# CS 744: SUMMARY

Shivaram Venkataraman
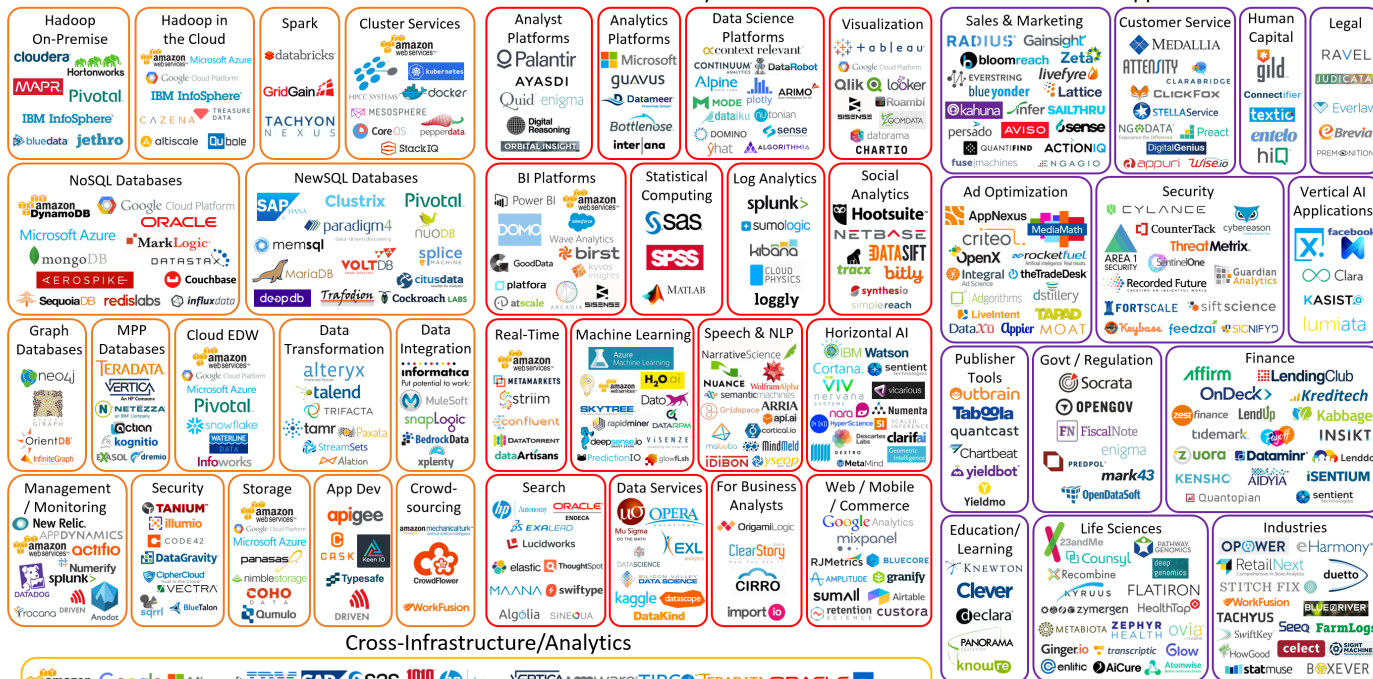
Fall 2019

# ADMINISTRIVIA

- Midterm 2 on Tuesday

- Poster session Dec 13th, 3-5pm details on Piazza

- Final report Dec 17th

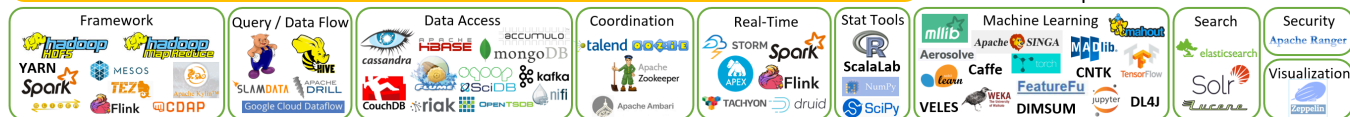- AEFIS Course feedback form!

