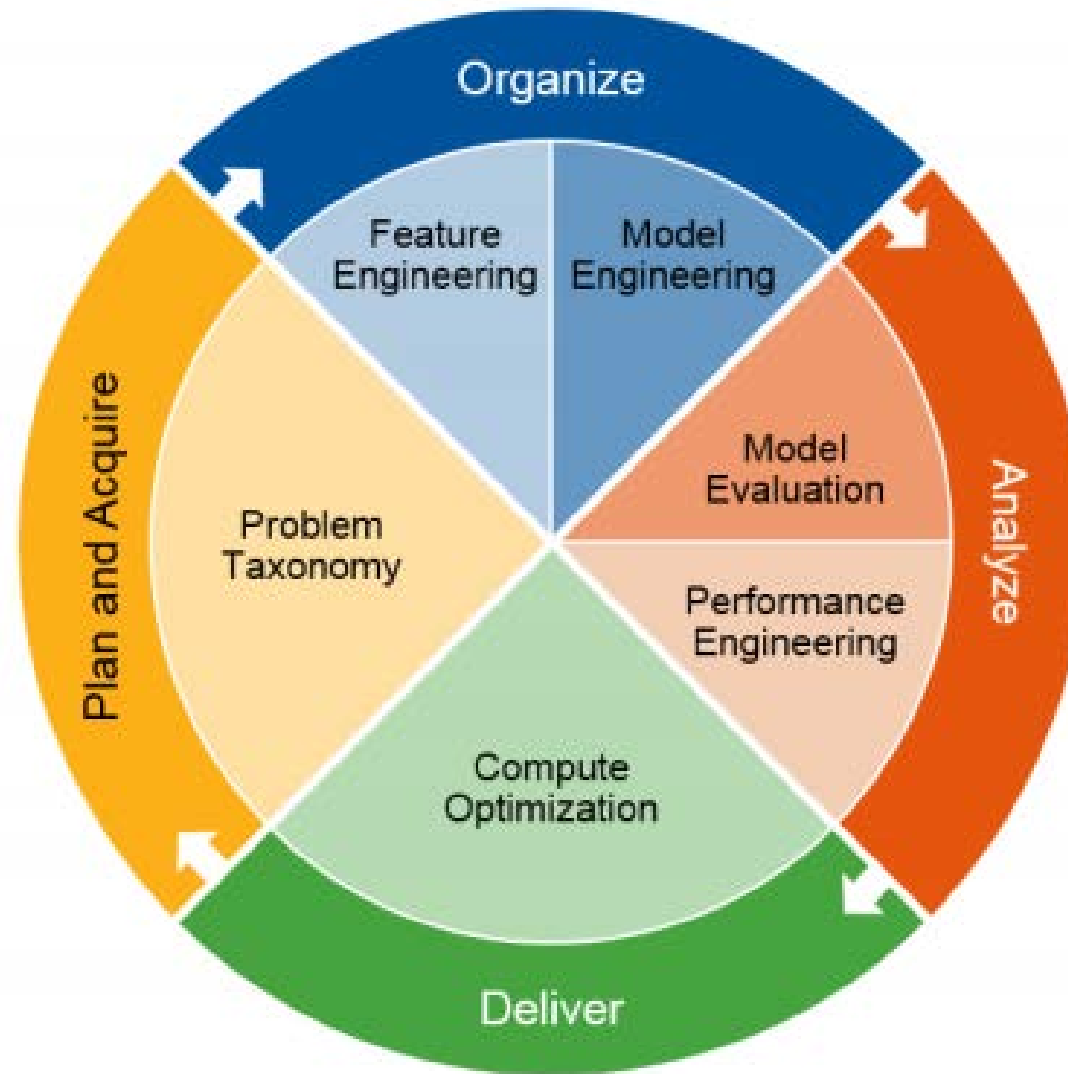# Core AI, Amazon Consumer



» Core tenet:  "Peak Jumping"

» Intersection of science, engineering and business

» Generating tangible value, and optimizing for Amazon's customers

» Partner with internal and external teams
  ▪ Supply Chain Optimization
  ▪ Pricing
  ▪ Amazon Fulfilment
  ▪ Finance teams
  ▪ among others

