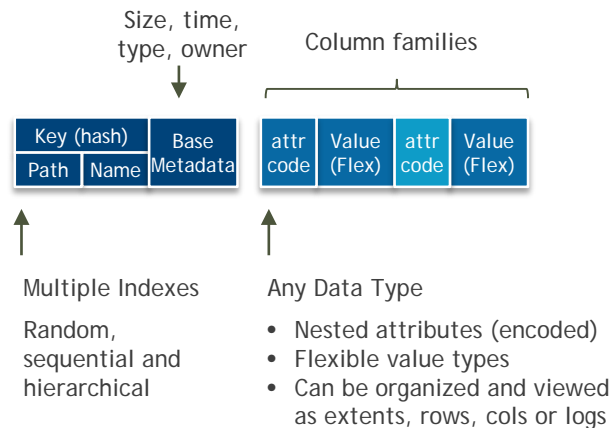
# Moving to Continuous Ingest + AI + Serve Flow



Real-time Layer — Unified and Real-time Database Engine — Serving Layer

Data Sources → Real-time Layer (Streaming and AI: Apache Spark; Serverless; Custom Docker) → Unified and Real-time Database Engine (Multi-indexed: Sequential, Random; Dynamic structure: Row, Col, Blob, Stream; f(x)) → Serving Layer (ANSI SQL, Time Series: presto; Analytics and AI: Jupyter; AWS APIs: S3 Kinesis DynamoDB) → Immediate Insights and Actions

External Data lakes

iguazio

# From Silos and ETLs to All-in-one DBs

## Traditional: Unique Model Per Store

## Multi-Model Store

| | Index | Metadata & data |
|---|---|---|
| File | Dir (tree) / Name (tree) | Simple Metadata / Data Extents |
| Object | Key (Random hash) | Extended Metadata / Data Blob (immutable) |
| K/V | Key (Random hash) | Simple Metadata / Value Blob (immutable) |
| Table (fixed) | Key (Seq tree) | Value (typed) / Value (typed) / Value (typed) |
| Document | Key (Seq tree) | attr / Value (Flex) / attr / Value (Flex) |
| Stream | Topic / Shard /Metric | ts / Value Blob / ts / Value Blob |

Size, time, type, owner

Column families

Key (hash) / Base Metadata

Path / Name / Base Metadata

attr code / Value (Flex) / attr code / Value (Flex)

**Multiple Indexes**

Random, sequential and hierarchical

**Any Data Type**

- Nested attributes (encoded)
- Flexible value types
- Can be organized and viewed as extents, rows, cols or logs

Independent tiering logic for indexes, metadata and data

iguazio

# Time Series Data Example

**Raw time series sample data**

Thousands of samples

```
{
  "metric" : "rx-bandwidth",
  "device": "xyz",
  "port": 1,
  "mac": "0123456...",
  "rack": "A13",
  "value": 77,
  "time": 1524690488000
}
```

Labels

Data

Ingest/compress
**In real-time**

**Optimized TSDB Layout** (per unique metric)

```
{
  "__name__" : "rx-bandwidth",
  "device": "xyz",
  "port": 1,
  "mac": "0123456...",
  "rack": "A13",

  "_v_count": [...],
  "_v_sum": [...],
  ...

  "_v0": <compressed blob>,
  "_v1": <compressed blob>,
  ...
}
```

Filter based on labels

Pre-aggregation arrays:
(to accelerate queries)

T/V chunks with 10:1
Gorilla compression

| Real-time Consistency | 50 : 1 Compression | 10–100x Faster Queries |
|---|---|---|

iguazio

# Serverless, The New Stored Procedure

## Traditional Dev and Ops Model

- **Write code + local testing**
- Build code and Docker image
- CI/CD pipeline
- Add logging and monitoring
- Harden security
- Provision servers + OS
- Handle data/event feed
- Handle failures/auto-scaling
- Handle rolling upgrades
- Configuration management

**80%**

## "Serverless" Development Model

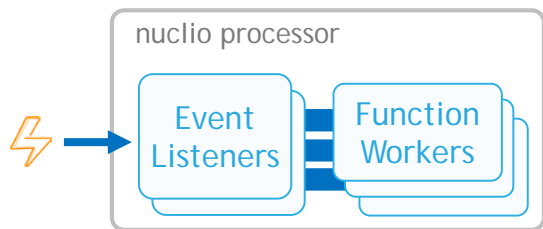- **Write code + local testing**
- **Provide spec, push deploy**