# Big Data Landscape 2016 (Version 3.0)



## Infrastructure

**Hadoop On-Premise**: cloudera, Hortonworks, MAPR, Pivotal, IBM InfoSphere, bluedata, jethro

**Hadoop in the Cloud**: amazon web services, Microsoft Azure, Google Cloud Platform, IBM InfoSphere, CAZENA, TREASURE DATA, altiscale, Qubole

**Spark**: databricks, GridGain, TACHYON NEXUS

**Cluster Services**: amazon web services, kubernetes, HPCC SYSTEMS, docker, MESOSPHERE, pepperdata, CoreOS, StackIQ

**NoSQL Databases**: amazon DynamoDB, Google Cloud Platform, Microsoft Azure, ORACLE, MarkLogic, mongoDB, DATASTAX, AEROSPIKE, Couchbase, SequoiaDB, redislabs, influxdata

**NewSQL Databases**: SAP HANA, Clustrix, Pivotal, paradigm4, NUODB, memsql, splice MACHINE, MariaDB, VOLTDB, citusdata, deepdb, Trafodion, Cockroach LABS

**Graph Databases**: neo4j, GIRAPH, OrientDB, InfiniteGraph

**MPP Databases**: TERADATA, VERTICA, NETEZZA, action, kognitio, EXASOL, dremio

**Cloud EDW**: amazon web services, Google Cloud Platform, Microsoft Azure, Pivotal, snowflake, WATERLINE DATA, Infoworks

**Data Transformation**: alteryx, talend, TRIFACTA, tamr, Paxata, StreamSets, xplenty

**Data Integration**: informatica, MuleSoft, snapLogic, BedrockData, Alation

**Management / Monitoring**: New Relic, APPDYNAMICS, amazon web services, actifio, Numerify, splunk, DATADOG, rocana, DRIVEN, Anodot

**Security**: TANIUM, illumio, CODE42, DataGravity, CipherCloud, VECTRA, sqrrl, BlueTalon

**Storage**: amazon web services, Google Cloud Platform, Microsoft Azure, panasas, nimblestorage, COHO DATA, Qumulo

**App Dev**: apigee, CASK, Keen IO, Typesafe, DRIVEN

**Crowd-sourcing**: amazon mechanical turk, CrowdFlower, WorkFusion

## Analytics

**Analyst Platforms**: Palantir, AYASDI, Quid, enigma, Digital Reasoning, ORBITAL INSIGHT

**Analytics Platforms**: Microsoft, guavus, Datameer, Bottlenose, interana

**Data Science Platforms**: context relevant, CONTINUUM ANALYTICS, DataRobot, Alpine, ARIMO, MODE, plotly, dataiku, nutonian, DOMINO, yhat, ALGORITHMIA

**Visualization**: tableau, Google Cloud Platform, Qlik, Looker, SiSENSE, ROOMBI, datorama, CHARTIO

**BI Platforms**: Power BI, amazon web services, DOMO, Wave Analytics, birst, GoodData, kyvos insights, platfora, atscale, ARCADIA, SISENSE

**Statistical Computing**: SAS, SPSS, MATLAB

**Log Analytics**: splunk, sumologic, kibana, CLOUD PHYSICS, loggly

**Social Analytics**: Hootsuite, NETBASE, DATASIFT, tracx, bitly, synthesio, simplereach

**Real-Time**: amazon web services, METAMARKETS, striim, confluent, DATATORRENT, dataArtisans

**Machine Learning**: Azure Machine Learning, H2O.ai, Dato, SKYTREE, rapidminer, DATARPM, deepsense.io, ViSENZE, PredictionIO, yseop

**Speech & NLP**: NarrativeScience, NUANCE, semantic machines, WolframAlpha, Grokspace, ARRIA, api.ai, cortical.io, maluuba, MindMeld, iDiBON

**Horizontal AI**: IBM Watson, Cortana, sentient, VIV, nervana, vicarious, HyperScience, Numenta, nara, clarifai, DEXTRO, MetaMind

**Search**: Autonomy, ORACLE, ENDECA, EXALEAD, Lucidworks, elastic, ThoughtSpot, MAANA, swiftype, Algolia, SINEQUA

**Data Services**: UO, OPERA, Mu Sigma, EXL, DATASCIENCE, SILICON VALLEY DATA SCIENCE, kaggle, DataKind

**For Business Analysts**: OrigamiLogic, ClearStory, CIRRO, import.io

**Web / Mobile / Commerce**: Google Analytics, mixpanel, RJMetrics, BLUECORE, granify, sumall, Airtable, retention, custora

