# Storage and Data Challenges for Production Machine Learning
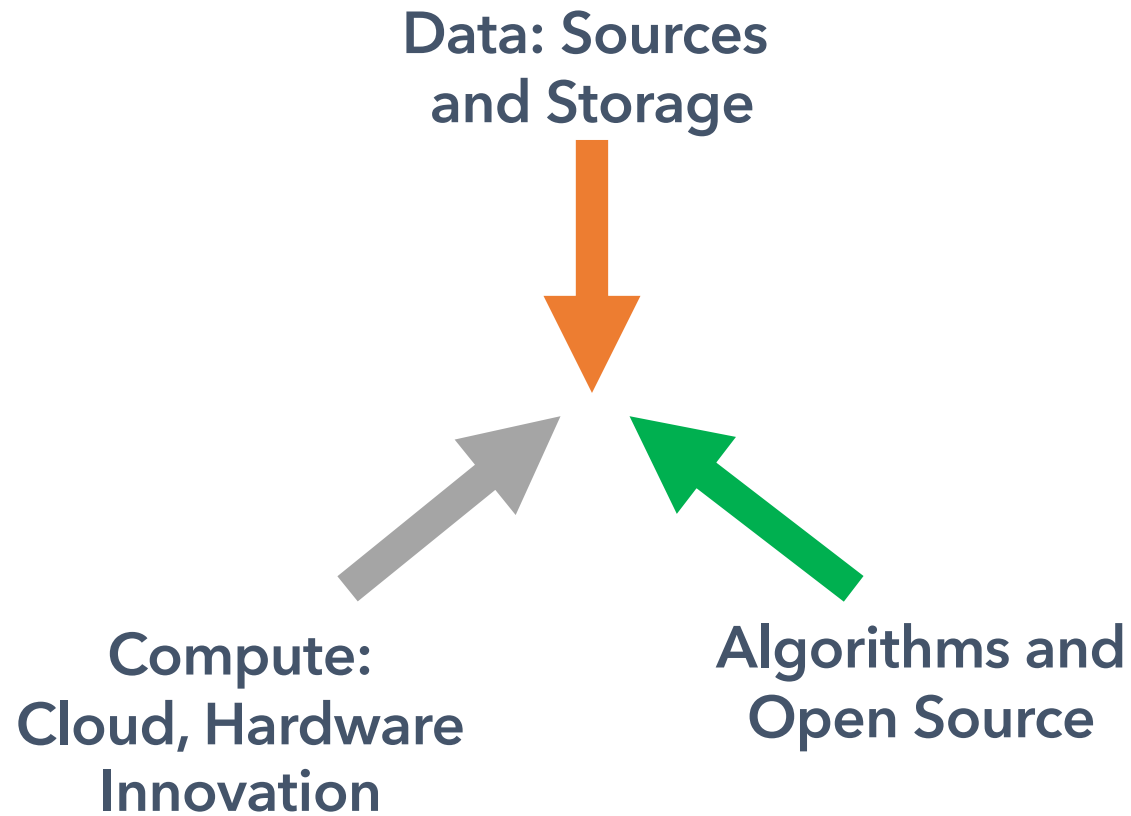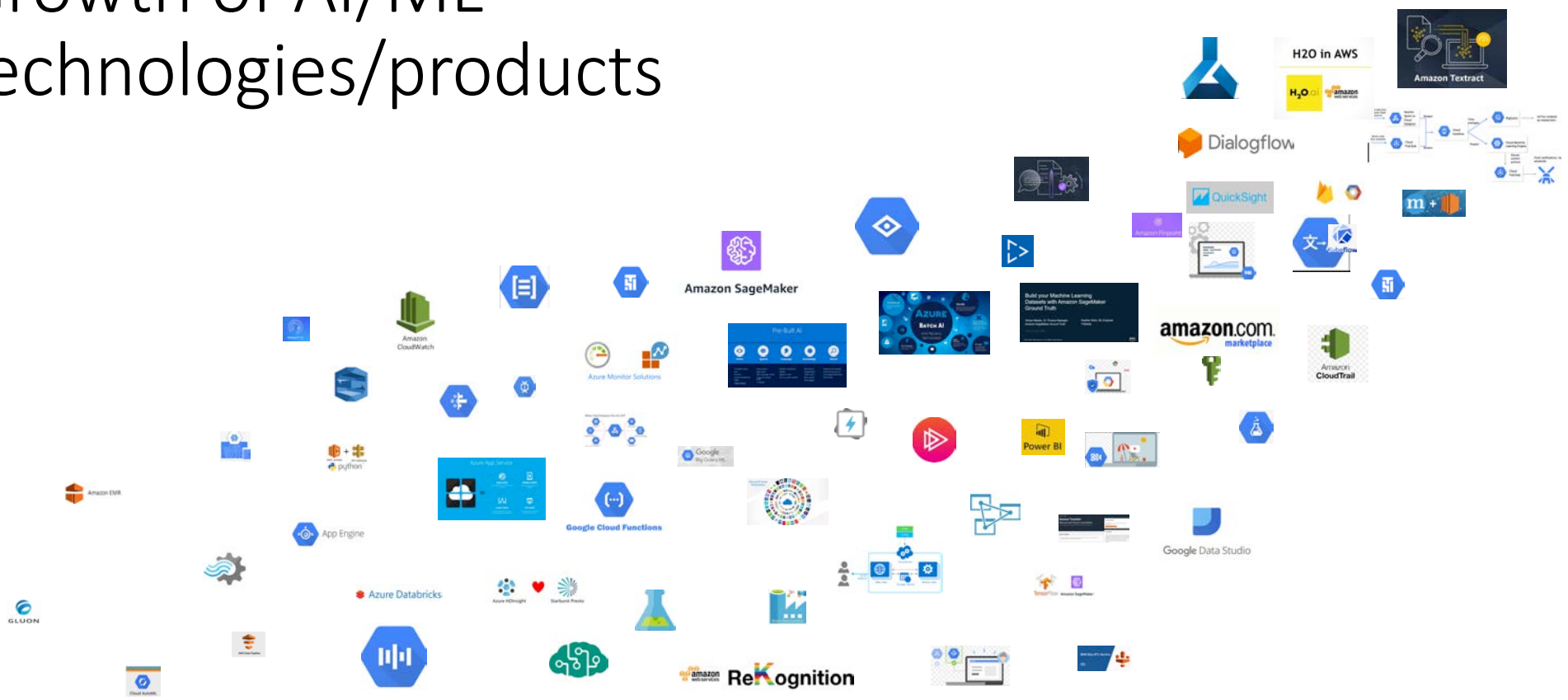
Nisha Talagala
CEO, Pyxeda AI

