# The **RISE** Lab: Real-time Intelligent Secure Execution

Krste Asanovic, Ali Ghodsi, Ken Goldberg, **Joey Gonzalez**, **Joe Hellerstein**, Michael Jordan, Randy Katz, Dave Patterson, **Raluca Ada Popa**, **Ion Stoica**, …

# Berkeley's lab tradition



- Working for 5-6 years on a new major problem

- Bringing faculty from different areas

# AMPLab (2010—2016)

Created popular open-source big data analytics:
  Spark, Mesos, Tachyon..

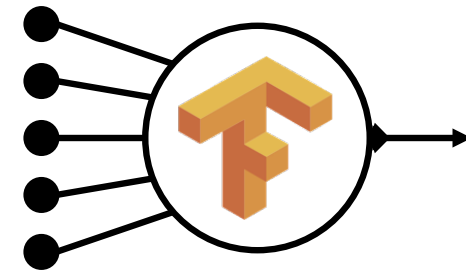AMPLab is coming to an end (December 2016)

**What is the next vision?**
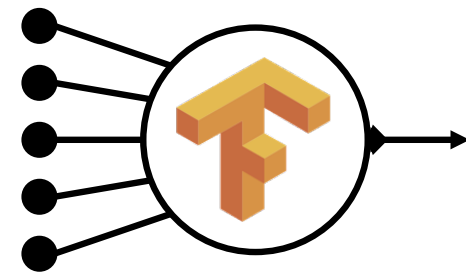
**RISE** Lab

From live data to real-time decisions
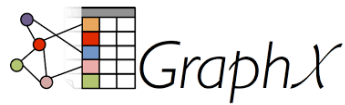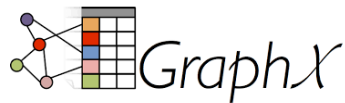
**AMP** Lab

From batch data to advanced analytics

Big Data → Training → Big Model

amplab

**Learning**

Big Data

Training

Big Model

?

**Learning**

Big Data → Training → Big Model → Conference Papers

**Learning**

Big Data

Training

Big Model

Conference Papers

Dashboards and Reports

**Learning**

Big Data → Training → Big Model

Conference Papers

Dashboards and Reports

**Drive Actions**

**Learning**

**Drive Actions**

Big Data

Training

Big Model

Hi, I'm Cortana.

# Learning

# Inference

Big Data

Training

Big Model

**Learning** | **Inference**

Big Data → Training → Big Model ← Query ← Application

Big Model → Decision → Application

**Learning**

**Inference**

Big Data

Training

Big Model

Query

Decision

Application

Often **overlooked**
Timescale: ~10 milliseconds
*An area of focus in the RISELab*

# why is **Inference** challenging?

Need to render **low latency** (< 10ms) predictions for **complex**

**Models**     **Queries**     **Features**



Top K

SELECT * FROM
users JOIN items,
click_logs, pages
WHERE …

under **heavy load** with system **failures**.

# Inference is moving beyond the cloud



Augmented Reality

Home Security

Home Automation

Mobile Assistants

Self Driving Cars

Personal Robotics

# **Inference** is moving beyond the cloud



## **Opportunities**

➤ Reduce latency and improve privacy

➤ Address network partitions

## **Research Challenges**

➤ Minimize **power consumption**

➤ **Limited hardware** & long life-cycles

➤ Develop new **hybrid models** to leverage cloud and devices

# Robust **Inference** is critical

Self "*Parking*" Cars

Self "*Driving*" Cars

Chat AIs

# Learning                          # Inference

Big Data $\xrightarrow{\text{Training}}$ Big Model

Big Model $\xleftarrow{\text{Query}}$ Application

Big Model $\xrightarrow{\text{Decision}}$ Application

Application $\xrightarrow{\text{Feedback}}$ Big Data

**Learning**

**Inference**

Big Data

Training

Decision

Application

**Timescale:** hours to weeks
*Often re-run training*
*Another area of focus in RISE*

Feedback

# Why is **Closing the Loop** challenging?



**Implicit** and **Delayed** Feedback

Self Reinforcing **Feedback Loops**

World Changes **at varying rates**

# Learning

# Inference

**Adaptive (~1 seconds)**

Training

**Responsive (~10ms)**

Query

Decision

Big Model

Feedback

Application

**Learning**
Adaptive
(~1 seconds)

**Inference**
Responsive
(~10ms)

Secure

# Intelligence in **Sensitive Contexts**

| AR/VR Systems | Home Monitoring | Voice Technologies | Medical Imaging |
|---|---|---|---|