» Growing the scientist community in consumer business

# Model Application

# Why does Cloud Computing matter?

# Big Data

*"There is little doubt, at least in our own minds, that over the next decades "big data" will change the landscape of economic policy and economic research. As we emphasized throughout, we don't think that big data will substitute for common sense, economic theory, or the need for careful research designs. Rather, it will complement them. How exactly remains to be seen."*
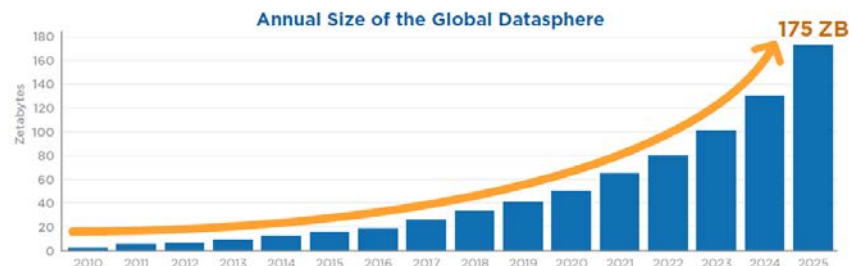
Liran Einav and Jonathan Levin, *The data revolution and economic analysis*
Technical report, NBER Innovation Policy and the Economy

amazon

## » Scalability

Figure 1 – Annual Size of the Global Datasphere

**Annual Size of the Global Datasphere**

175 ZB

- Applications of Data Science growing exponentially
- Dedicated hardware is inelastic, and has a short shelf-life

## » Technical Innovation

- Qualitatively different models may be more effective in this space

## » Collaboration

- Dispersed groups of people can interact virtually
- Easily share code and data in real-time
- Always-on availability and mobility
- Cloud-based workflow and sharing apps also help interaction with business

amazon

# Case Study

» Hedonic Price Regression to construct price indices

$$ln(p_{i,t}) = f(\alpha, x'_{i,t}\beta_t, \gamma'_i f_t, \varepsilon_{i,t})$$

- $\beta$: time varying implicit prices
  $\gamma$: latent product attributes
  $f$: time specific loadings

» After estimating the hedonic models, build Laspeyres, Paasche, and chain weighted indices

# Case Study: Hedonic Price Indices

» Single Business Unit
  - 4 years of history
  - Over 1Bn products

» Feature engineering
  - Structured text parsing
  - NLP of product attributes
  - Image processing (thumbnails)
  - Behavioral: Clicks, adds, purchases
  - Sales across Amazon.com, FBA, 3$^{rd}$-party merchants
  - Product search metadata

» Over 200Tb of feature data

amazon

# Case Study: Hedonic Price Indices

» Not feasible to manually encode regressors

» Complexity of models
  ▪ Need to go through hundreds of iterations of input data engineering, parameter tweaking, and experimentation with the algorithms themselves

» Integration with downstream consumer teams, operational cost

» This is a typical problem-space

amazon

# Functionality

- Cloud Storage

- Cloud Compute

- Development

- Collaboration



Cloud computing

# Storage (Data)

# Aspects of Cloud Storage

» Durability
  - Replication
  - Backup and Disaster Recovery

» Multi-tenancy

» Access Control

» Security
  - Can include lifecycle policy, etc.

» Structured vs unstructured data

» File-system based vs tabular

» Relational (SQL) vs NoSQL

» Serial vs Random Access

**amazon**

# Common Offerings



| Amazon EFS | Amazon EBS | Amazon EC2 Instance Store | Amazon S3 | Amazon Glacier |
|------------|------------|---------------------------|-----------|----------------|
| File | Block | | Object | |

| Amazon DynamoDB | Amazon RDS | Amazon Redshift | Amazon Redis & Memcached | Amazon SQS | Amazon Kinesis |
|-----------------|------------|-----------------|--------------------------|------------|----------------|
| Database | | | Other Storage | | |