# Machine Learning Growth

**Data: Sources and Storage**

**Compute: Cloud, Hardware Innovation**
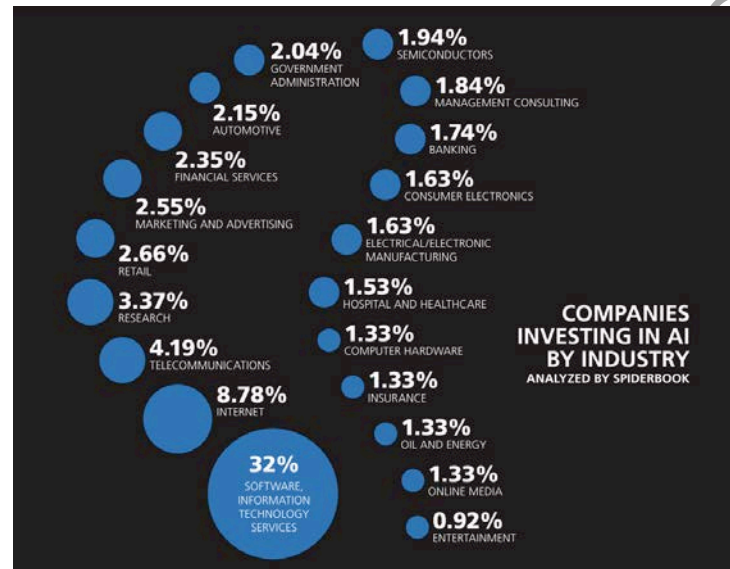
**Algorithms and Open Source**

# Growth of AI/ML technologies/products



Each logo is a (separate) service offered by GCP, AWS or Azure for part of an AI workflow

# Realities of Production Use



Companies Investing in AI by Industry
Analyzed by Spiderbook

- 2.04% Government Administration
- 1.94% Semiconductors
- 1.84% Management Consulting
- 2.15% Automotive
- 1.74% Banking
- 2.35% Financial Services
- 1.63% Consumer Electronics
- 2.55% Marketing and Advertising
- 1.63% Electrical/Electronic Manufacturing
- 2.66% Retail
- 1.53% Hospital and Healthcare
- 3.37% Research
- 1.33% Computer Hardware
- 4.19% Telecommunications
- 1.33% Insurance
- 8.78% Internet
- 1.33% Oil and Energy
- 32% Software, Information Technology Services
- 1.33% Online Media
- 0.92% Entertainment

*There are only 1,500 companies in North America that are doing anything related to AI today, even using its narrow, task-based definition. That means less than one percent of all medium-to-large companies across all industries are adopting AI.*