## Protect the **data**, the **model**, and the **query**

# Protect the **data**, the **model**, and the **query**

## High-Value **Data** *is Sensitive*



- Medical Info.
- Home video
- Finance

## **Models** capture **value** in data

- Core Asset
- Sensitive



## **Queries** can be as sensitive as the data

**Adaptive**

**Responsive**

**Secure**

# Goal

Real-time decisions

*decide in ms*

on live data

with strong security

# Goal

Real-time decisions

*decide in ms*

on live data

*the current state of the environment*

with strong security

# Goal

Real-time decisions

> *decide in ms*

on live data

> *the current state of the environment*

with strong security

> *privacy, confidentiality, integrity*

| Challenges | RISE Lab |
|---|---|
| Automated decisions on live data are hard | Real-time, sophisticated decisions that guarantee worst-case behavior on noisy and unforseen live data |
| Poor security: exploits are daily occurrences | Ensure privacy and integrity without impacting functionality |
| One-off solutions, expensive and slow to build | General platform: Secure Real-time Decision Stack |

# Example: Zero-time defense

**Problem**: zero-day attacks can compromise millions of sites in seconds

**Solution**: analyze network flows to detect attacks and patch sites/software in real-time

- **Intermediate data**: create attack model
- **Decision**: detect attack, patch

| Quality | sophisticated, accurate, robust |
|---|---|
| Latency | sec (decision) / sec (update) |
| Security | privacy (encourage users to share logs), integrity |

# Example: "Fleet" driving

**Problem**: suboptimal driving decisions
**Solution**: collect & leverage info from
other cars and drivers in <span style="color:orange">real-time</span>

- **Intermediate data**: automatically annotate maps, actions of other drivers
- **Decision**: avoid obstacles, congestions



| Quality | sophisticated, accurate, noise tolerant |
|---|---|
| **Performance** | sec (decision) / sec (update) |
| **Security** | privacy, data integrity |

33

# Example: Infectious disease discovery

**Problem**: infectious diseases spread quickly (Zika), may need quarantine (Ebola)

**Solutions**: real-time DNA seq. & analysis to identify pathogens

- Rapid analysis to trace evolution, source
- 100x faster → 100x people tested
  - **Intermediate data**: evolution, spread, symptoms
  - **Decision**: quarantine or not, diagnosis

| Quality | sophisticated, accurate |
|---------|-------------------------|
| Latency | min (decision) / hour (update) |
| Security | privacy, integrity |



MinION Nanopore
(Dr. Charles Chiu UCSF using it to identify Zika virus)

34

| Applications | Quality | Latency | | Security |
| --- | --- | --- | --- | --- |
| | | Decision | Update | |
| Zero-time defense | sophisticated, accurate, robust | sec | sec | privacy/confidentiality, integrity |
| Parking assistant | sophisticated, robust | sec | sec | |
| Disease discovery | sophisticated, accurate | sec/min | hours | |
| IoT (smart buildings) | sophisticated, robust | sec | min/hour | |
| Earthquake warning | sophisticated, accurate, robust | ms | min | |
| Chip manufacturing | sophisticated, accurate, robust | sec/min | min | |
| Fraud detection | sophisticated, accurate | ms | min | |
| "Fleet" driving | sophisticated, accurate, robust | sec | sec | |
| Virtual companion | sophisticated, robust | sec | min/hour | |
| Video QoS at scale | sophisticated | ms/sec | min | |

