



# Data Science on TAP

Kyle H. Ambert, PhD

Intel Big Data Solutions, Datacenter Group

# About me.

---



# About me.

---



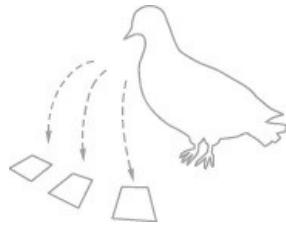
# About me.

---



# About me.

---

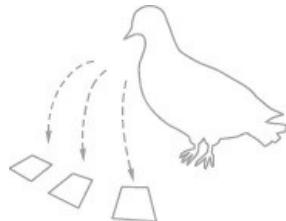


# About me.

---



10001001010  
0111000100111001  
01010010110100101001  
10110100011101001001010  
01101001001010100110101  
1001001010010111000100  
101101100101001000101  
011010010010010010  
10010010010  
01010101  
110010  
0

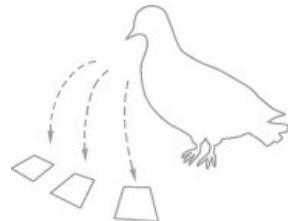


# About me.

---



10001001010  
0111000100111001  
01010010110100101001  
10110100011101001001010  
01101001001010100110101  
1001001010010111000100  
101101100101001000101  
011010010010010010  
10010010010  
01010101  
110010  
0

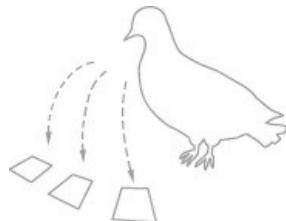


# About me.

---



10001001010  
0111000100111001  
01010010110100101001  
1011010011101001001010  
01101001001010100110101  
1001001010010111000100  
101101100101001000101  
011010010010010010  
10010010010  
011010101  
110010  
0

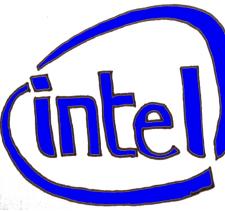
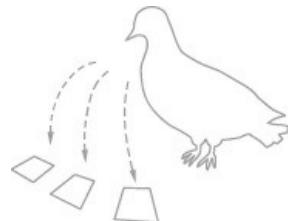


# About me.

---



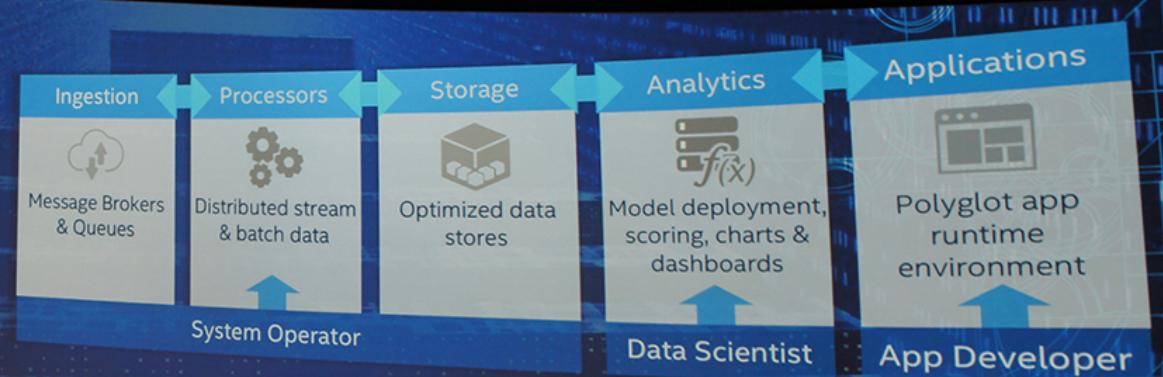
10001001010  
0111000100111001  
01010010110100101001  
1011010011101001001010  
01101001001010100110101  
1001001010010111000100  
101101100101001000101  
011010010010010010  
10010010010  
011010101  
110010  
0



# RELEASING DISCOVERY PEAK

## OPEN SOURCE ANALYTICS PLATFORM

for Data Scientists &  
Application Developers



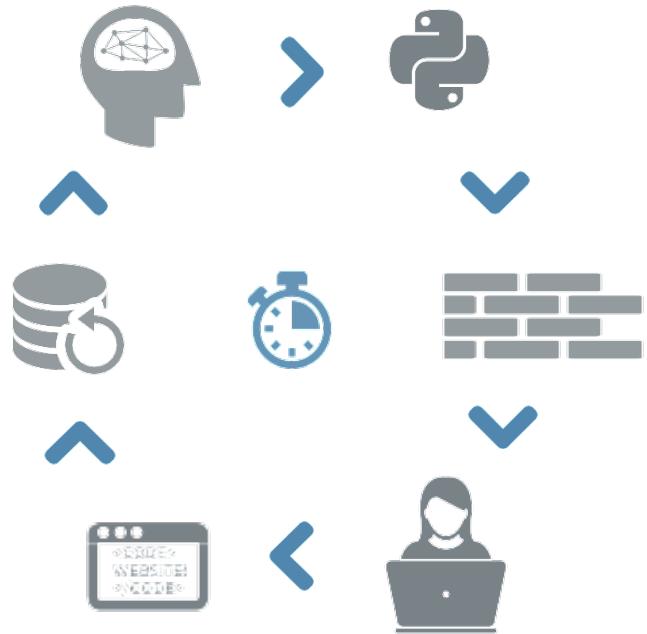
IDF15

FORUM

A woman in a striped dress is speaking on stage.