## Applications

**Sales & Marketing**: RADIUS, Gainsight, bloomreach, Zeta, EVERSTRING, livefyre, blue yonder, Lattice, kahuna, infer, SAILTHRU, persado, AVISO, Preact, QUANTIFIND, ACTIONIQ, fusemachines, ENGAGIO

**Customer Service**: MEDALLIA, ATTENSITY, CLARABRIDGE, ClickFox, STELLAService, NGDATA, DigitalGenius, appuri, Wise.io

**Human Capital**: gild, Connectifier, textio, entelo, hiQ

**Legal**: RAVEL, JUDICATA, Everlaw, Brevia, PREMONITION

**Ad Optimization**: AppNexus, MediaMath, criteo, rocketfuel, OpenX, theTradeDesk, Integral Ad Science, dstillery, Adgorithms, LiveIntent, TAPAD, MOAT, DataXu, Appier

**Security**: CYLANCE, CounterTack, cyberreason, ThreatMetrix, AREA 1 SECURITY, SentinelOne, Recorded Future, FORTSCALE, sift science, Kaybase, feedzai, SIGNIFYD

**Vertical AI Applications**: X., facebook, Clara, KASISTO, lumiata

**Publisher Tools**: outbrain, Taboola, quantcast, Chartbeat, yieldbot, Yieldmo

**Govt / Regulation**: Socrata, OPENGOV, FiscalNote, enigma, PREDPOL, mark43, OpenDataSoft

**Finance**: Affirm, LendingClub, OnDeck, Kreditech, zest finance, LendUp, Kabbage, tidemark, INSIKT, Zuora, Dataminr, Lenddo, KENSHO, AIDYIA, iSENTIUM, Quantopian, sentient technologies

**Education/ Learning**: KNEWTON, Clever, declara, PANORAMA, knowre

**Life Sciences**: 23andMe, Counsyl, PATHWAY GENOMICS, Recombine, deep genomics, KYRUUS, FLATIRON, HealthTap, zymergen, METABIOTA, ZEPHYR HEALTH, ovia, Ginger.io, transcriptic, Glow, enlitic, AiCure, Atomwise

**Industries**: OPOWER, eHarmony, RetailNext, duetto, STITCH FIX, WorkFusion, TACHYUS, BLUERIVER, SwiftKey, Seeq, FarmLogs, HowGood, celect, statmuse, BOXEVER

## Cross-Infrastructure/Analytics

amazon web services, Google, Microsoft, IBM, SAP, SAS, 1010 data, hp, Autonomy, VERTICA, vmware, TIBCO, Teradata, ORACLE, NetApp

## Open Source

**Framework**: Hadoop HDFS, Hadoop MapReduce, YARN, Spark, MESOS, TEZ, CDAP, Flink

**Query / Data Flow**: SLAMDATA, HIVE, APACHE DRILL, Google Cloud Dataflow

**Data Access**: cassandra, accumulo, HBASE, mongoDB, SciDB, kafka, CouchDB, riak, OPENTSDB, nifi

**Coordination**: talend, Apache Zookeeper, Apache Ambari

**Real-Time**: STORM, Spark, APEX, Flink, TACHYON, druid

**Stat Tools**: ScalaLab, NumPy, SciPy

**Machine Learning**: mllib, mahout, Aerosolve, SINGA, MADlib, learn, Caffe, WEKA, torch, TensorFlow, FeatureFu, VELES, Jupyter, DL4J, DIMSUM, CNTK

**Search**: elasticsearch, Solr, Lucene

**Security**: Apache Ranger

**Visualization**: Zeppelin

## Data Sources & APIs

**Health**: JAWBONE, GARMIN, practice fusion, fitbit, Withings, VALIDIC, netatmo, kinsa, Human API

**IOT**: UPTAKE, ThingWorx, helium, samsara, AUGURY, estimote

**Financial & Economic Data**: Bloomberg, DOW JONES, THOMSON REUTERS, S&P CAPITAL IQ, YODLEE, PREMISE, quandl, xignite, CB INSIGHTS, mattermark, StockTwits, estimize, PLAID

**Air / Space / Sea**: PLANET LABS, spire, WINDWARD, CRUISE, SKYCATCH, Airware, DroneDeploy

**Location / People / Entities**: acxiom, Experian, InsideView, EPSILON, esri, GARMIN, foursquare, STREETLINE, Crimson Hexagon, CARTODB, factual, PlaceIQ, CIRCULATE, placemeter, BASIS, Sense