# Challenges

# RISE Lab

| | |
|---|---|
| Automated decisions on live data are hard | Real-time, sophisticated decisions that guarantee worst-case behavior on noisy and unforseen live data |
| Poor security: exploits are daily occurrences | Ensure privacy and integrity without impacting functionality |
| One-off solutions, expensive and slow to build | General platform: **Secure Real-time Decision Stack** |

36

# A bird's eye view



decision

query

push decision

end-point info

Secure Real-time Decision Stack

end-points
(e.g., users, devices)

# What exists today?

# Pull decisions

**Example: recommendation system**



**recommendation**

Service Layer
(WebApp, Cassandra, MySQL)

**per-user recommendations**

request

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark, Cloud Dataflow)

**collaborative filtering**

Storage (HDFS)

**user purchase and rating history**

39

# Pull decisions

**Problem: cannot have sophisticated decisions**



Cannot pre-compute all decisions
Cannot touch much data

decision

ms

request

Service Layer
(WebApp, Cassandra, MySQL)

min/day

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark,
Cloud Dataflow)

Storage (HDFS)

40

# Pull decisions

Cannot pre-compute all decisions
Cannot touch much data

decision

ms

request

Service Layer
(WebApp, Cassandra, MySQL)

min/day

insecure

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark,
Cloud Dataflow)

Storage (HDFS)

secure
(encryption at rest)

41

# Solution scorecard

| Solution | Decision Quality | Latency | | Security |
| --- | --- | --- | --- | --- |
| | | Decision | Update | |
| Pull decisions | simple | ms | min/day | weak |

# Pull decisions: contextual decisions



decision

request

Prediction Serving Layer
(WebApp, Velox, on-line MWT)

light-weight, ensemble
and correction models,
policies

Ingestion
(Kafka)

ML
(Hadoop, Spark,
Cloud Dataflow)

off-line, large
models

Storage (HDFS)

# Solution scorecard

| Solution | Decision Quality | Latency | | Security |
| --- | --- | --- | --- | --- |
| | | Decision | Update | |
| Pull decisions | simple | ms | min/day | weak |
| Pull decisions: prediction service | sophisticated, specialized | ms | min* | weak |

*light-weight, ensemble + correction models, policies

# Push Decisions
**Example: anomaly detection**



alert

Streaming
(Storm, Spark, IBM Streams)

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark)

Storage (HDFS)

**outlier
detection**

**support vector
machine (SVM)**

**machine logs**

# Push Decisions
**Example: anomaly detection**



alert
ms/sec

Streaming
(Storm, Spark, IBM Streams)

hour

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark)

Storage (HDFS)

**outlier detection**

**support vector machine (SVM)**

**machine logs**

# Solution scorecard

| Solution | Decision Quality | Latency | | Security |
| --- | --- | --- | --- | --- |
| | | Decision | Update | |
| Pull decisions | simple | ms | min/day | weak |
| Push decisions | simple | ms/sec | hour | weak |

*light-weight, ensemble + correction models, policies

# Security tools

- Computation on encrypted data
- Hardware enclaves

# State-of-the-art security solutions:
## Computation on encrypted data

# State-of-the-art security solutions:
## Computation on encrypted data



decision

request

Web App

Key-value store
(Cassandra, MySQL)

Ingestion
(Kafka)

ETL/ML
(Hadoop, Spark)

user info

Storage (HDFS)

SQL (CryptDB), ML classification, web serving (Mylar), etc.

Ensure security, but relatively simple algorithms

# State-of-the-art security solutions: Hardware enclaves



Intel SGX, ARM Trustzone, RISC V (e.g., Haven, VC3)

decision

request

Confid

Confid

user info

Web App

Key-value stpre (Cassandra, MySQL)

Ingestion (Kafka)

ETL/ML (Hadoop, Spark)

TCB, side-channel leakage

51

# Solution scorecard

| Solution | Decision Quality | Latency | | Security |
|---|---|---|---|---|
| | | Decision | Update | |
| Pull decisions | simple | ms | min/day | weak |
| Push decisions | simple | ms/sec | hour | weak |
| State-of-the-art security | simple | ms | min/hour | strong |