1. Automated by the serverless platform

2. Pay for what you use

iguazio

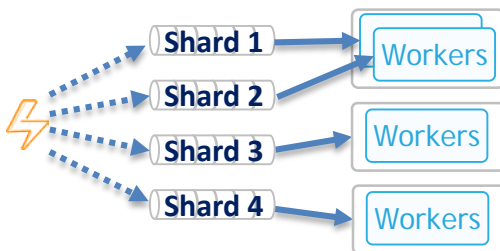# Addressing Serverless Limitations With Nuclio
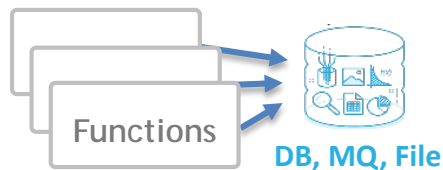
### Performance



- Non-blocking, parallel
- Zero copy, buffer reuse
- Up to 400K events/sec/proc

### Streaming and Batch



- Auto-rebalance, checkpoints
- Any source: Kafka, NATS, Kinesis, event-hub, iguazio, pub/sub, RabbitMQ, Cron
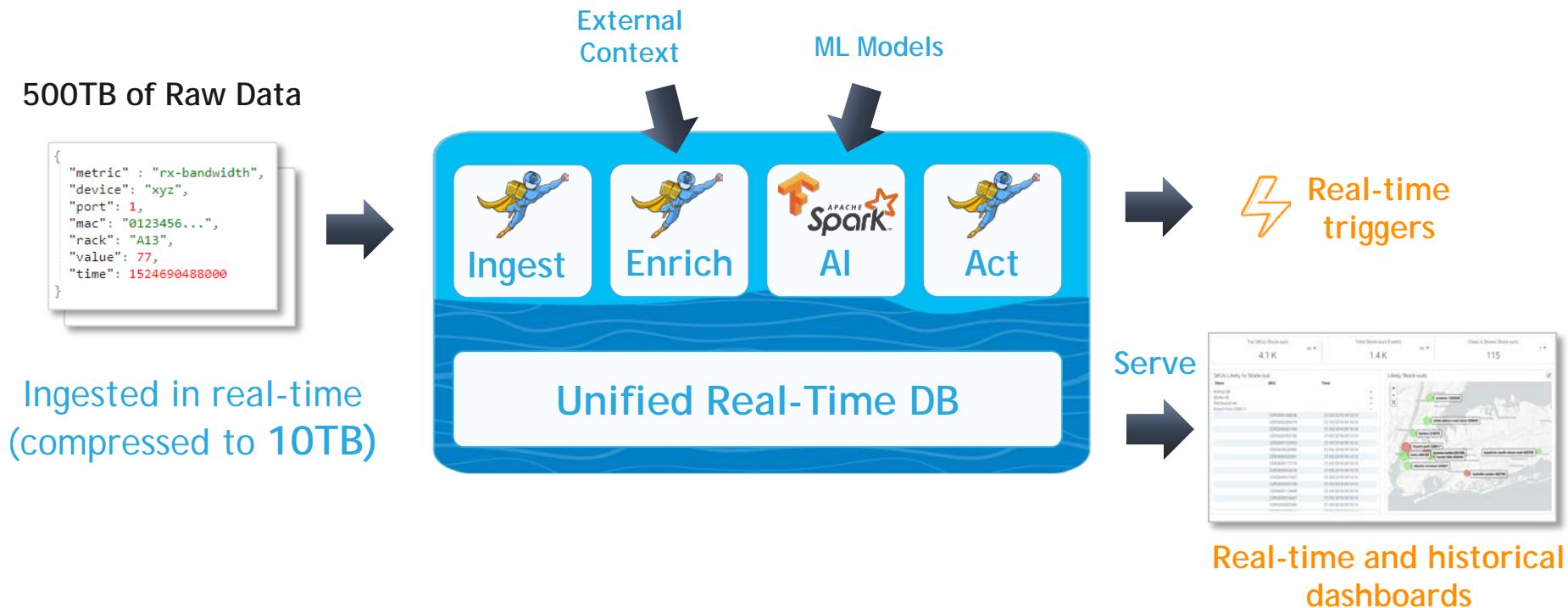
### Statefulness