IDF15

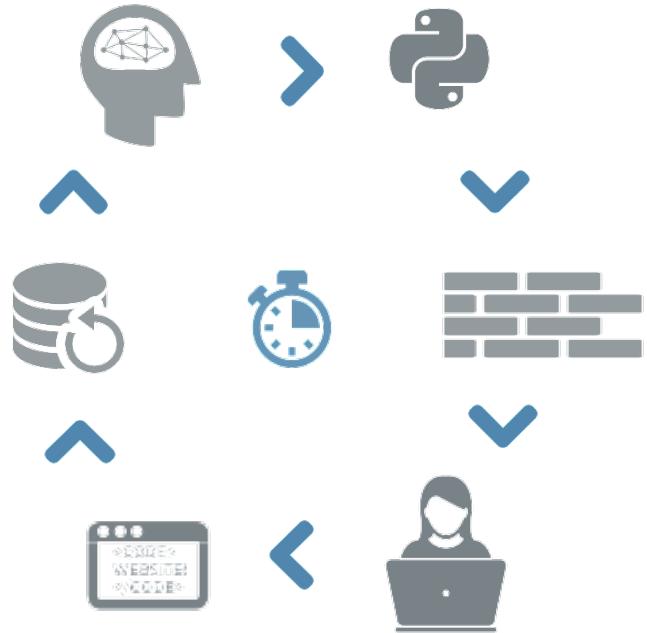
# Problem Statement



From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

Data Science:

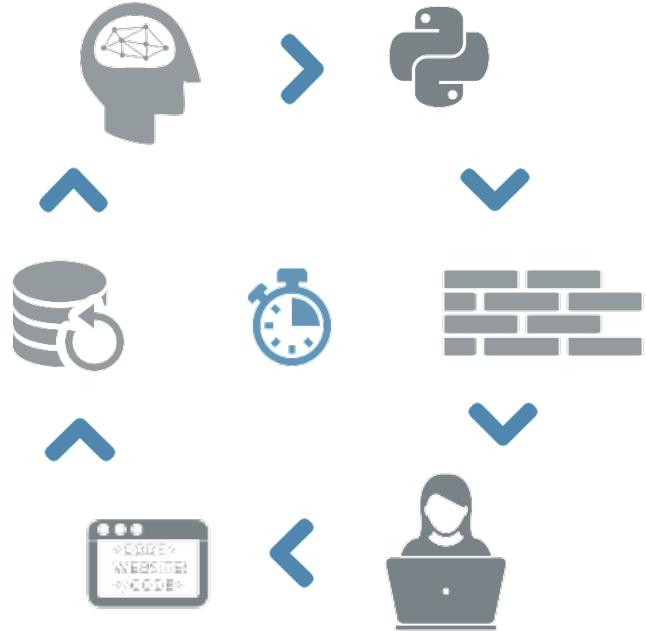


From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

Data Science:

- Iterative error-prone drudgery

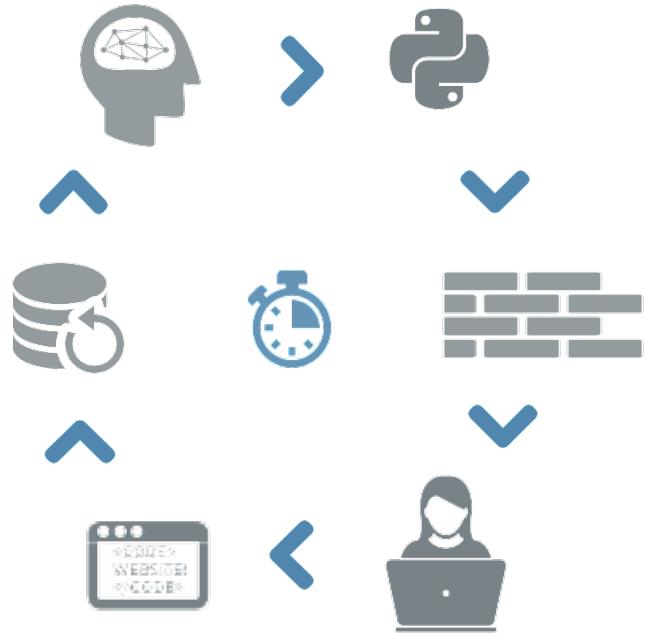


From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

Data Science:

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation



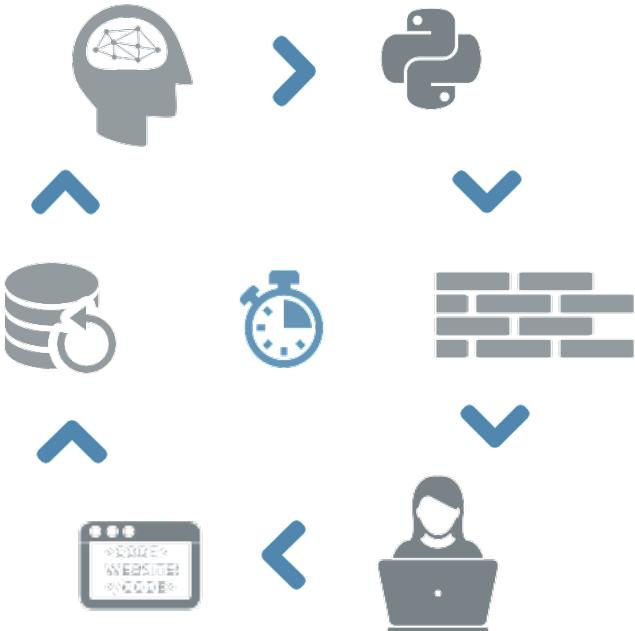
From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

