

H₂O

WORLD
2 0 1 7



Sergei Izrailev

Chief Data Scientist @ Beeswax

@sizrailev / bit.ly/MLatScale / sergei@beeswax.com / www.beeswax.com

BEESWAX 



H₂O
WORLD
2017

Design Patterns for Machine Learning in Production

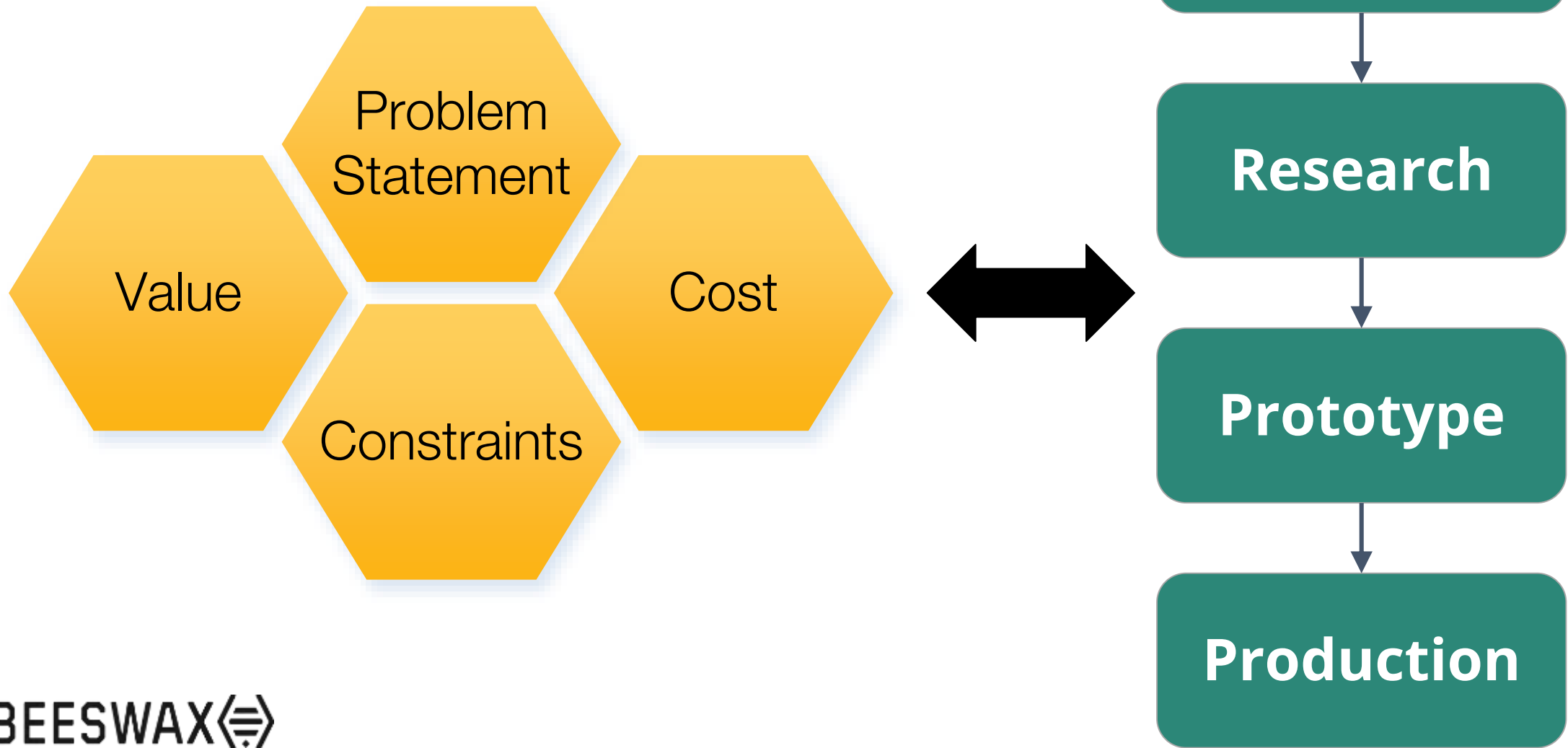
Motivation

- A widespread need to leverage in-house data
- Data science expertise is available
- And yet, it seems to be too hard to extract value from ML