- Data bindings
- Shared volumes
- Context cache

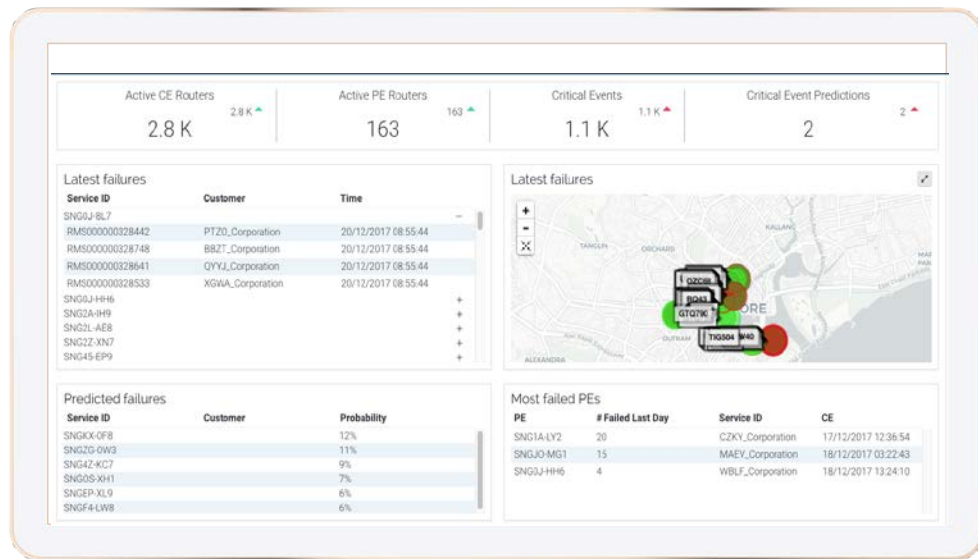**Serverless for compute and data intensive tasks
100x faster than AWS Lambda !**
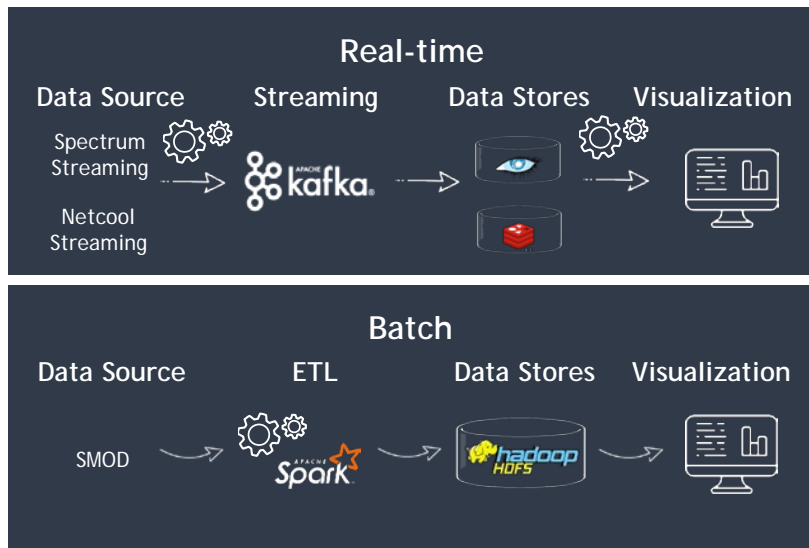
# Delivering Intelligent Decisions in Real-Time

**500TB of Raw Data**

```
{
  "metric" : "rx-bandwidth",
  "device": "xyz",
  "port": 1,
  "mac": "0123456...",
  "rack": "A13",
  "value": 77,
  "time": 1524690488000
}
```

Ingested in real-time
(compressed to **10TB**)

**External Context**

**ML Models**

Ingest

Enrich

AI

Act

**Unified Real-Time DB**

**Real-time triggers**

**Serve**

**Real-time and historical dashboards**

iguazio

# Cyber and Network Ops

## A leading telco needs to predict network behavior in real-time:

- Processing high message throughput from multiple streams at the rate of > 50K events/sec