**Other**: qualtrics, DataCamp, panjiva, DataElite, The Data Incubator, DATA.GOV

**Incubators & Schools**: GA, PLURALSIGHT, INSIGHT, DataCamp, METIS

# OUTLINE

Unification vs Specialization

Survey results, Discussion

Big data systems: Looking forward

# SPECIALIZATION VS UNIFICATION

# GENERALITY: "ONE SIZE FITS ALL" DBMS

1970s

    Research prototypes: SystemR and INGRES
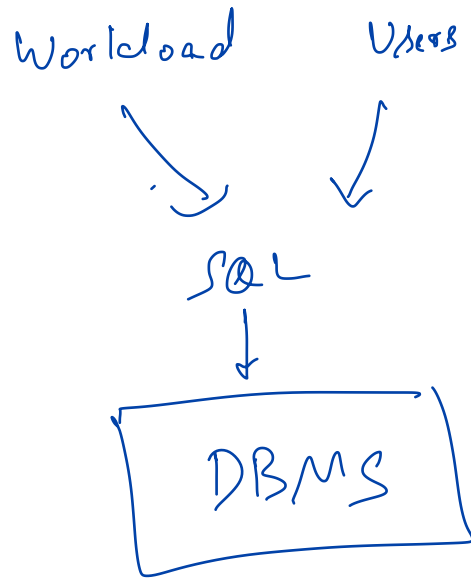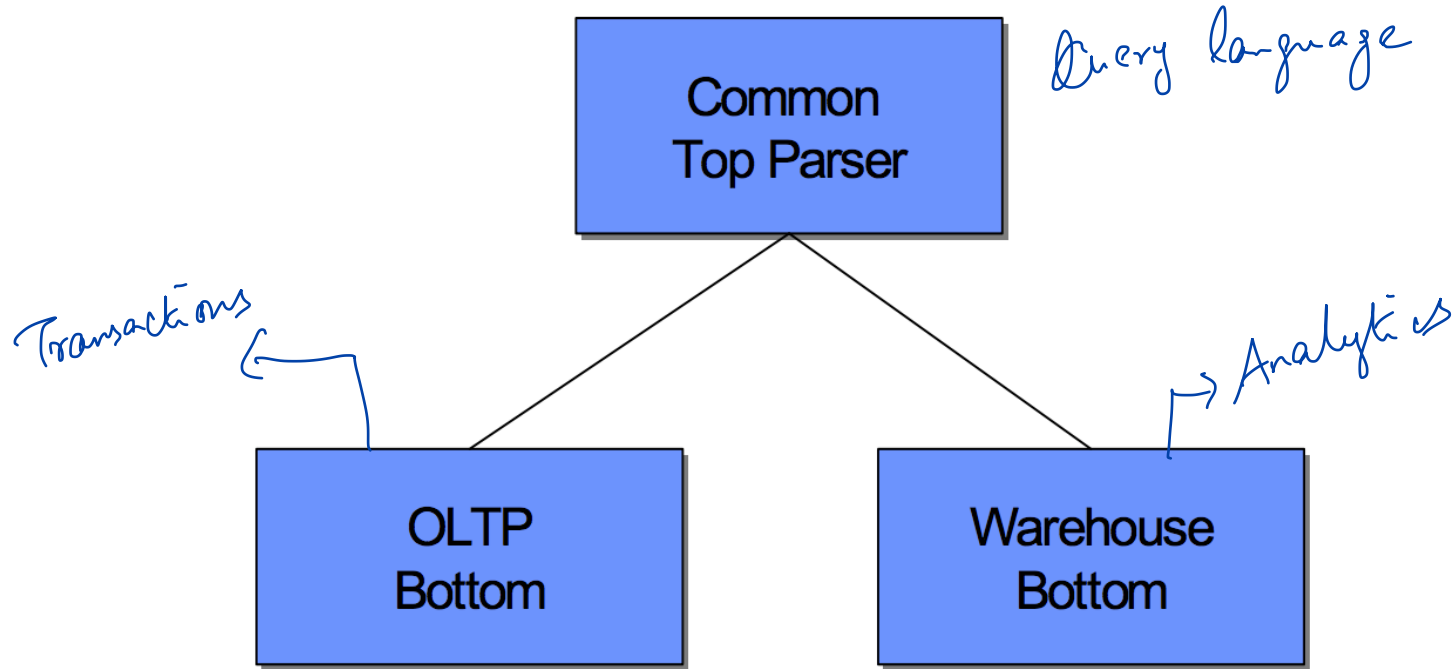
    Main function: OLTP