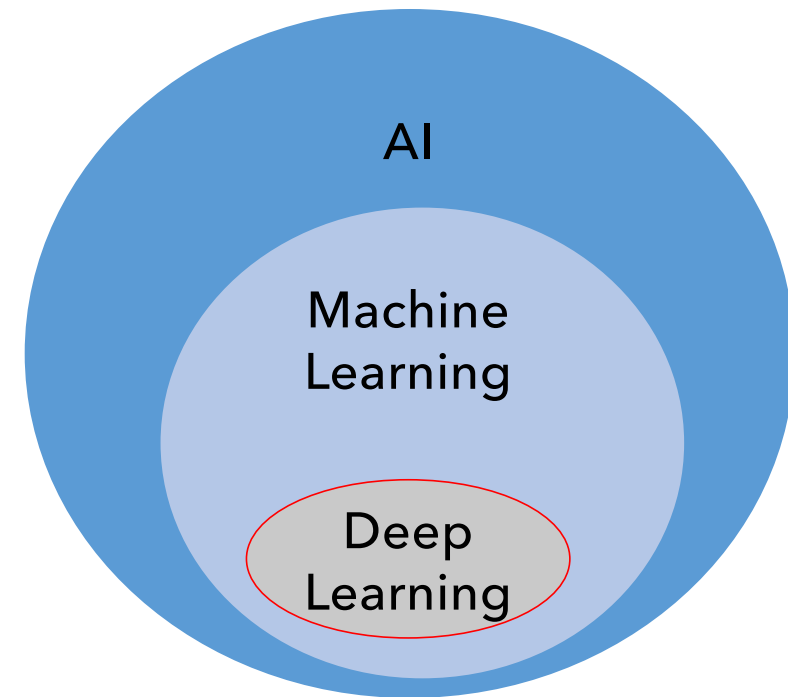**Despite the advanced services available, AI usage still minimal**

https://www.oreilly.com/library/view/the-new-artificial/9781492048978/
https://emerj.com/ai-sector-overviews/valuing-the-artificial-intelligence-market-graphs-and-predictions/

**Pyxeda**

# In This Talk:

- AI and ML: A quick overview

- Trends as relevant for Storage

  - Workloads

  - Trust, Governance and Data Management
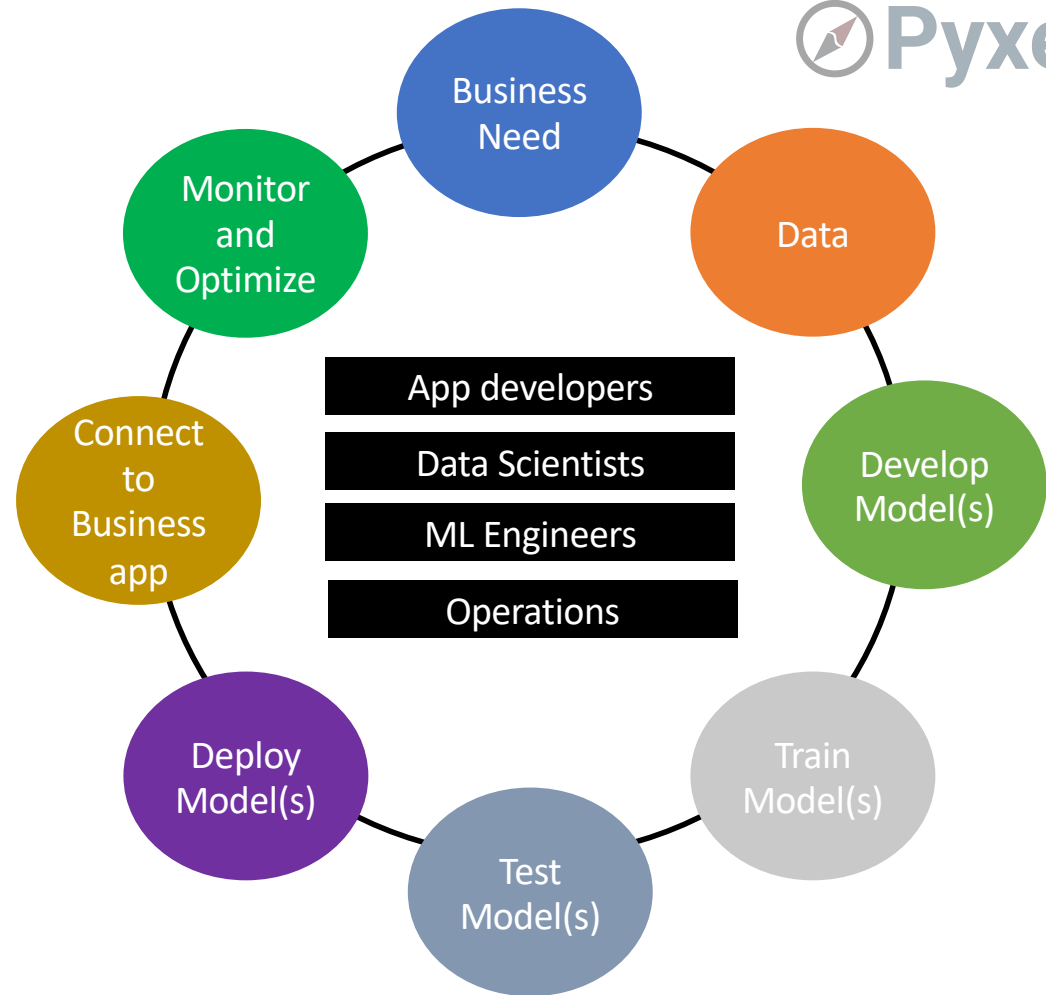
  - Edge