Data Science:

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

Analytics Processing:



From data science to big data analytics: Less alchemy, more chemistry

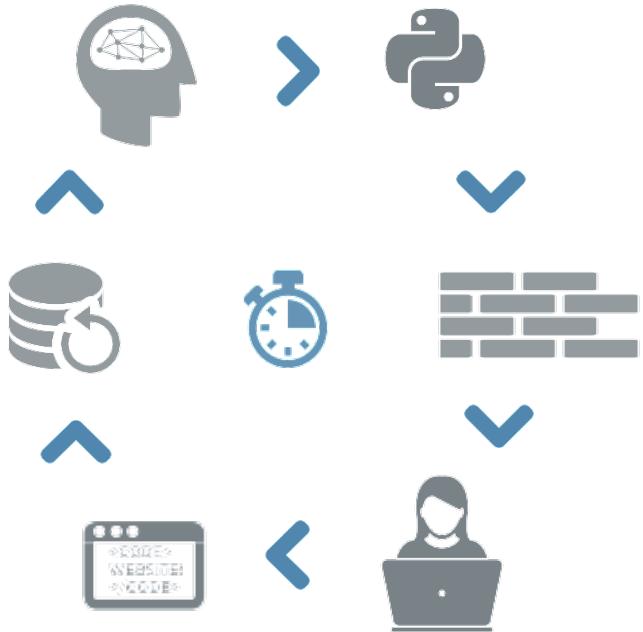
# Problem Statement

Data Science:

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

Analytics Processing:

- Single-threaded, single-node processing



From data science to big data analytics: Less alchemy, more chemistry

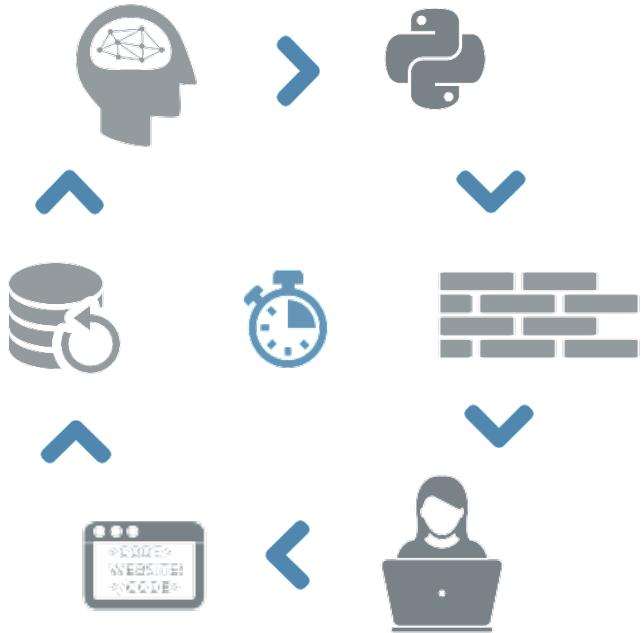
# Problem Statement

Data Science:

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

Analytics Processing:

- Single-threaded, single-node processing
- Proprietary, fixed-function solutions



From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

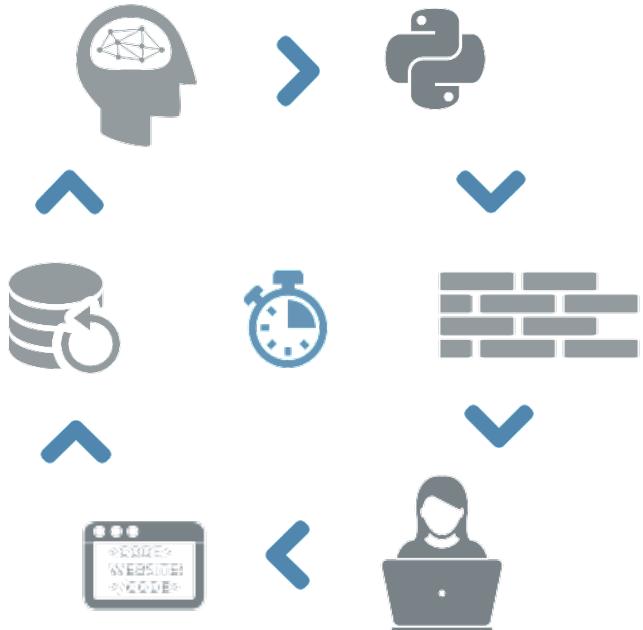
Data Science:

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

Analytics Processing:

- Single-threaded, single-node processing
- Proprietary, fixed-function solutions

Application Code:



From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

## Data Science:

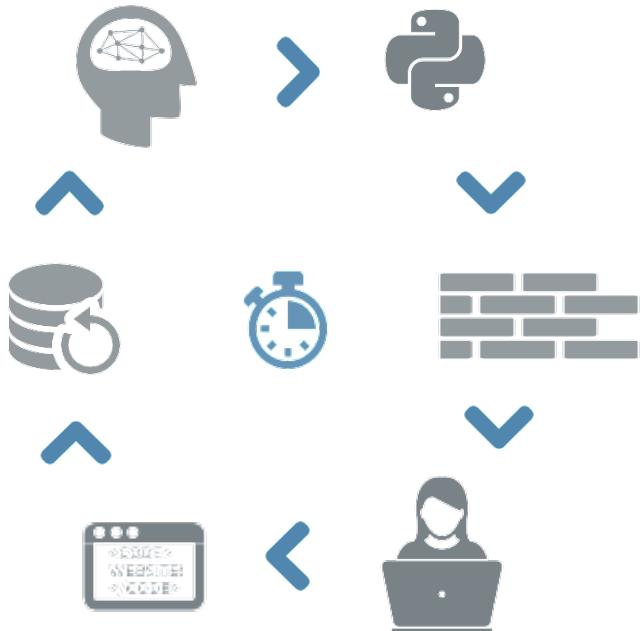
- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

## Analytics Processing:

- Single-threaded, single-node processing
- Proprietary, fixed-function solutions

## Application Code:

- Monolithic architecture



From data science to big data analytics: Less alchemy, more chemistry

# Problem Statement

## Data Science:

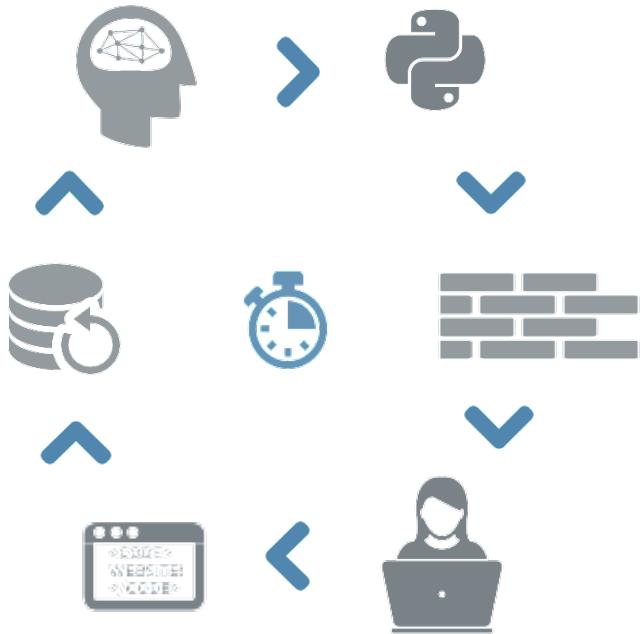
- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

## Analytics Processing:

- Single-threaded, single-node processing
- Proprietary, fixed-function solutions

## Application Code:

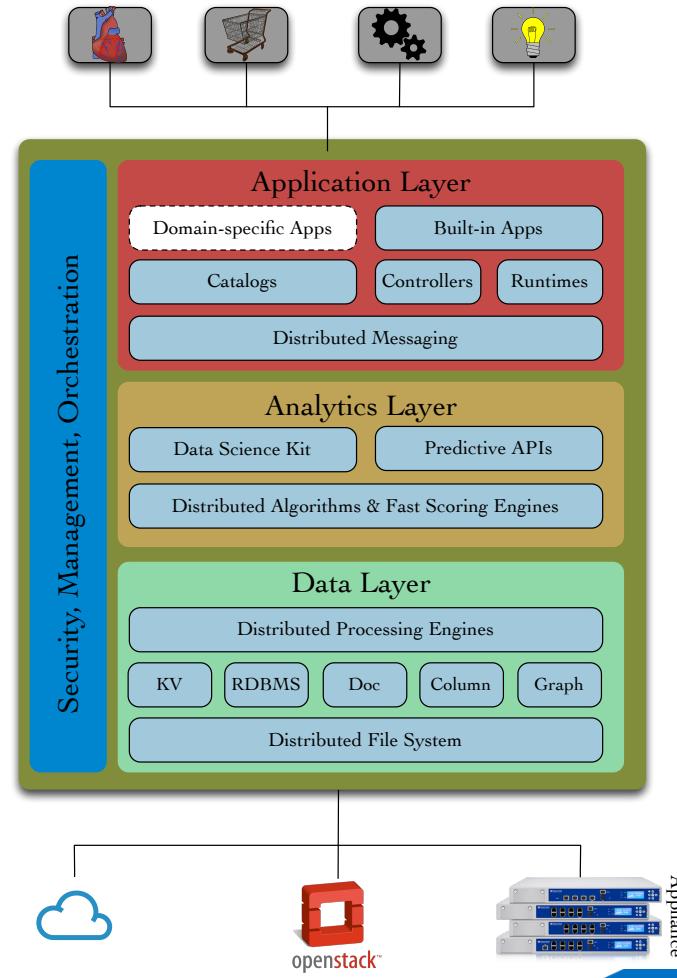
- Monolithic architecture
- Legacy components