# Cloud Storage for Econometric Datasets

» Characteristics

  ▪ Generally quantitative, semi-structured data

  ▪ Batch access

» Object-store with HDFS layer works well

» Storage Format

  ▪ Columnar, Binary formats:  Parquet or ORC

» Organization:  Data Catalog

amazon

# Data Discovery

# Cloud Data Catalog

» Discover, annotate, and share metadata about your datasets

» Automatic crawling and classification

# Cloud Data Catalog

# Cloud Compute

# Cloud Compute

» Typically separate from storage

» Cluster compute:  massive acceleration of most analytic jobs

» Parallelism has a different set of bottlenecks

- Network
- Disk
- Memory
- Threading

**amazon**

# Cloud Computing frameworks

» MapReduce (original)

» Apache Flink, Beam

» Spark
  - High level of community input
  - Supports SQL, streaming, and analytics
  - Interactive shells for code in Python, R, Scala
  - Supported by most cloud providers

» Other libraries with cluster support
  - e.g. TF-Cluster, R.parallel, etc.

# Breadth of Models with Spark/MLLib

» Classification

- Logistic regression
  - Binomial logistic regression
  - Multinomial logistic regression
- Decision tree classifier
- Random forest classifier
- Gradient-boosted tree classifier
- Multilayer perceptron classifier
- One-vs-Rest classifier (a.k.a. One-vs-All)
- Naive Bayes

» Regression

- Linear regression
- Generalized linear regression
- Decision tree regression
- Random forest regression
- Gradient-boosted tree regression
- Survival regression
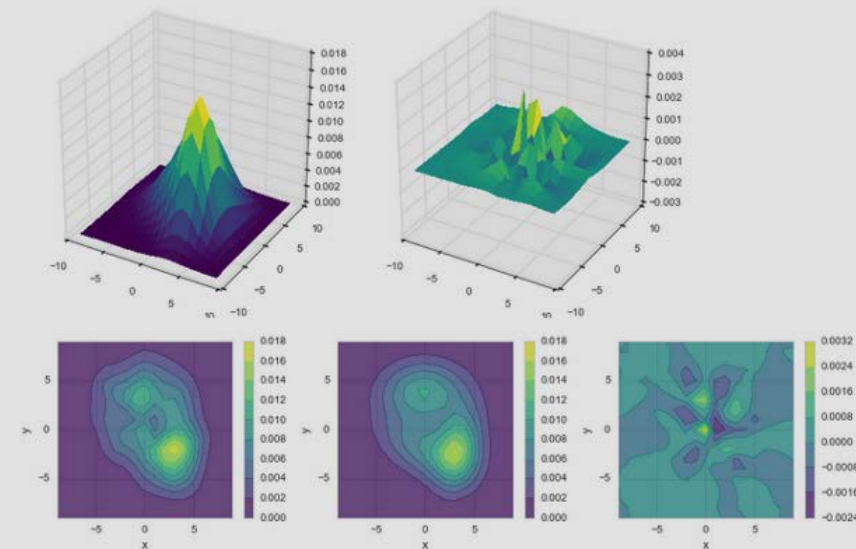- Isotonic regression

amazon

# Development

# Cloud Notebooks

» Fully managed online compute instances
  - Kernel runs on cloud servers
  - Minimal setup effort

» Portable, highly available

» Blend together live code, equations, visualizations, and explanatory text

» New generation
  - Full-featured IDE support
  - Launch into automated pipelines