  - The users

# What is Machine Learning and AI?



- AI: Natural Language Processing, Image Recognition, Anomaly Detection, etc.
- Machine Learning: Supervised, Unsupervised, Reinforcement, Transfer, etc.
- Deep Learning: CNNs, RNNs etc.
- Common Threads
  - Training
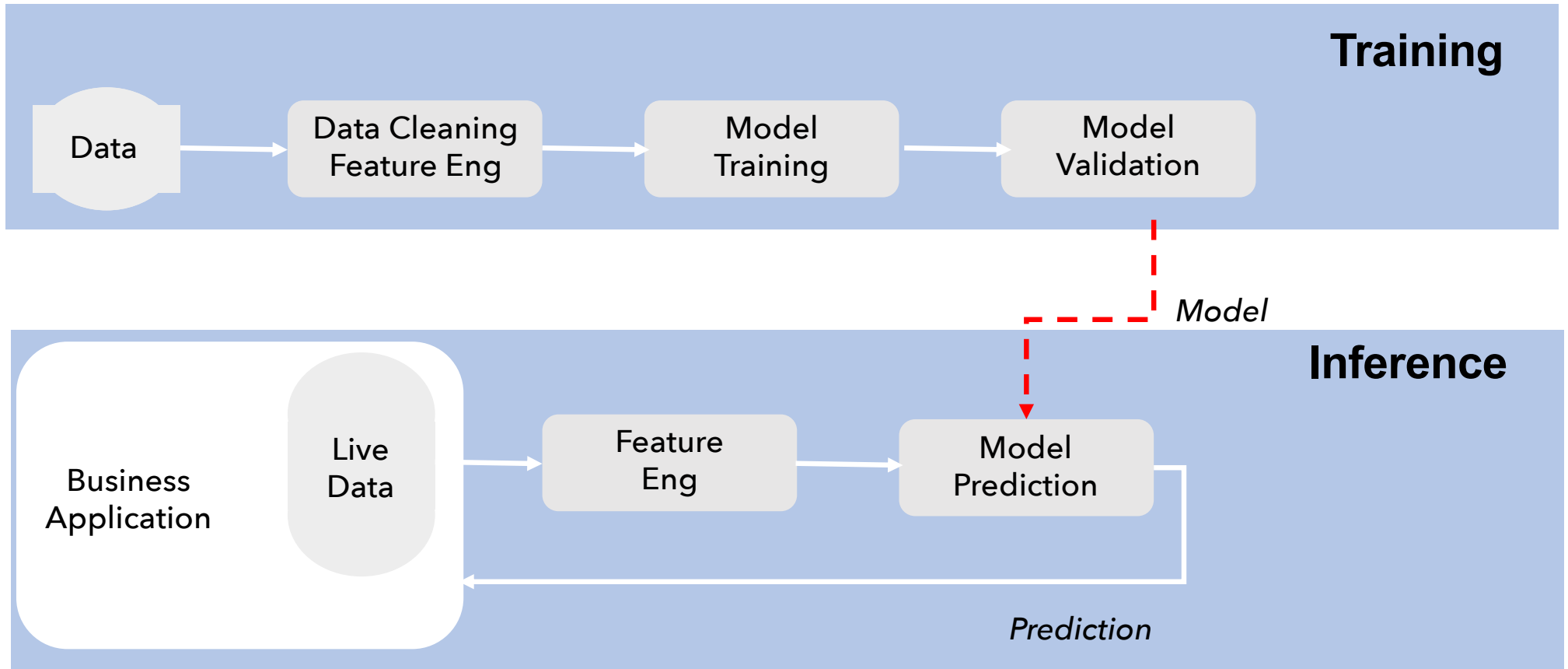  - Inference (aka Scoring, Model Serving, Prediction)

# A typical flow

- Use case definition
- Data prep
- Modeling
- Training
- Deploy
- Integrate
- Monitor/Optimize
- Iterate

**Pyxeda**

Business Need

Data

Develop Model(s)

Train Model(s)

Test Model(s)

Deploy Model(s)