*light-weight, ensemble + correction models, policies

# Solution scorecard: **RISE**

| Solution | Decision Quality | Latency | | Security |
| --- | --- | --- | --- | --- |
| | | Decision | Update | |
| Pull decisions | simple | ms | min/day | weak |
| Push decisions | simple | ms/sec | hour | weak |
| Security | simple | ms | min/hour | strong |
| RISE | sophisticated, accurate, robust | ms | sec | strong |

*light-weight, ensemble + correction models, policies

53

# Research areas

Systems: Spark-like functionality with 100x lower response time, and 1000x higher job throughput

Machine Learning:
- On-line ML algorithms
- Robust algorithms: handle noisy data, guarantee worst-case behavior

Security: achieve privacy, confidentiality, and integrity without impacting performance

54

# Early Projects



SQL    ML    Graph

Opaque

query optimization

( o-filter ) ( o-groupby ) ( o-join )

Catalyst

Spark Execution



Analytics & Vis

Reference Data

Wrangling

Data Quality

Catalog & Discovery

Reproducibility

Parsing & Featurization

Model Serving

COMMON GROUND      ABOVEGROUND API TO APPLICATIONS

CONTEXT MODEL

UNDERGROUND API TO SERVICES

Scavenging and Ingestion    Versioned Storage    Search & Query    ID & Auth    Scheduling & Workflow

http://ground-context.org

ground

# Research area: Systems

☐ on-going work    ⬚ future work

**Apache Spark**

**Lattice Flow**

**Ray**

**Clipper**

...

shim layer    shim layer    shim layer    shim layer

**RISE μkernel**

scheduler    optimizer    ...

in-memory obj

unified model
(rich experience w/ both models)

**IndexedRDDs**

GPU/ASICs algos

**Drizzle**

sharded driver,
in-memory processing,
per-core NIC, HBM

system-state store

- Support task-graph & BSP execution models
- Support fine grain updates
- Support heterogeneous hardware
- Millisecond level parallel jobs
- Handle 10K-100K jobs/sec
- Ability to faithfully replay jobs