# Cloud Notebooks

» Effortless to spin up more powerful instances or clusters

» Sharing the compute environment along with code

» Supported languages: Python, R, Scala…
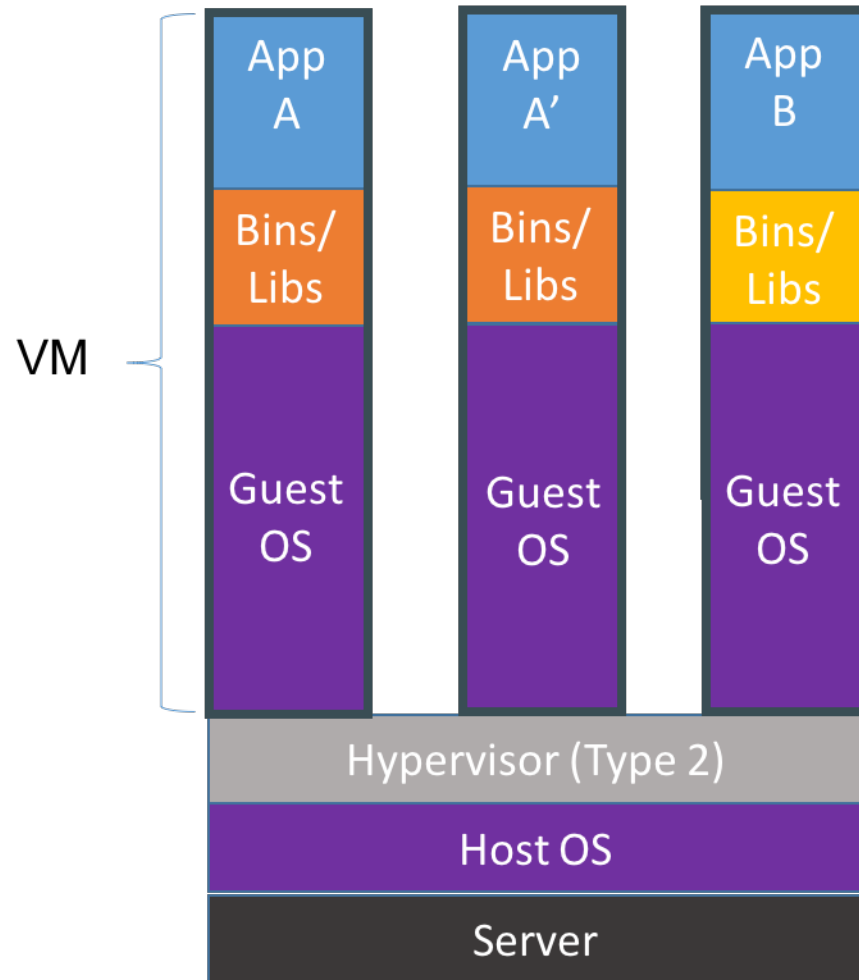
# Bundling other code packages

» Most common for cloud deployments:  Docker containers
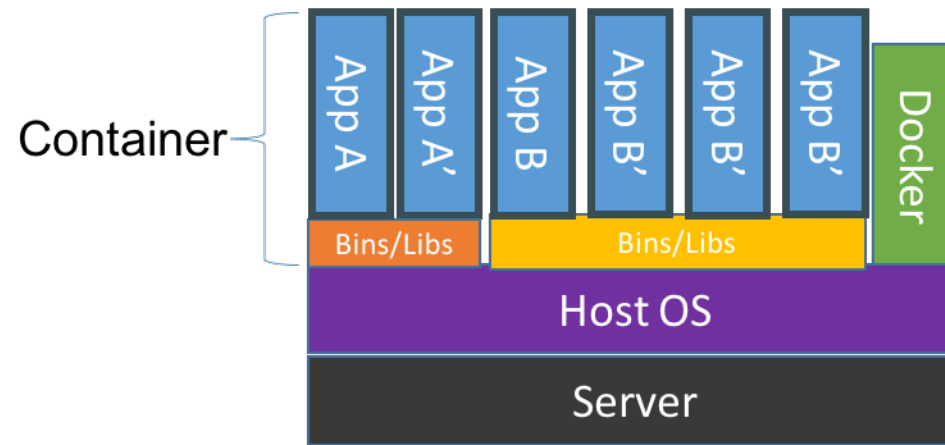
» Industry standard, lightweight, application isolation



Containerized Applications

App A | App B | App C | App D | App E | App F

Docker

Host Operating System

Infrastructure

# Containers



VM

App A
Bins/Libs
Guest OS

App A'
Bins/Libs
Guest OS

App B
Bins/Libs
Guest OS

Hypervisor (Type 2)

Host OS

Server

Containers are isolated, but share OS and, where appropriate, bins/libraries

…result is significantly faster deployment, much less overhead, easier migration, faster restart