Connect to Business app

Monitor and Optimize

App developers

Data Scientists

ML Engineers

Operations

# A Typical ML Operational Pipeline

**Pyxeda**

## Training

Data → Data Cleaning Feature Eng → Model Training → Model Validation

*Model*

## Inference

Business Application

Live Data → Feature Eng → Model Prediction

*Prediction*

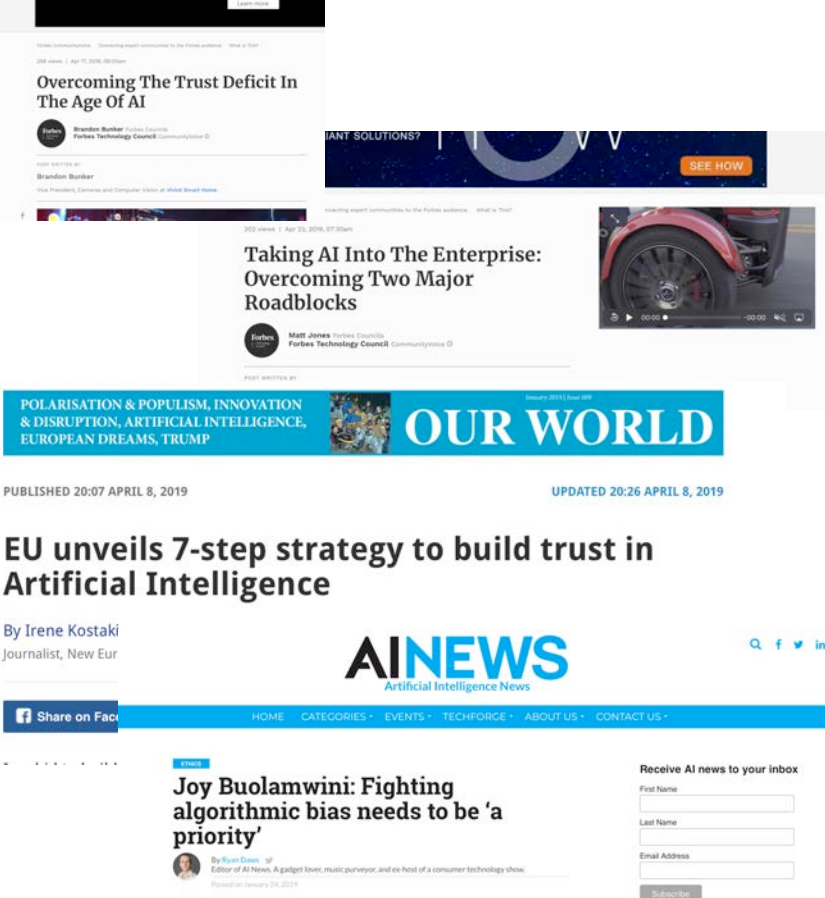# Trend 1: How ML/DL Workloads Think About Data **Pyxeda**

- Data Sizes
  - Incoming datasets can range from MB to TB
  - Statistical ML Models are typically small. Largest models tend to be in deep neural networks (DL) and range from 10s MB to GBs
- Common Structured Data Types
  - Time series and Streams
  - Multi-dimensional Arrays, Matrices and Vectors
- Common distributed patterns
  - Data Parallel, periodic synchronization
  - Model Parallel
  - Straggler performance issues can be significant

# Trend 1: How ML/DL Workloads Think About Data **Pyxeda**

- The older data gets – the more its "role" changes
  - Older data for batch- historical analytics and model reboots
  - Used for model training (sort of), not for inference
- Guarantees can be "flexible" on older data
  - Availability can be reduced (most algorithms can deal with some data loss)
  - A few data corruptions don't really hurt ☺
  - Data is evaluated in aggregate and algorithms are tolerant of outliers
  - Holes are a fact of real life data – algorithms deal with it
- Quality of service exists but is different
  - Random access is very rare
  - Heavily patterned access (most operations are some form of array/matrix)
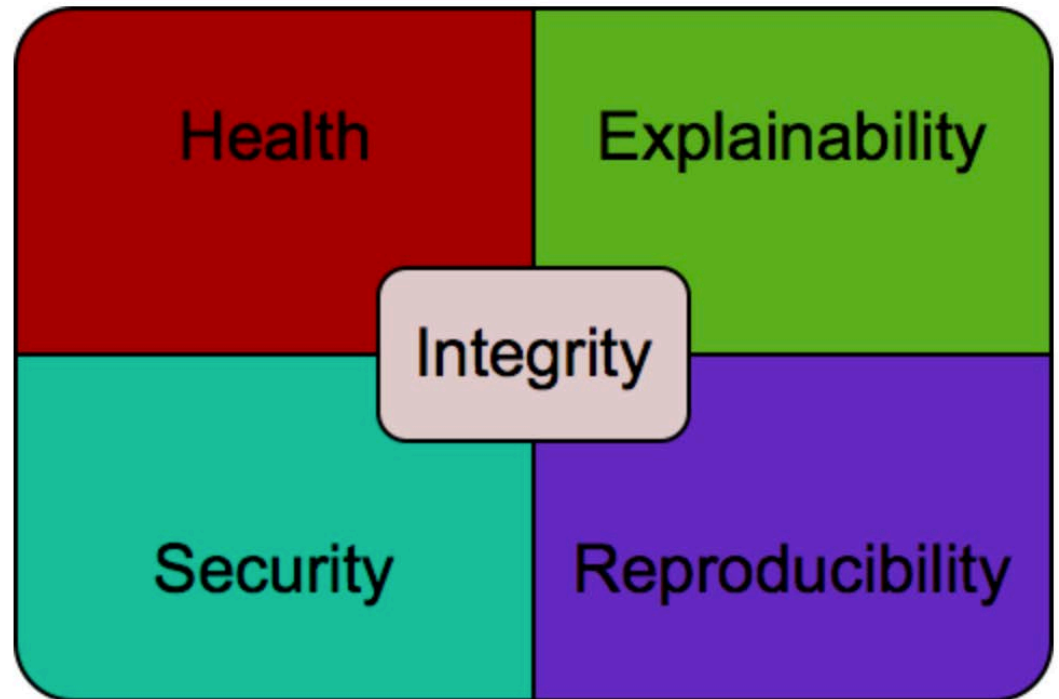  - Shuffle phase in some analytic engines