56

# THE MEANING AND VALUE OF DATA DEPENDS ON CONTEXT

**Application context**
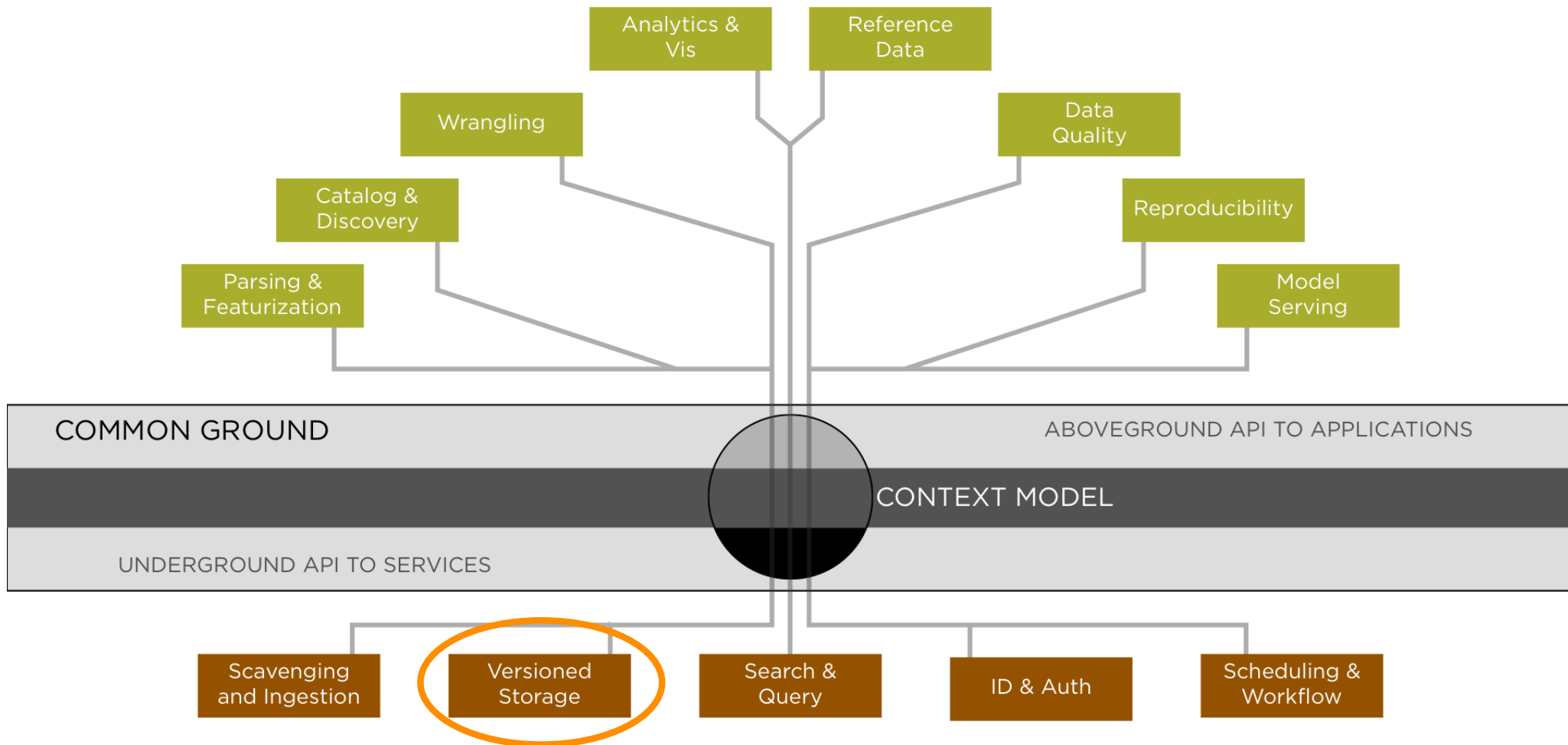Multiple data models

because truth
is subjective

**Behavioral context**
Logs and Lineage

because behavior
determines meaning

**Historical context**
Immutable versions
for code and data

because
things change

A broader context for big data

ground

http://ground-context.org

# LatticeFlow and Bedrock (working names)

**Driving Hypotheses for LatticeFlow**

- A core programming API for both real-time and scale
- Everything is (async) data: event dispatch, real-time data streams
- Coordination Avoidance: lattices + async dataflow = no locks/barriers