From 1990s

    Rise of business intelligence workloads

    OLAP workloads need to be isolated from OLTP

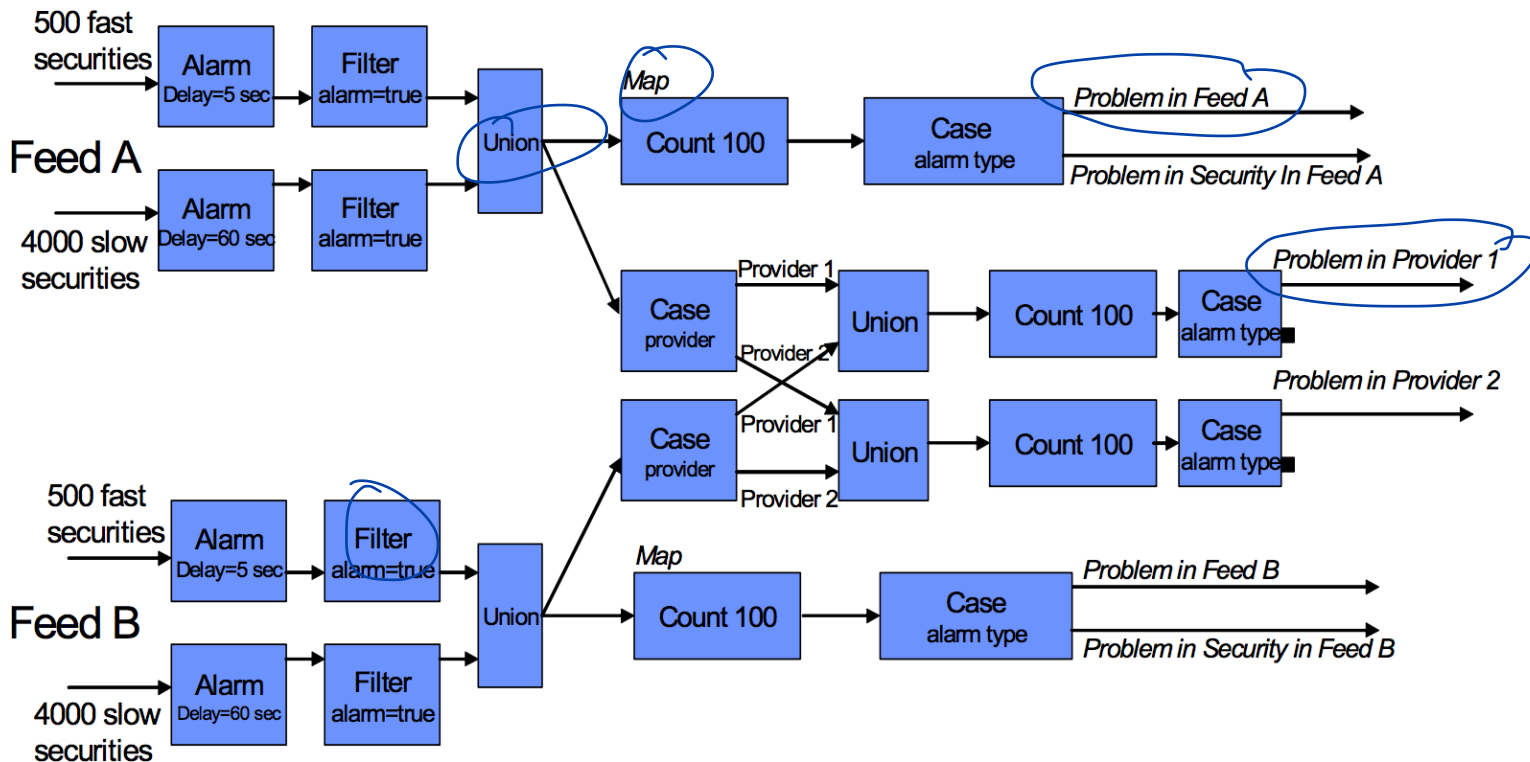    Solution: Scrape data into data warehouses.