# AI Trust

- Publicized "mistakes" that damage corporate brands and generate business risk
  - Example Racism in Microsoft Tay bot and Bias in Amazon HR hiring tool
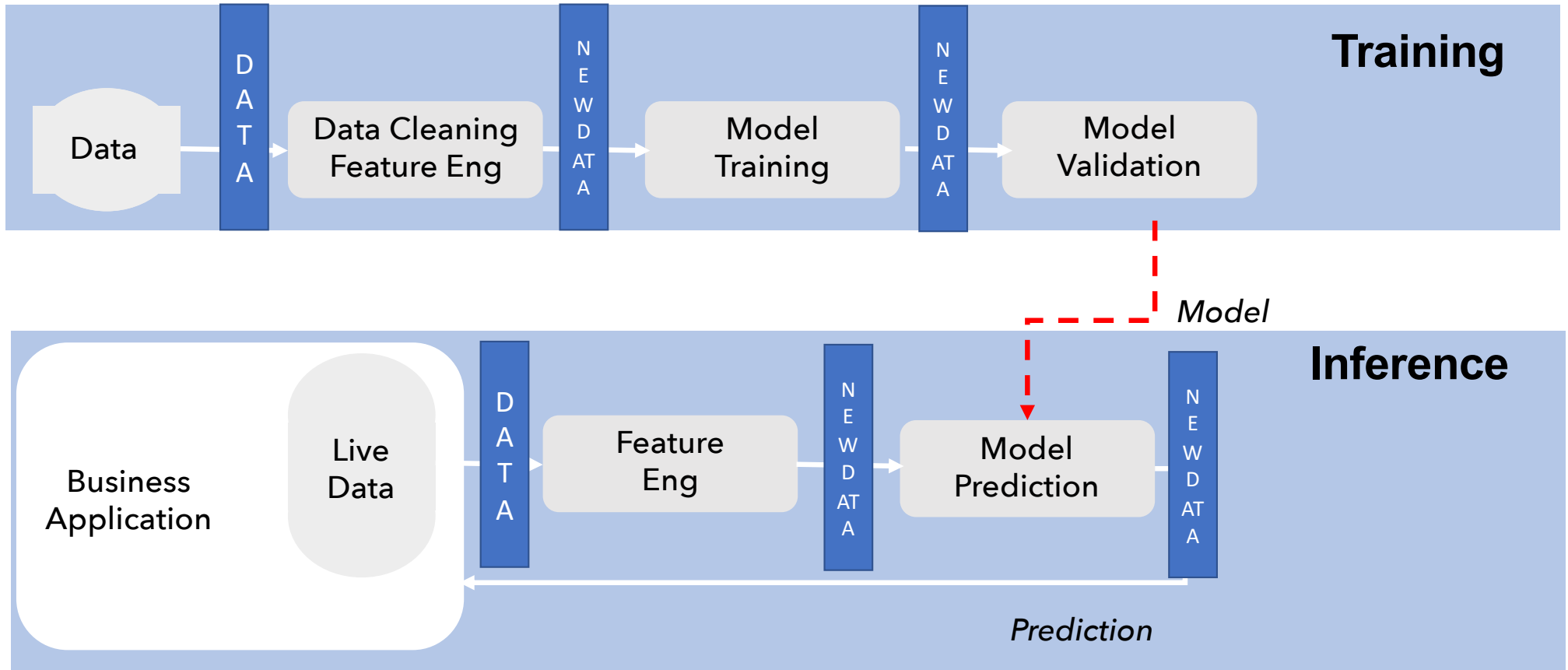- Intersection of AI decisions and human social values

# Pillars for AI Trust

- Together ensure that the ML is operating correctly and free from intrusion
- Details about how and why predictions and made
- Reproduce cases if needed

# What does this mean for data?



**Pyxeda**

**Training**

Data → DATA → Data Cleaning Feature Eng → NEW DATA → Model Training → NEW DATA → Model Validation

*Model*

**Inference**

Business Application — Live Data → DATA → Feature Eng → NEW DATA → Model Prediction → NEW DATA

*Prediction*

Access control, Lineage, Tracking of all data artifacts is critical for AI Trust
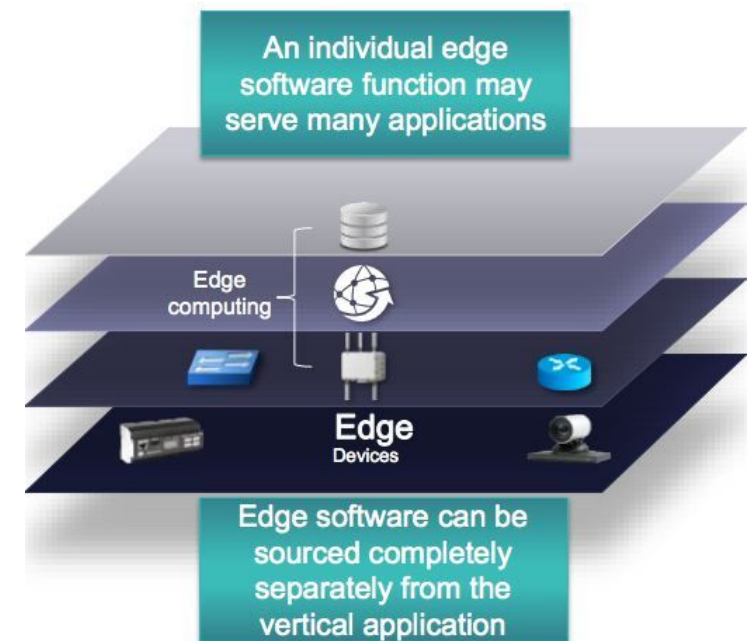
**Pyxeda**