Container

App A
App A'
App B
App B'
App B'
App B'
Docker

Bins/Libs
Bins/Libs

Host OS

Server

amazon

# Why Dockerize?

» Infrastructure
- Agility and scaling
- Standardized environments
- Portability
- Resource efficiency

» Application
- "Batteries included"
- Repeatable builds and orchestration
- Faster development cycles
- Lightweight

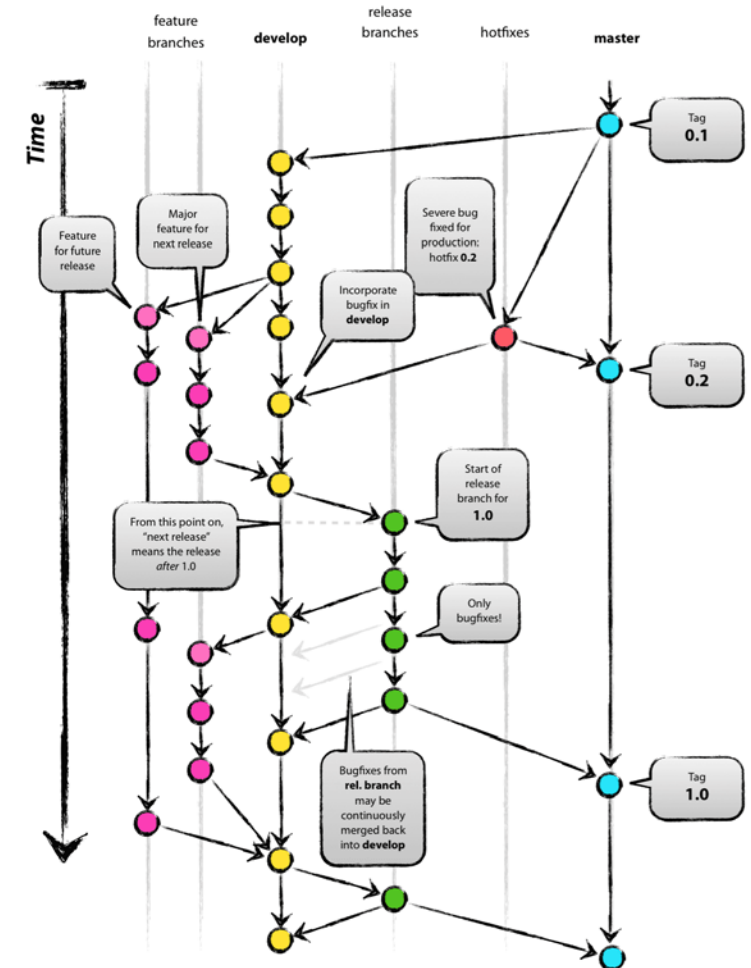**amazon**

# Many Docker image managers

# Collaboration

# Source Control System

» Benefits
- Long-term version history
  - Auditing, reversion, compare multiple versions
- Code Review
- Merging changes across contributors

» Cloud hosted git repo
- e.g. GitHub, CodeCommit, Cloud Source Repositories
- Scalable, secure, automatic backup
- Easy to share across a distributed team
- Code search
- Cloud notebook integration

# Model & experiment tracking

- » Models often have dozens of configurable parameters

- » Code changes over time

- » Need to explicitly track parameters, code, and input data that went into a given run

- » Hard to reproduce

- » Exponentially more challenging in a larger organization or team

**amazon**

# Hosted Model Registry

» e.g. MLflow Tracking

» Similar idea to VCS, but parameters and input data snapshots are first-class entities

# Recap

- Cloud Storage

- Cloud Compute

- Development

- Collaboration


Cloud computing

What's next?