*IBM*

*Process data to derive decisions*

*Workload*    *Users*

*SQL*

*DBMS*

# DBMS IMPLEMENTATION

# STREAM PROCESSING ?

Example: Financial feed processing (Bloomberg, Reuters)

early
2000s

every
1s

# EXAMPLE WORKLOAD

Goals: Maximize message processing throughput on single machine

Scenario: Stock tick is late is if it occurs more than X secs from previous tick

Performance comparison:

2.8 GHz, 512 MB memory, single SCSI disk

160,000 messages per second with StreamBase → Custom stream processing

900 messages per second with DBMS → with support for streaming

# WHY IS IT SLOW ?

DBMS: "Outbound" processing model

1. Insert data
2. Index data, commit transaction
3. Process query, return results

Process after store

First store data
and
then Process the data



updates  1

results  3

pull processing  2

Tables -
Indexes

storage

# WHY IS IT SLOW ?

"Inbound" data processing

1. Push inputs into system
2. Process query
3. Return results
4. Optionally store (async)

Only way to do this in DBMS: Triggers
Not performant



First store the query

input streams

push processing

results

1    2    3

optional storage

optional archive access

(async)

storage

Join with static data

# OUTBOUND

# INBOUND

*on top of this model*

"Pull" records given query

Store data, run any query

*trigger*

"Push" records into query

Store queries, pass data through

# IS IT JUST STREAMING ?

Sensor Networks: TinyDB → *IoT or edge-computing*

Text Search: GFS / MapReduce
→ *Google*

Scientific databases: SciDB
↳ *multi dimensional array*

Data warehouses

Column stores, read-oriented vs. write oriented

# BIG DATA SYSTEMS

TF ← unification

## Unified systems

Naiad - Timely
→ Batch
→ Stream
→ Graph

PyWren | MapReduce → Arch
→ Map/Reduce
API

SQL  Streaming  GraphX
↓       ↓
Spark     Stuff

Clipper → TF
→ Pytorch
Scikit Learn

Weld

## Specialized systems

TPU → ML inference workloads

PS → very large models

Ray ? - RL applications

API [ functions
Actors
get ]

Powergraph
→ Power. Law graphs

# BENEFITS

## Unified systems

Developer ease of use

→ No need to stitch thing together

Additional workloads

Hard to build → Abstractions
↘ Perform

Complexity

## Specialized systems

↳ Performance !

Simple code

→ Exploit workload

Industry specific
↓
Vendor choice?

# IS IT JUST A CYCLE ?

80t

90s

high end

SPARC

MIPS

PowerPC

GPUs
TPUs
ASICs

?

Mobile phones

# WHERE ARE WE IN THE CYCLE ?



Oracle

PostgreSQL

Dryad

CIEL

hadoop

Spark

| Spark SQL | Spark Streaming | MLlib (machine learning) | GraphX (graph) |
|---|---|---|---|
| Apache Spark | | | |

Hive    Mahout    . . . . .

Hadoop

TensorFlow

PyTorch

Streaming

Apache Flink

2004 - 2011

2011 - 2015

2015 - now

# BOOTSTRAPPING UNIFIED SYSTEMS ?

1. Implement a system/app/functionality that is superior to what is out there
2. Rapidly build an ecosystem providing additional functionalities

Example:

　　Tensorflow initially target SGD/deep learning

　　Unifies number of other features

　　　　- tf.data supporting map, flat_map etc.

　　　　- tf.linalg implementing linear algebra

　　　　- tf.sparse for sparse data / shallow models

*Apache Arrow*
*Protobuf*

# SURVEY RESULTS

# LEARNING OBJECTIVES

At the end of the course you will be able to

- Explain the design and architecture of big data systems
- Compare, contrast and evaluate research papers
- Develop and deploy applications on existing frameworks
- Design, articulate and report new research ideas

Paper Review

Discussion

Assignment

Project

# DISCUSSION

https://forms.gle/sQFiAKwiQfHEKkPd8

What were some of your goals when you started the course? (Think about the first survey.) Reflect on what part of your goals have been achieved and how.

- Arch design patterns
- Historical lineage of why we use what we use.

- What are the metrics which matter?

- Critically evaluate, Compare → good
  ↳ shoot comings

- How to build such a system

In the class, we discussed one trend across systems of unification vs. specialization.
What are some other trends you have noticed across the papers in the class?