# Trend 2: Need for Governance

- ML is only as good as its data
- Managing ML requires understanding **data provenance**
  - *How was it created? Where did it come from? When was it valid?*
  - *Who can access it? (all or subsets)? Which features were used for what?*
  - *How was it transformed?*
  - *What ML was it used for and when?*
- Solutions require both storage management and ML management

# Trend 2: Need for Governance

**Pyxeda**

- Examples
  - Established: Example: Model Risk Management in Financial Services
  - https://www.federalreserve.gov/supervisionreg/srletters/sr1107a1.pdf
- Example GDPR/CCPA on Data, Reproducing and Explaining ML Decisions
  - https://iapp.org/news/a/is-there-a-right-to-explanation-for-machine-learning-in-the-gdpr/
- Example: New York City Algorithm Fairness Monitoring
  - https://techcrunch.com/2017/12/12/new-york-city-moves-to-establish-algorithm-monitoring-task-force/

# Trend 3: The Growing Role of the Edge

**Pyxeda**

- Closest to data ingest, lowest latency.
  - Benefits to real time ML inference and (maybe later) training
- Varied hardware architectures and resource constraints
- Differs from geographically distributed data center architecture
- Creates need for cross cloud/edge data storage and management strategies



IoT Reference Model

# Trend 4: The Changing Role of Persistence

**Pyxeda**

- For ML functions, most computations today are in-memory
  - Data load and store are primary storage interaction
  - Intermediate data storage sometimes used
  - Tiered memory can be used within engines
- For in-memory databases, persistence is part of the core engine
  - Log based persistence is common
- Loading & cleaning of data is still a very large fraction of the pipeline time
  - Most of this involves manipulating stored data

# Trend 5: Who accesses the data

**Pyxeda**

- Multiple ML roles interact with data
  - Data Scientist
  - Decision Scientist, Decision Intelligence
  - Data Engineer / ML Engineer
- ML roles need to collaborate with Operations roles for successful Operational ML.
- Requires data access controls, access management to ensure ML consistency and governance

# Storage for ML: Challenges and Opportunities

**Pyxeda**

- Data access Speeds (Particularly for Deep Learning Workloads)

- Data Management

- Reproducibility and Lineage

- Governance and the Challenges of Regulation, Data Access Control and Access Management

- The Edge

- The new data managers

# Storage for ML: Example systems

- Databricks Delta

- Apache Atlas

- RDMA data acceleration for Deep Learning (Ex. from Mellanox)

- Time series optimized databases (Ex. BTrDB, GorrillaDB)

- API pushdown techniques and Native RDD Access APIs (Ex. Iguaz.io)

- Lineage: Link data and compute history (Ex. Alluxio/formerly Tachyon)

- Memory expansion (Ex. Many studies on DRAM/Persistent Memory/Flash tiering for analytics)

# Takeaways

**Pyxeda**

- The use of ML/DL in enterprise is at its infancy

- The first and most obvious storage challenge is performance
- The larger challenge is likely data management and governance
- Edge and distribution are also emerging challenges

- Opportunities exist to significantly improve storage and memory for these use cases