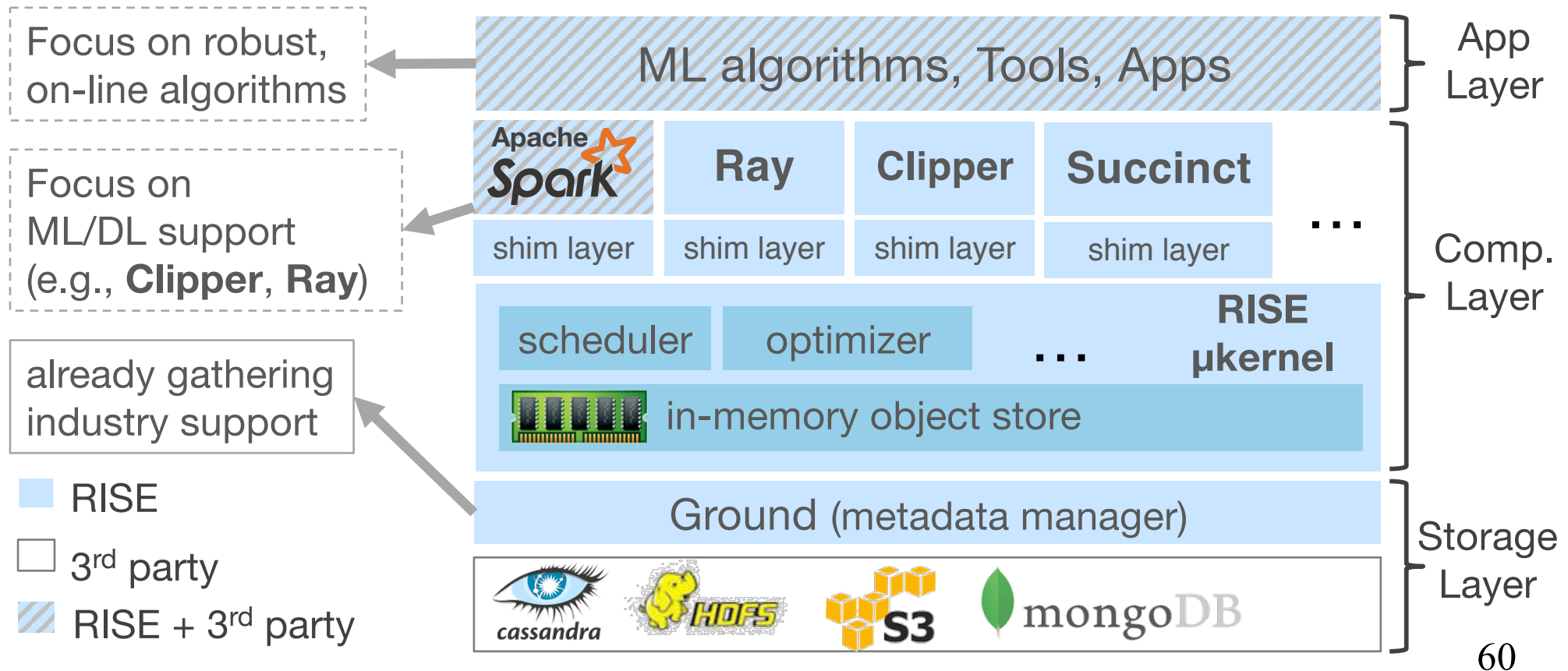**Initial proof point this fall: LatticeKVS**

MultiScale, MultiConsistency key-value store

- Same system beats Redis on one node, Cassandra on scale-out…
- …while providing family of rich consistency and transactional isolation options
- …with lean codebase, derived from a core LatticeFlow library in C++

**Prototype toward "Project Bedrock"**

- Immutable, never-forget versioned storage under ground

# Research area: Systems & ML

Focus on robust, on-line algorithms

ML algorithms, Tools, Apps

App Layer

Focus on ML/DL support (e.g., **Clipper**, **Ray**)

Apache **Spark**  **Ray**  **Clipper**  **Succinct**  . . .

shim layer | shim layer | shim layer | shim layer

Comp. Layer

already gathering industry support

scheduler  optimizer  . . .  **RISE µkernel**

in-memory object store

RISE

3rd party

RISE + 3rd party

Ground (metadata manager)