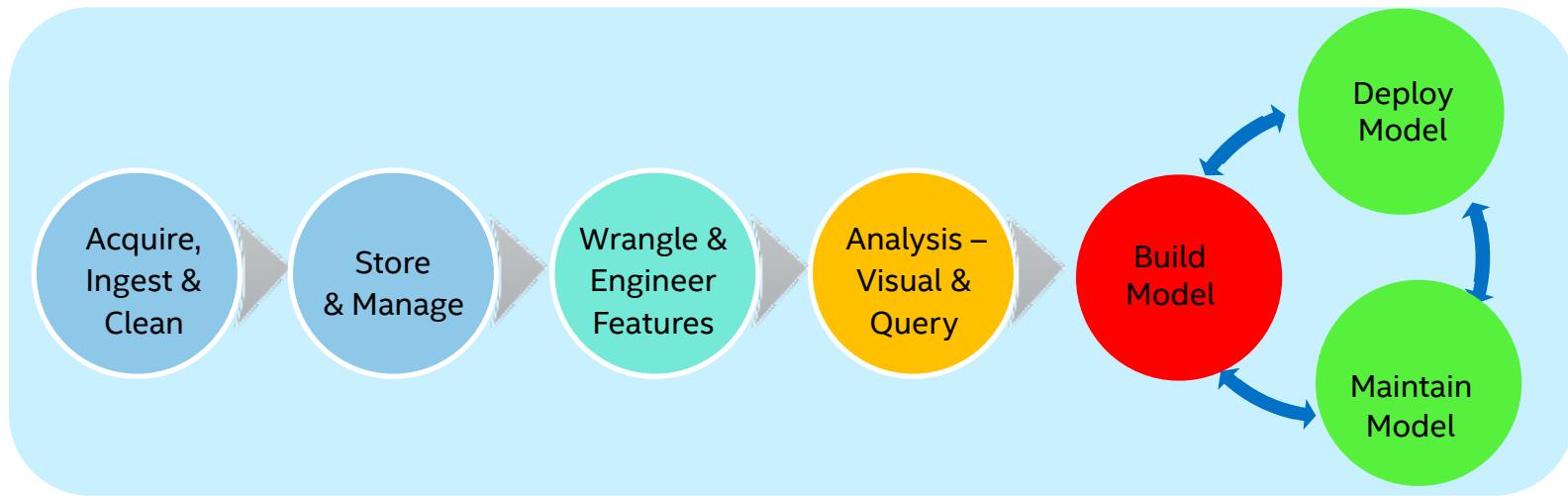
From data science to big data analytics: Less alchemy, more chemistry

# Trusted Analytics Platform (TAP)

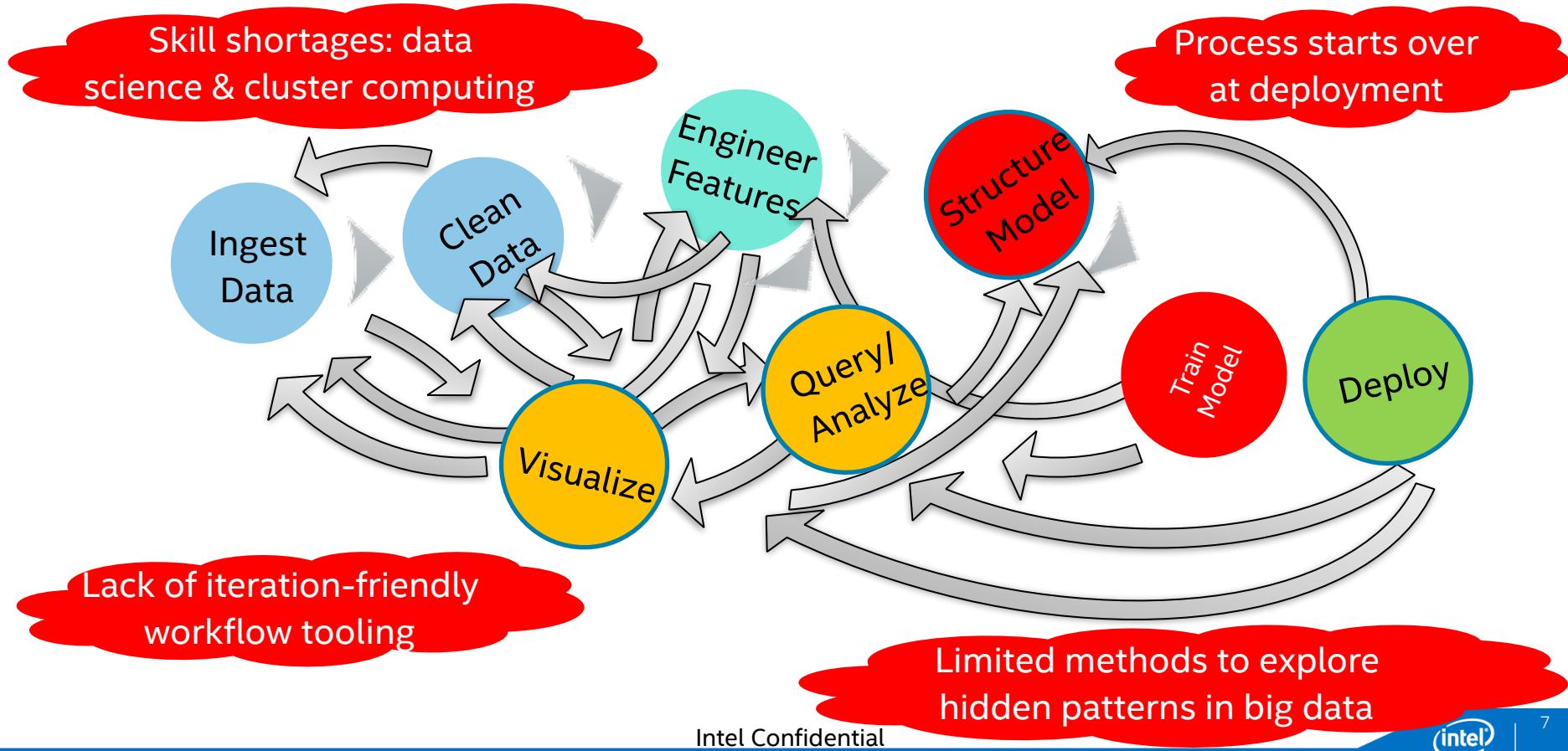
Open source software project that accelerates creation of **Cloud-native** apps driven by big data analytics. TAP makes it easier for developers to collaborate with data scientists by providing a shared environment for advanced analytics on big data.