About

- Beeswax
 - Beeswax is an ad tech startup; 40 employees in NYC and London
 - Founded by three ex-Googleers
 - Real-time bidding (RTB) platform for buying online ads (1M+ QPS)
 - Platform tailored for customers to leverage in-house data science
- Myself
 - Production AI systems in Pharma, Finance and Ad Tech
 - Interested in both technology and organizations
 - bit.ly/MLatScale

Overall process



Start with defining the problem

Problem statement

- Is this the right problem to solve?
- Suppose, we've solved the stated problem - what's the value?
- Is ML the right tool to solve the problem?
- What are the constraints?

Define constraints

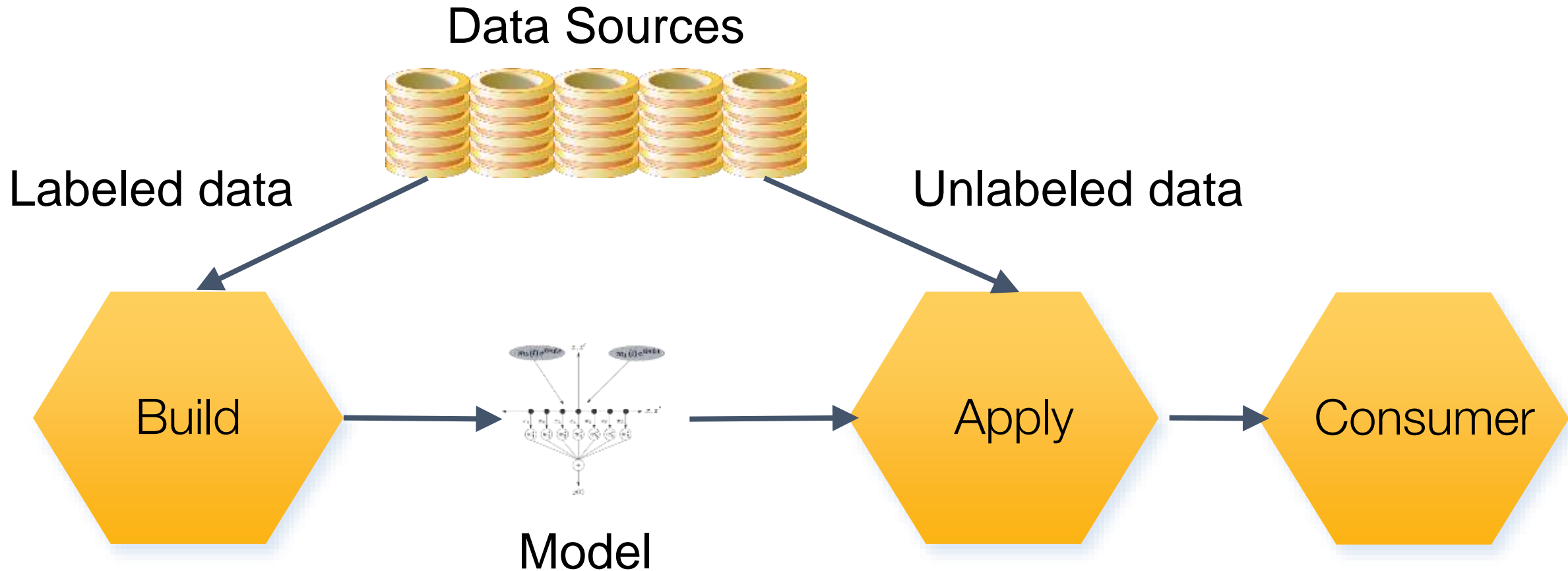
- Existing production environment architecture
- Technology stack
- Available people and their skills
- Requirements for scale