# Additional Resources

- NFS Vision report on Storage for 2025
  - See Storage and AI track

- Proceedings/Slides of USENIX OpML 2019

- Research at HotStorage, HotEdge, FAST, USENIX ATC

**Pyxeda**

# Thank You

Nisha Talagala

nisha@pyxeda.ai

# A Sample Analytics Stack: (Partial) Ecosystem

**Pyxeda**

**Algorithms and Libraries**

SparkML, TensorFlow

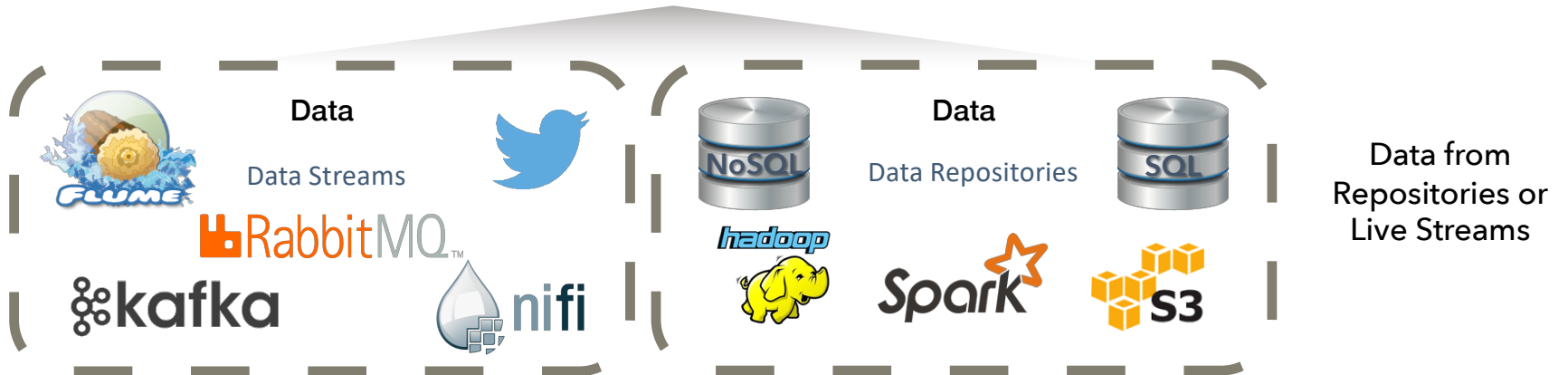**Processing Engines**

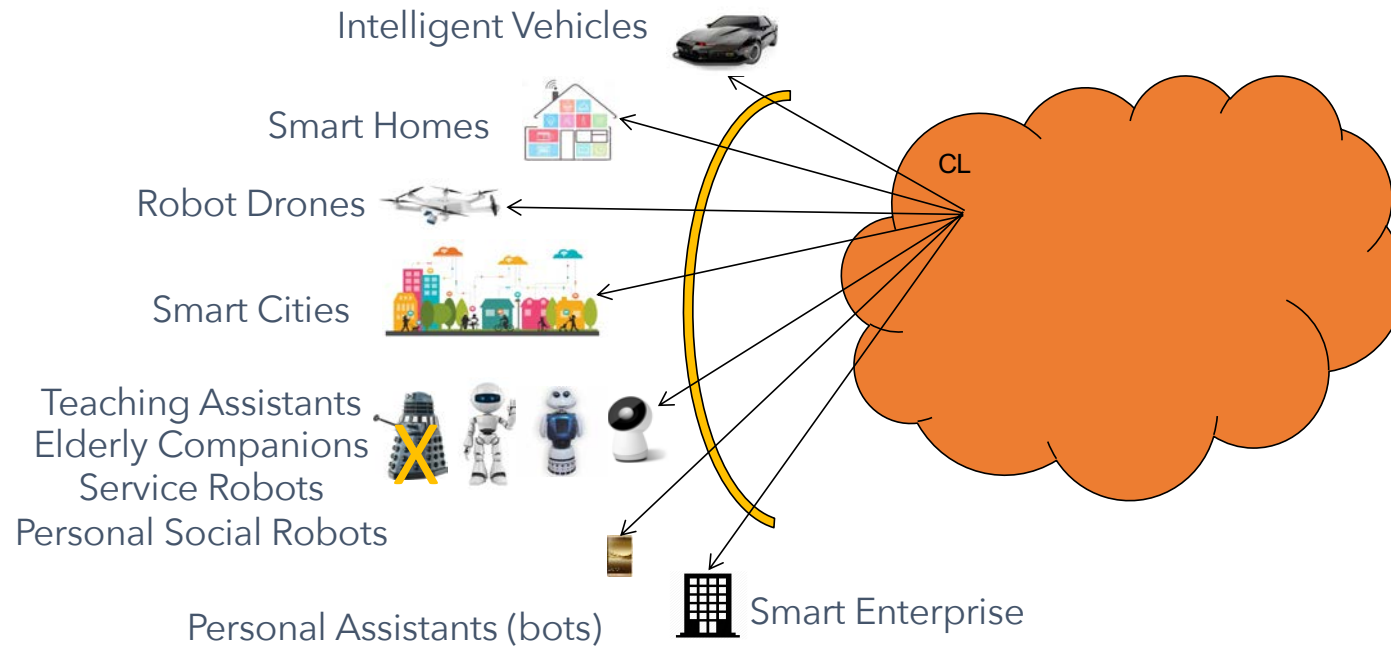| | | | |
|---|---|---|---|
| Hadoop<br>Spark<br>Tensor Flow | Flink / Apex<br>Spark Streaming<br>Storm / Samza / NiFi | Caffe<br>Tensor Flow<br>Pytorch | Containerized Models (Python etc.) |

**Data**

Data Streams

RabbitMQ

kafka · nifi

**Data**

NoSQL · SQL

Data Repositories

hadoop · Spark · S3

Data from Repositories or Live Streams

# Growing Sources of Data

Intelligent Vehicles

Smart Homes

Robot Drones

Smart Cities

Teaching Assistants
Elderly Companions
Service Robots
Personal Social Robots

Personal Assistants (bots)

CL

Smart Enterprise

**Edge** ⟷ **Cloud**