Storage Layer

cassandra  HDFS  S3  mongoDB

60

# Research area: ML

**Robust optimization methods:**

- noise tolerant and parallelizable

**Handle uncertainty**:

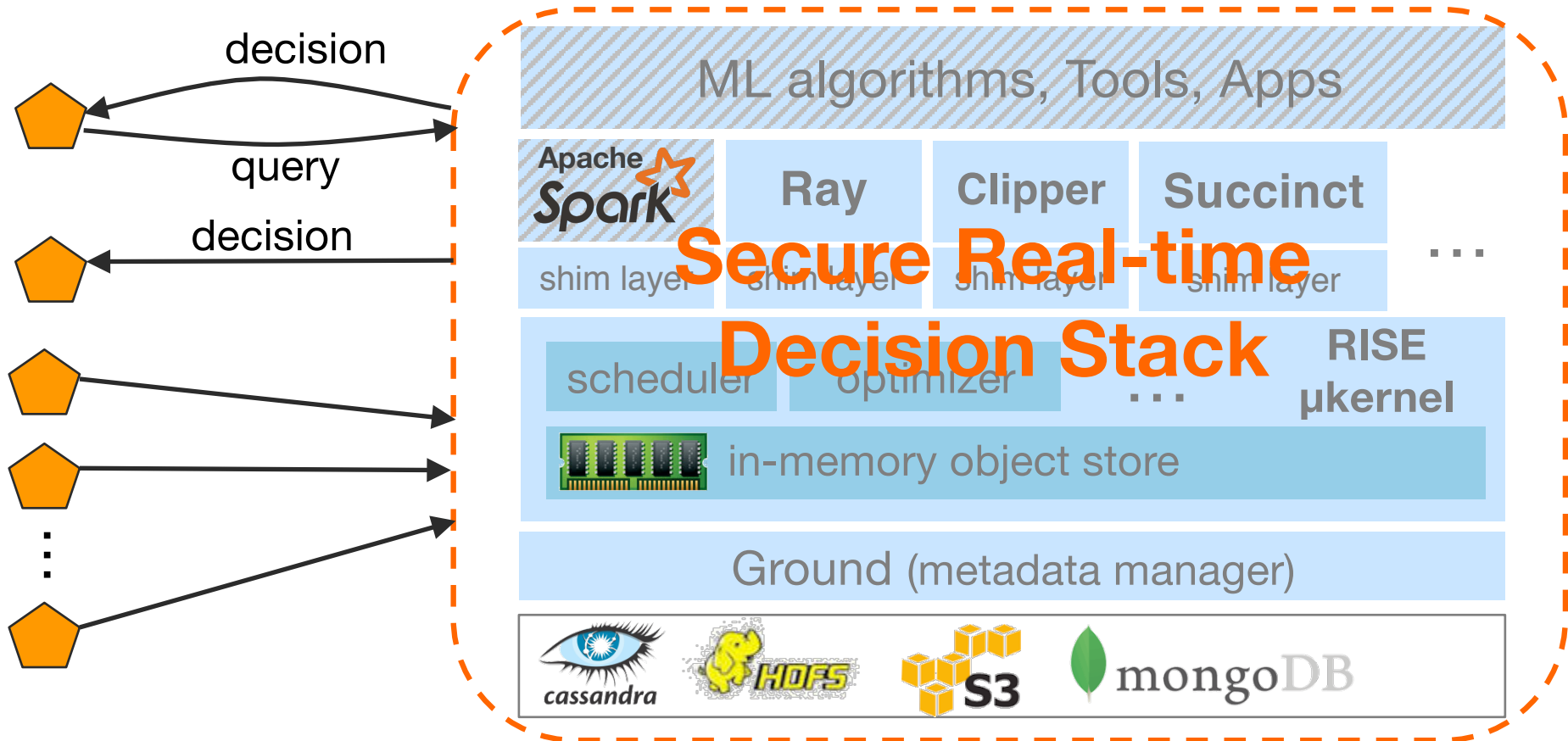- robust control techniques to handle unforseen real world situations

Quantify decision accuracy:

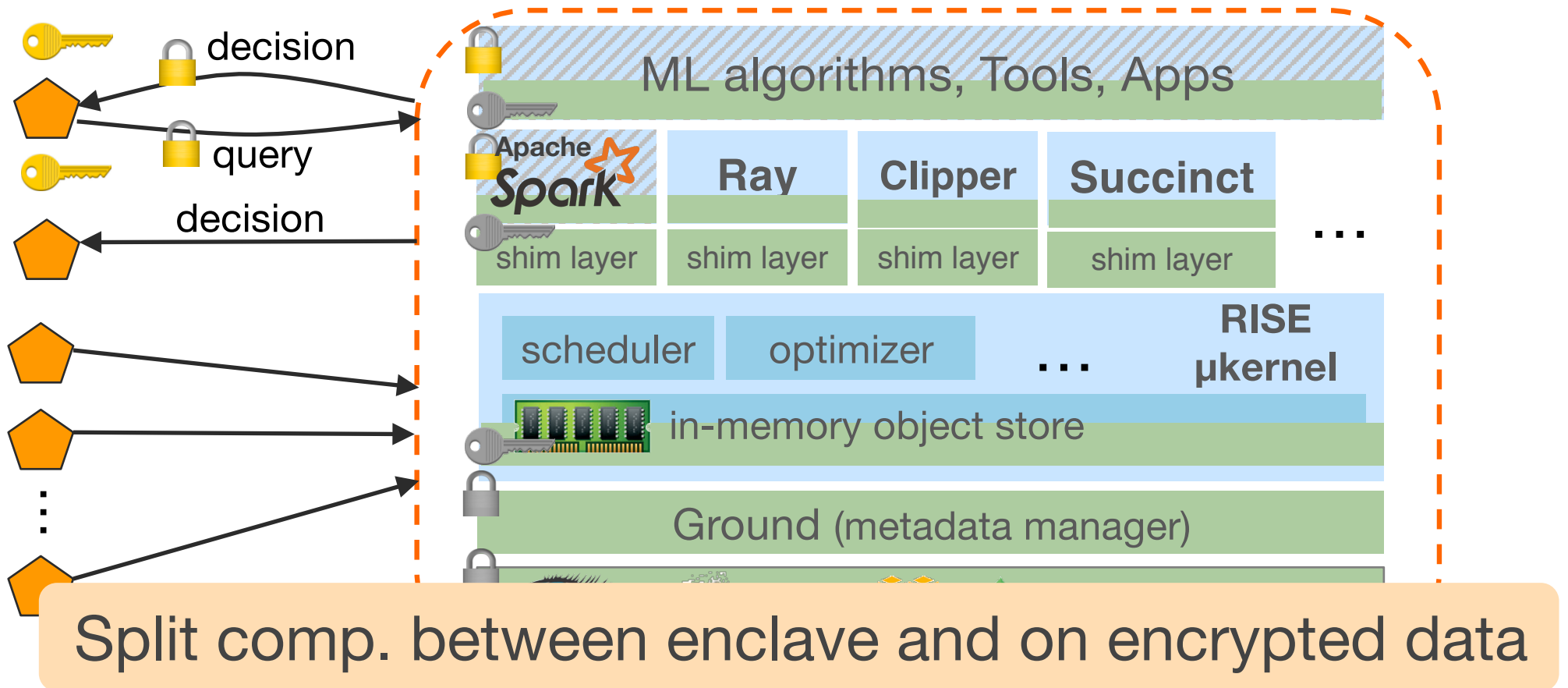- confidence and credible intervals on outputs of ML systems

On-line ML algorithms

- handle time-varying models with high performance and reliability

# End-to-end system

# Research area: Security



decision

query

decision

ML algorithms, Tools, Apps

Apache Spark | Ray | Clipper | Succinct | ...

shim layer | shim layer | shim layer | shim layer

scheduler | optimizer | ... | **RISE μkernel**

in-memory object store

Ground (metadata manager)

**Split comp. between enclave and on encrypted data**

# Research area: Security

ML (decisions) on encrypted data

**Opaque**: oblivious & encrypted analytics

New hardware support (RISC V, HBM)

Fine grained access control

**Arx**: queries on encrypted data

ML algorithms, Tools, Apps

Apache **Spark** | **Ray** | **Clipper** | **Succinct** | ...

shim layer | shim layer | shim layer | shim layer

scheduler | optimizer | ... | **RISE μkernel**

in-memory object store

Ground (metadata manager)

cassandra | HDFS | S3 | mongoDB

# Why now, why us?

Latency

Quality

Security

## Hardware Trends

**Processing**
- RISC V*, GPUs, ASICs
- Built in security support (enclaves)

**Storage & networking**
- 3D Xpoint

**Next gen rack designs**
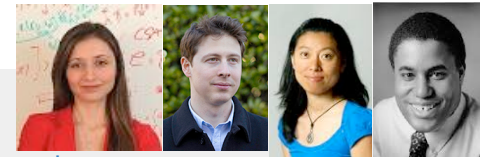- FireBox*

## Systems

- Berkeley's BDAS*

## Machine Learning

- Robust algorithms*
- On-line ML* and DL

## Security

- Code verification tools*
- Comp. on encrypted data*

# Summary

**Goal**: develop Secure Real-time Decision Stack,
an open source platform, tools and algorithms
for real-time decisions on live data with strong security

Five year project, similar to AMPLab
We are uniquely positioned to tackle this challenge

Looking to partner with companies