- Cross correlating with historical and external data in real-time

- AI predictions/inferencing conducted on live data
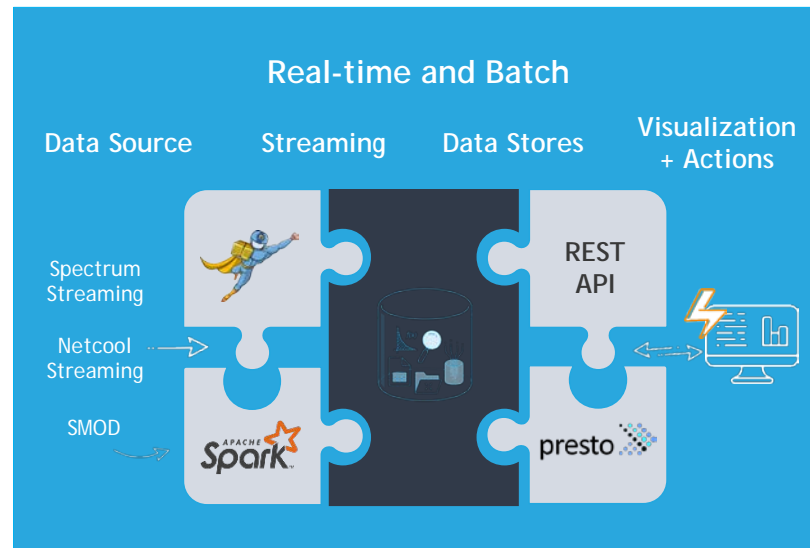
- Small footprint to fit network locations