# On the horizon
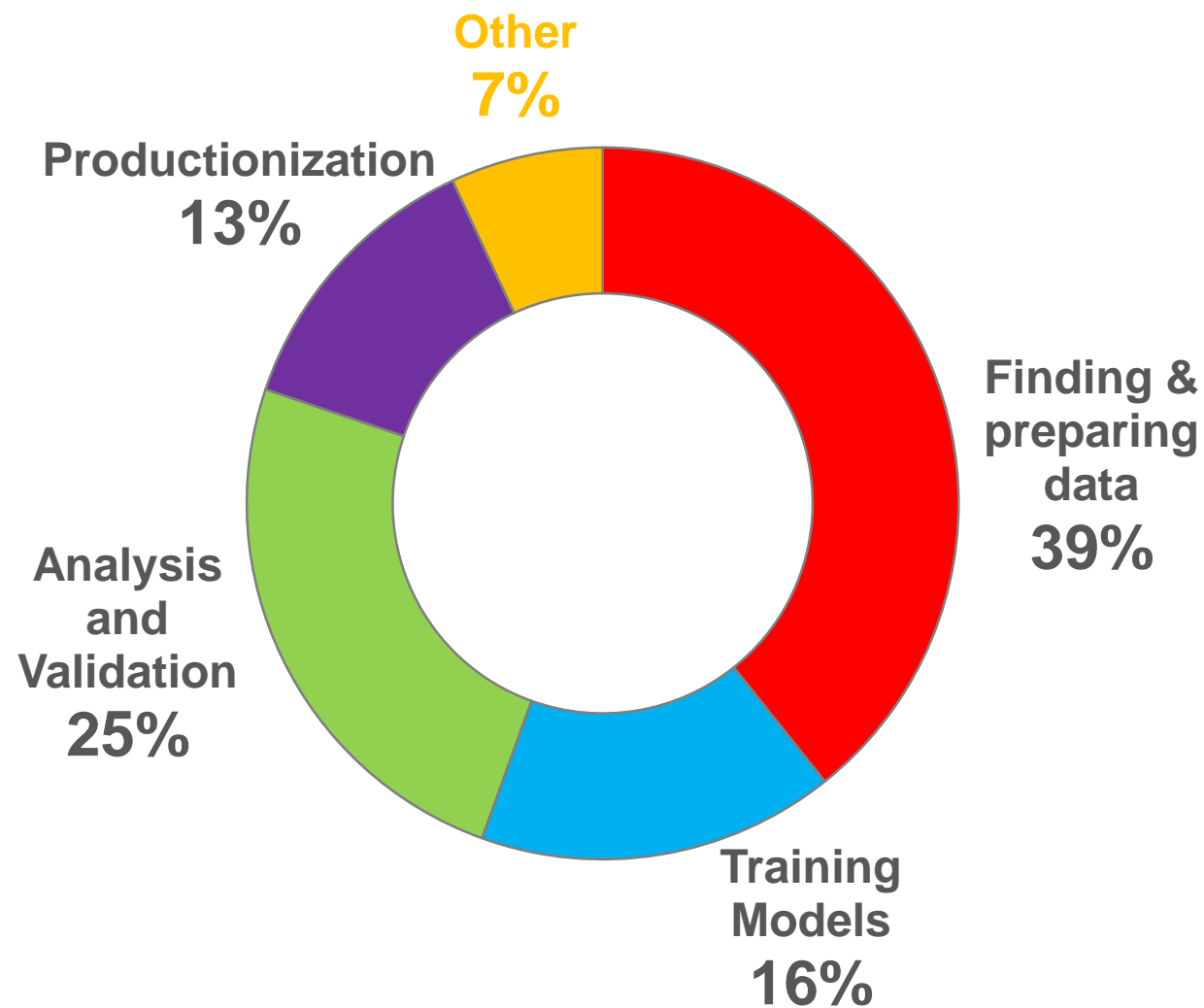
» Beyond Spark
- Ray
- Avoids block-synchronous compute paradigm
- Supported libraries including Modin, RLLib, TUNE

» ML-Economics integration
- Cloud support for DL-type workloads

» Deeper model dependency/lineage

*https://arxiv.org/abs/1712.05889*

# Environment Improvements



How Economists Spend their time? – Economist Bi-annual Survey 2019

- Other 7%
- Productionization 13%
- Finding & preparing data 39%
- Training Models 16%
- Analysis and Validation 25%

amazon

# Questions?

coreai-org@amazon.com

amazon