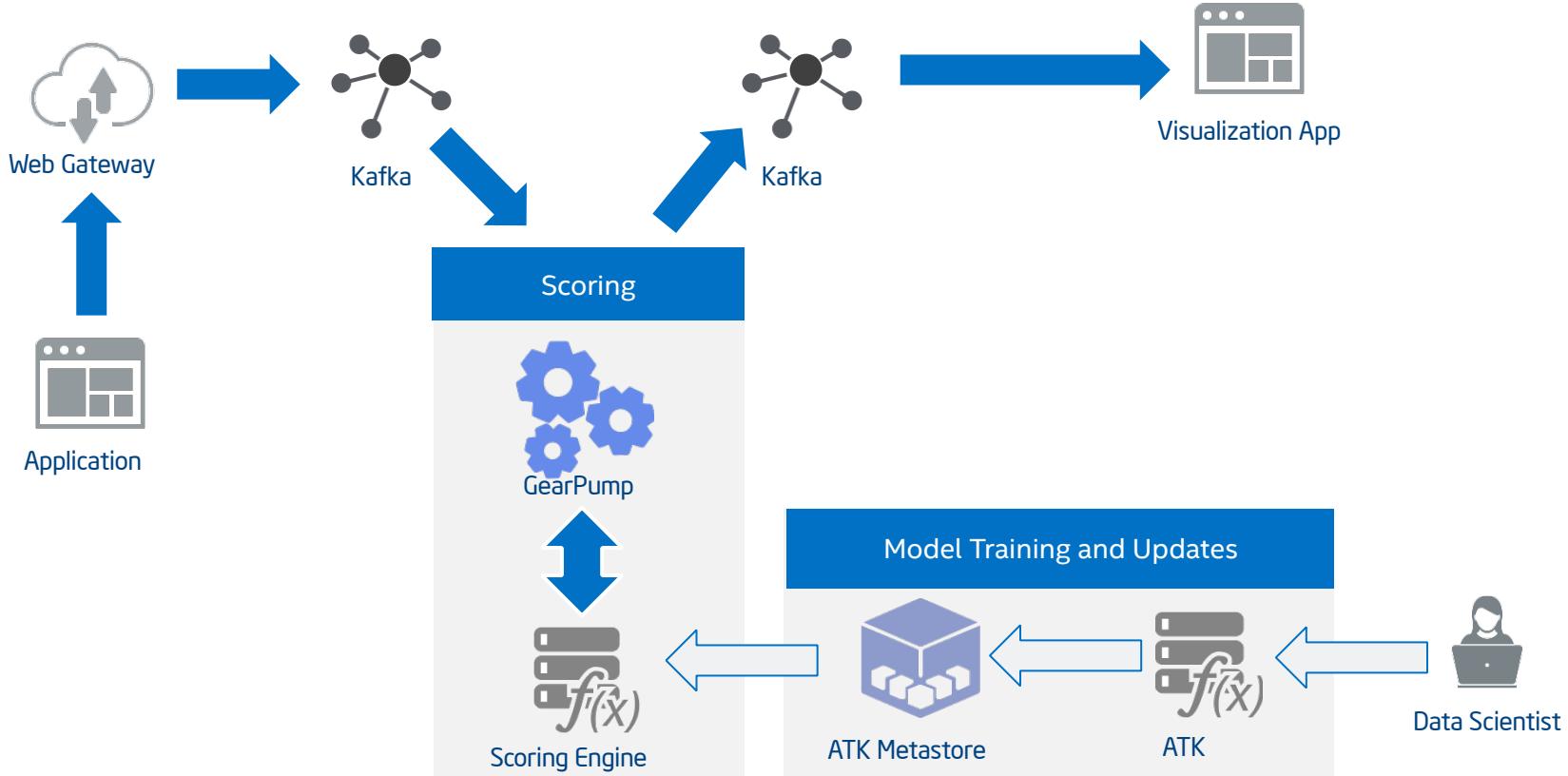
# The Data Science Workflow



# Challenges Across The Workflow



# Example: App With End to End Scoring



# Our Analytics Toolkit

A unifying framework for analytics pipelines

An extensible platform for data science

- Development
  - Interactive
  - Language integration
  - Familiar data structures
  - Extensible data science pipeline (methods, functions, engines)
  - Separates the role of Data Scientist from Hadoop Expert
- Production
  - Interactive solution deployed to production as prediction service / scoring engine