# Build and Operationalize Proactive Systems Faster
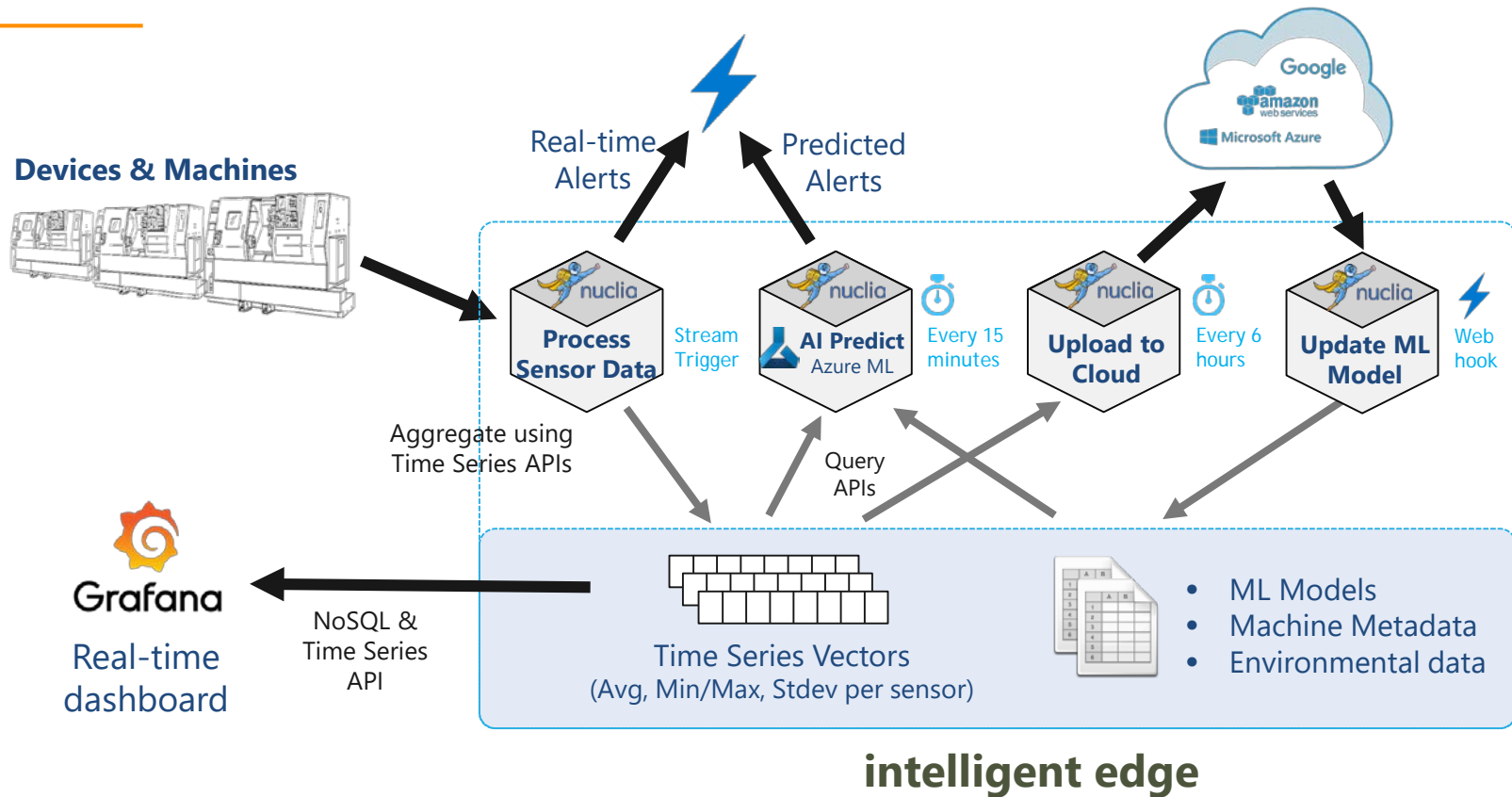
## Traditional



- Complex, skill gaps, slow to productize
- No single view of ops, real-time, history
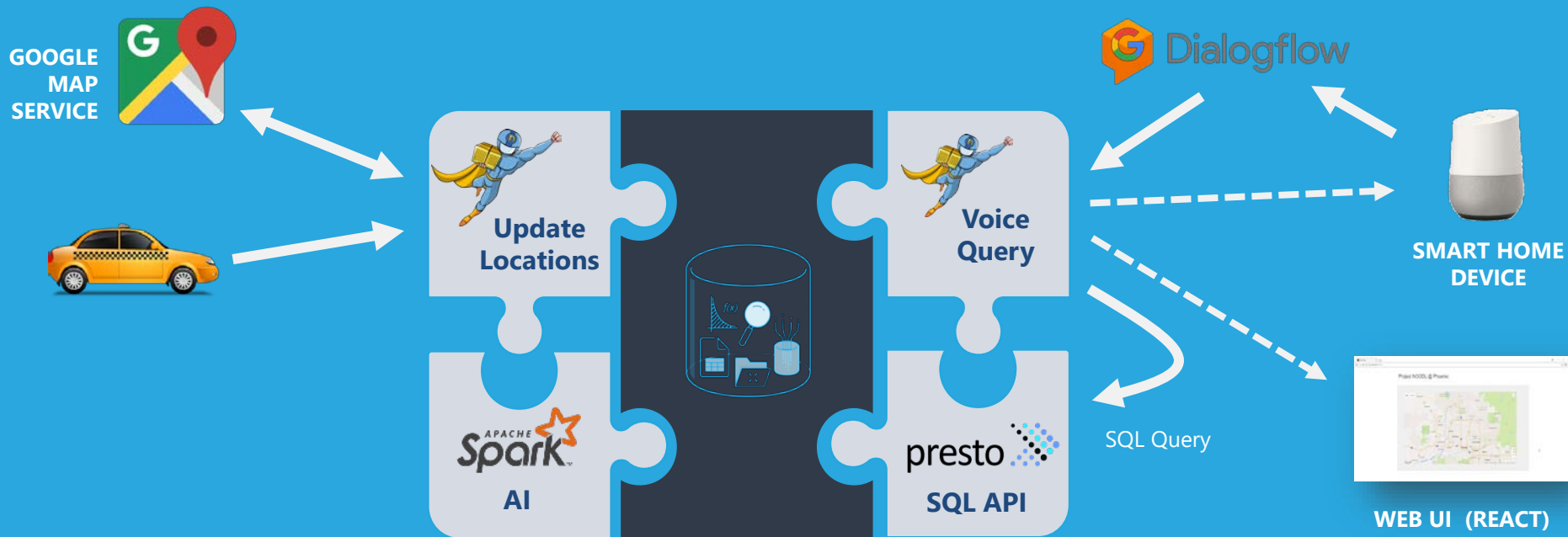- Reactive (no actions)

## Continuous Analytics



- Simple, just a few weeks to a working app
- Unified view across ALL data
- AI driven, proactive

# Predictive Maintenance Based on Real-time + Historical + Ops Data

Demo: Voice Driven Real-Time Analytics

# Summary

## Build continuous, data-driven and proactive apps

- Deliver real-time analytics on fresh, historical and operational data

- Optimize Flash usage to deliver in-memory speed at much lower costs

- Create a unified data layer for stream processing, AI and serving

- Adopt cloud-native and serverless approaches to gain agility

iguazio

# Thank You

info@iguazio.com  |  www.iguazio.com

iguazio