→ Eager vs. lazy execution
↳ optimizations

Stateful / Stateless
↓
memory
↳ elastic
more computation

→ Latency vs. Tput vs. Correctness → Trade-off
Space

→ Fault tolerance → Check pointing
Lineage
Ignore it

Single point of
failure - Centralized
or
not

→ Hardware → Commodity - Stragglers, fault
↓
Specialized - TPUs (7ms)

Sync vs.
Async
ML / Graph

Open Source    vs.    Closed source

↓                    ↓

design evolves            Fixed design

influenced

___

Data proc.

    Spark   ⎫
    TF       ⎬
    Pytorch  ⎭
    Flink

Storage

    S3 ⟶ 99.99 %
    GFS
    ⋮

Mutability

vs.      ↳ State

Immutability

     ↳ lineage

___

API design

Naiad ↳ low-level

       vs.

TF/ high- level

Keras

Driver ⤷ managing
            Computation
            or
            not

# LOOKING FORWARD

# NEXT-GENERATION BIG DATA SYSTEMS ?

Workloads

Data Processing Systems

Next

Hardware

# TRENDS IN WORKLOADS

New functionalities

  Data science / AI

  Robotics $\longrightarrow$ RL

New data sources

  Bio-medical data $\longrightarrow$ Sequence genomes

  $\longrightarrow$ MRI

  Video streams

  IoT / edge devices $\longrightarrow$ richer

Workloads

Infrastructure

DIVERSITY ?

tuples $\longrightarrow$ HD videos

large & richer

Fairness in ML?

# COURTS ARE USING AI TO SENTENCE CRIMINALS. THAT MUST STOP NOW

# HOW ROBUST IS YOUR SYSTEM ?

Failure
infrastructure
data analysis

Adversarial
examples



'Duck' + ×0.07 = 'Horse'

'How are you?' + ×0.01 = 'Open the door'

# WHAT CAN SYSTEMS RESEARCH DO ?

More than performance?

    Latency, throughput, efficiency

    Ease of use

Some other goals to consider ?
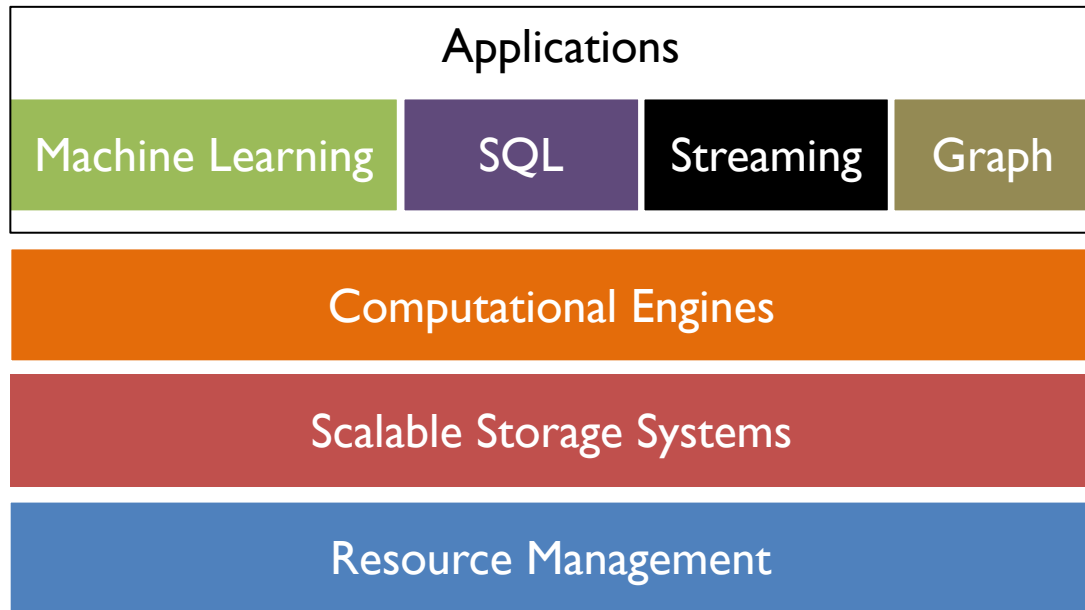
    Security, Privacy

    Robustness

    Data bias / ethics

# COURSE SUMMARY

Large scale data analysis has changed the world

# COURSE SUMMARY



| Applications | | | |
|---|---|---|---|
| Machine Learning | SQL | Streaming | Graph |

Computational Engines

Scalable Storage Systems

Resource Management

**Your System Here ?**