Data Scientist & Developer Tools

Analytics APIs

Data Wrangling &  
Feature Creation  
Tools

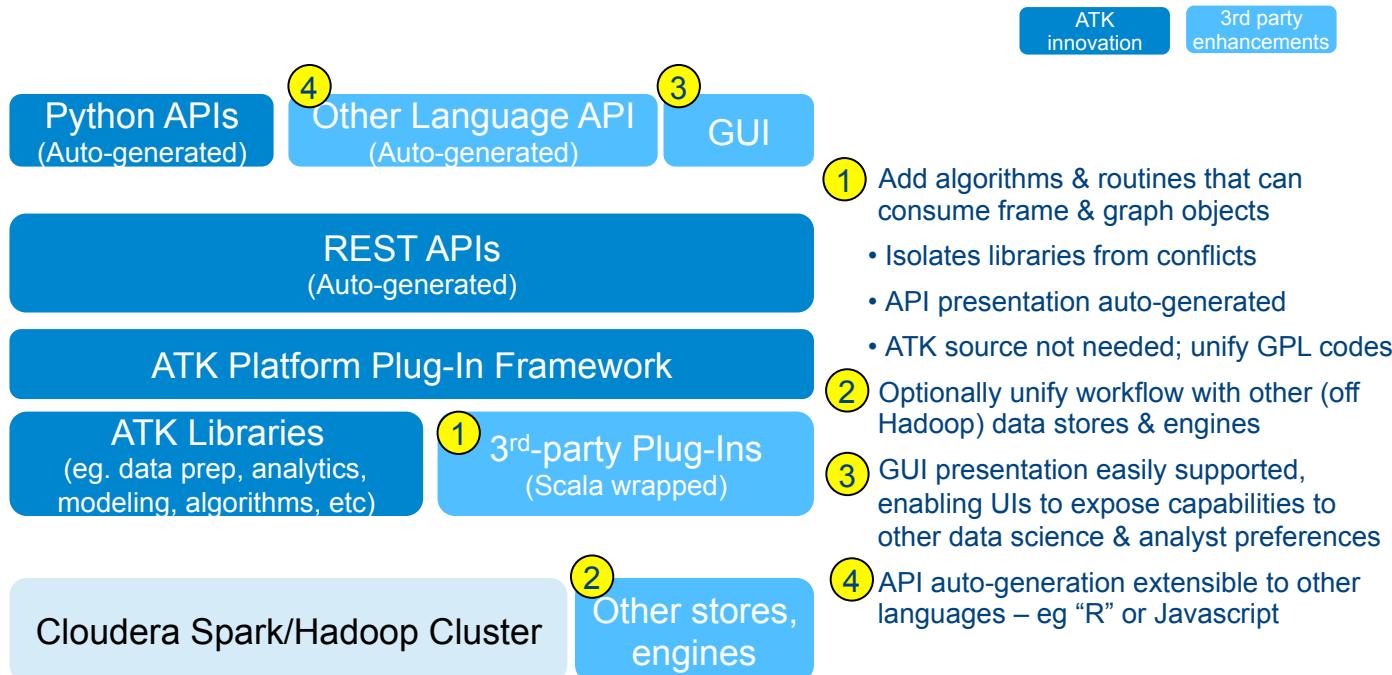
Graph Analysis &  
Machine Learning  
Algorithms

Cloudera® (IA Optimized Hadoop & SPARK)

Falcon Peak (IA Optimized Machine Learning)

Hardware (IA)

# Our Analytics Toolkit: Extensibility



# Machine Learning for Predicting Risk for Emergency Room Readmittance

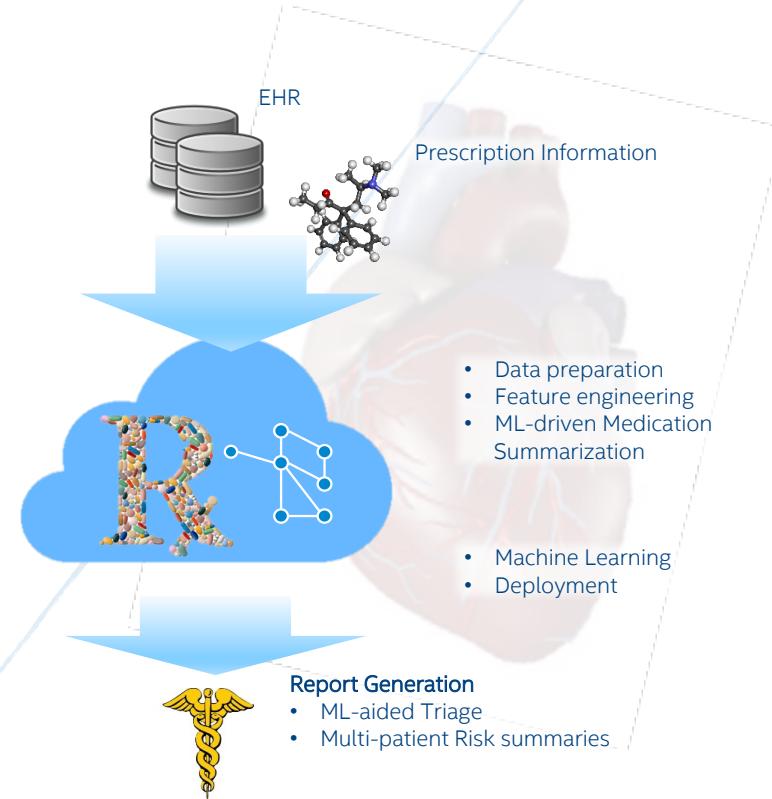
## Patient Outcomes

- At-risk patients can be identified on admit
- Enables more personalized treatment