Dimensions of ML system scalability

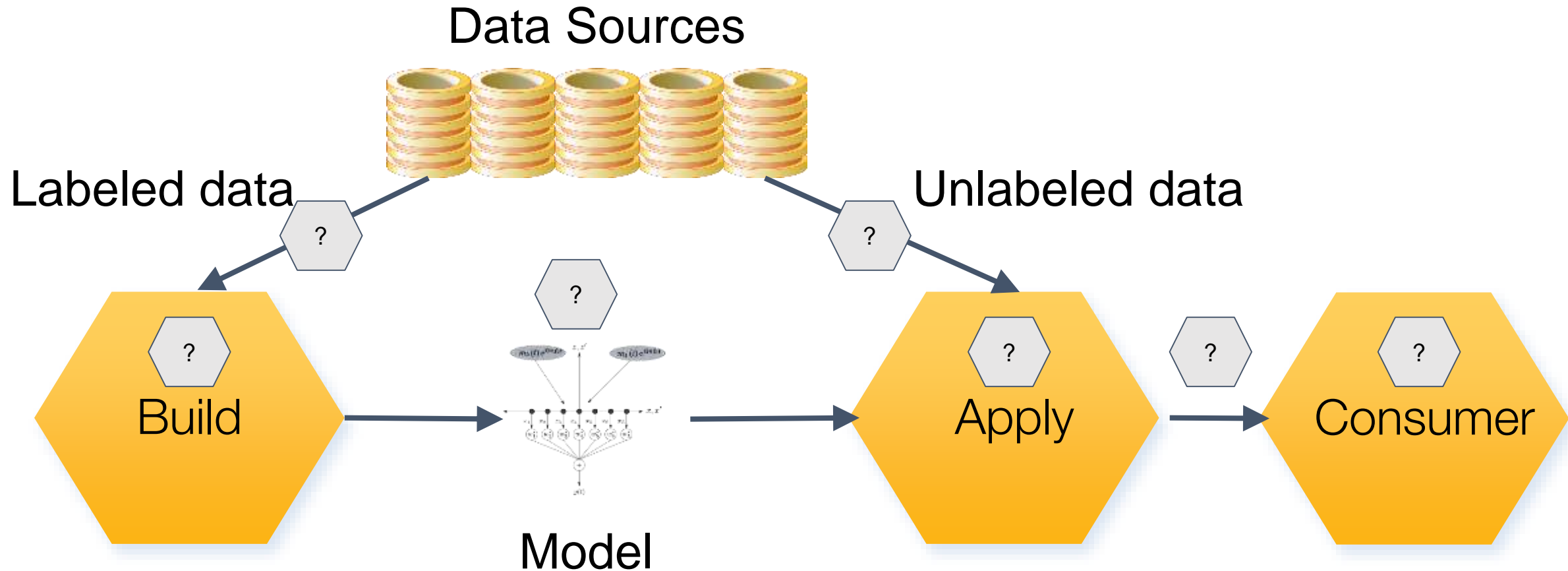
- **Volume:** how much data do we need to process?
- **Velocity:** how quickly does the data change?
- **Variety:** what are the types of data, models, and applications?
- **Veracity:** how accurate are our models?
- **Value:** how does it matter to the ML consumer?
- **Viability:** do the benefits outweigh the costs?