## Hospital Outcomes

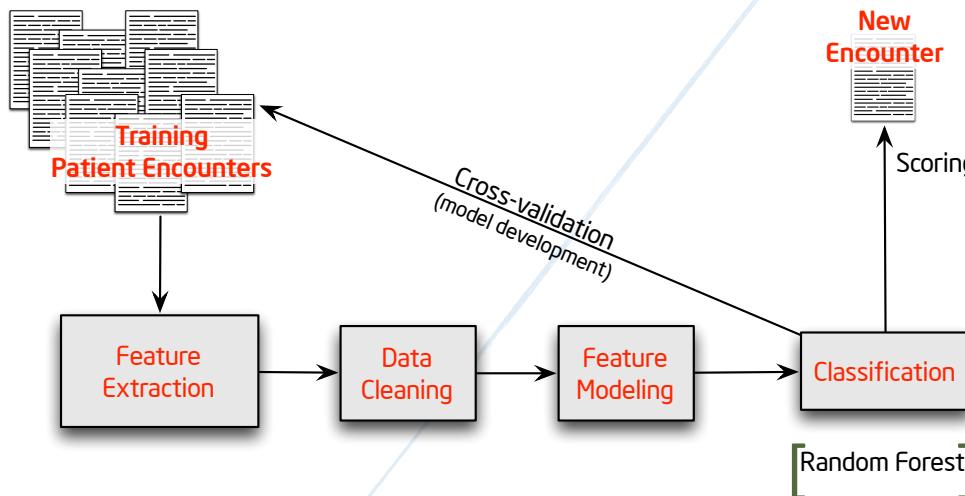
- Decreased costs
- More efficient use of resources

By leveraging big data in health analytics, patients at risk for negative outcomes can be identified early on, enabling better care



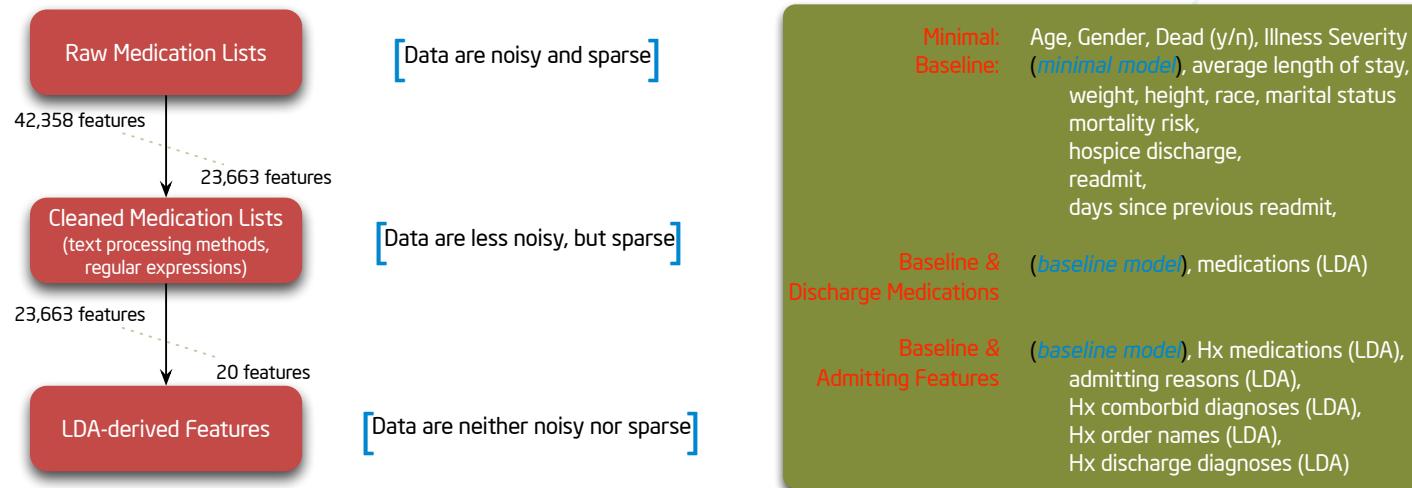
# Iterative Model Building in the Trusted Analytics Platform (TAP)

- De-identified electronic health record data was obtained for one year of patient encounters at Penn Med
- A cohort of emergency-admitted cardiac-related patients was identified ( $\approx 140K$ )
- We conducted feature selection and modeling experiments to iteratively construct a model useful for scoring previously-unseen encounters
- We built a classifier to identify which patients would be readmitted within **30** and **90** days of being discharged



→ TAP enabled us to efficiently conduct iterative experiments in a distributed environment!

# Contribution of Medication Modeling with Topic Modeling



Including LDA-derived medication features resulted in a 15% improvement!