Technical Design of ML Systems

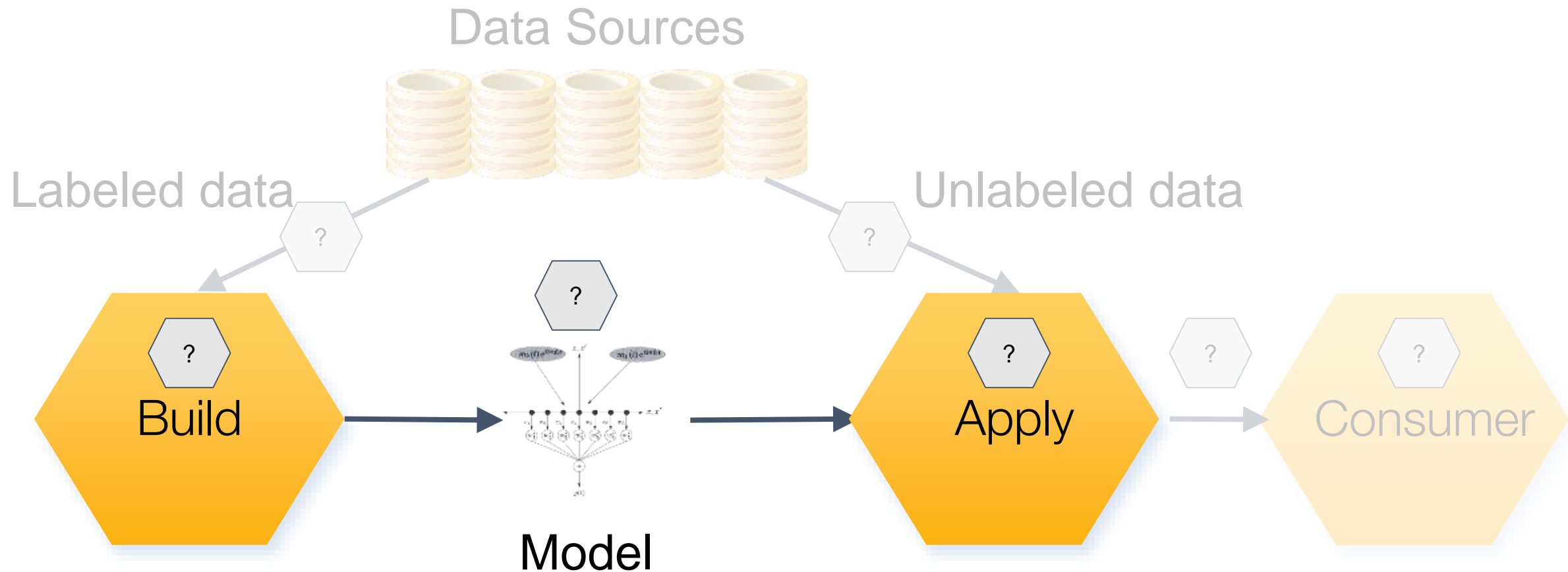
Machine learning systems



ML system design

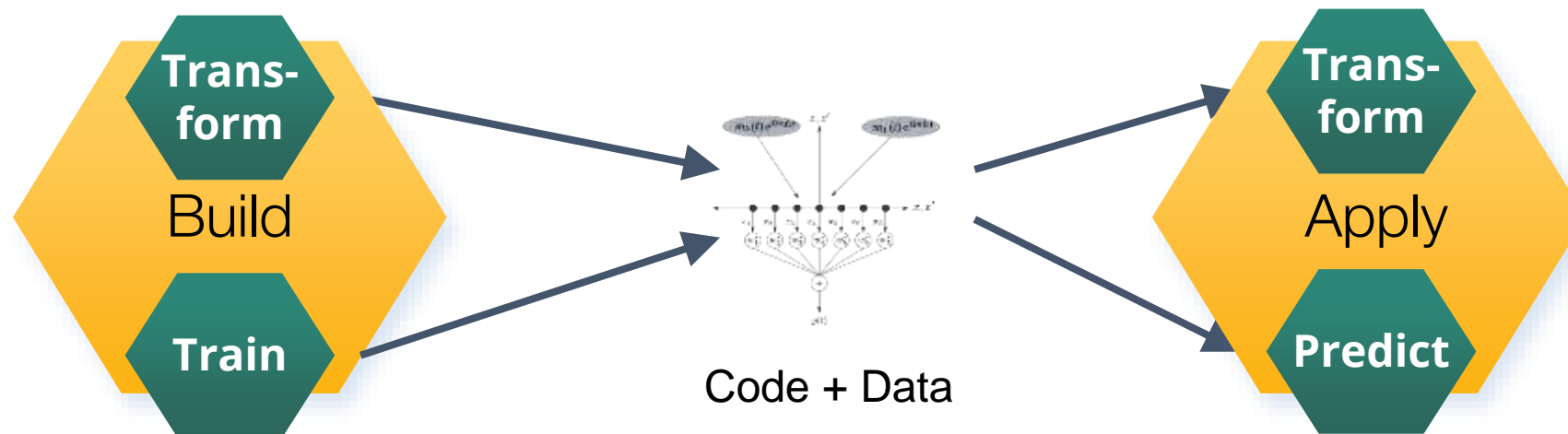


Model deployment



Model deployment

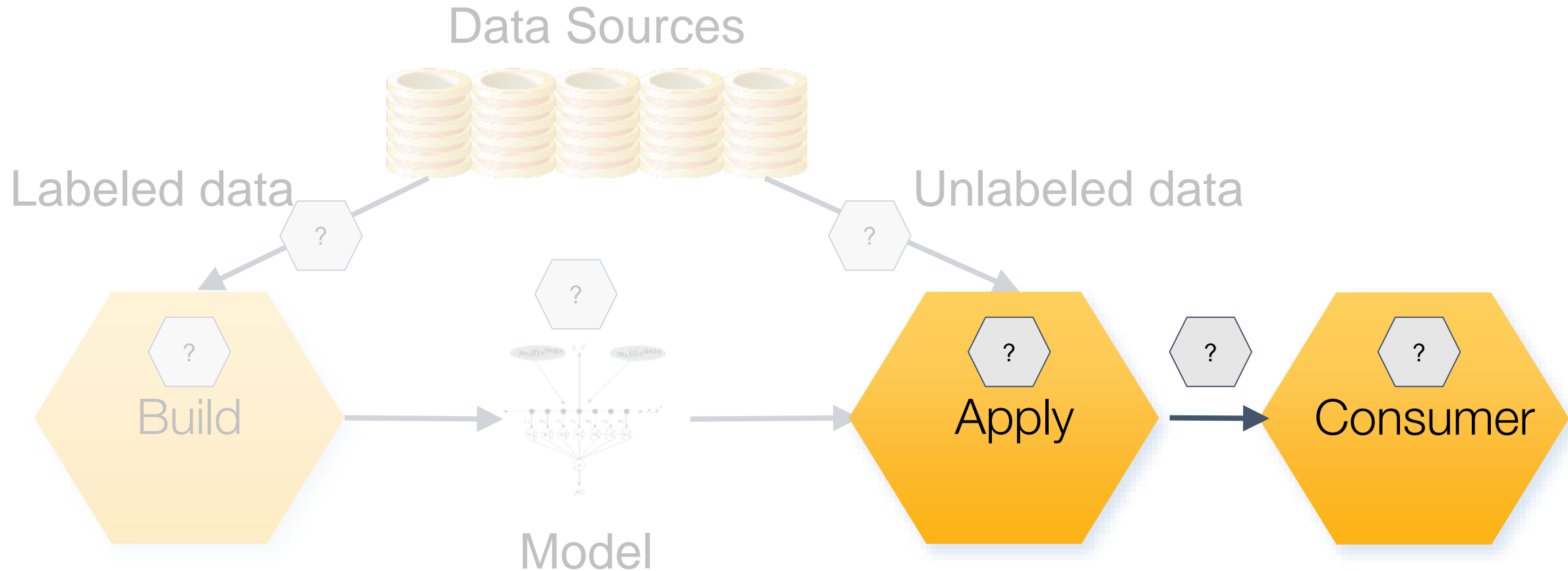
- Data transformations must be the same in training and scoring
- Some transformations are “models” (PCA, top N, TF-IDF)
- Hence, most ML pipelines are DAGs
- These DAGs must be reproduced in production scoring



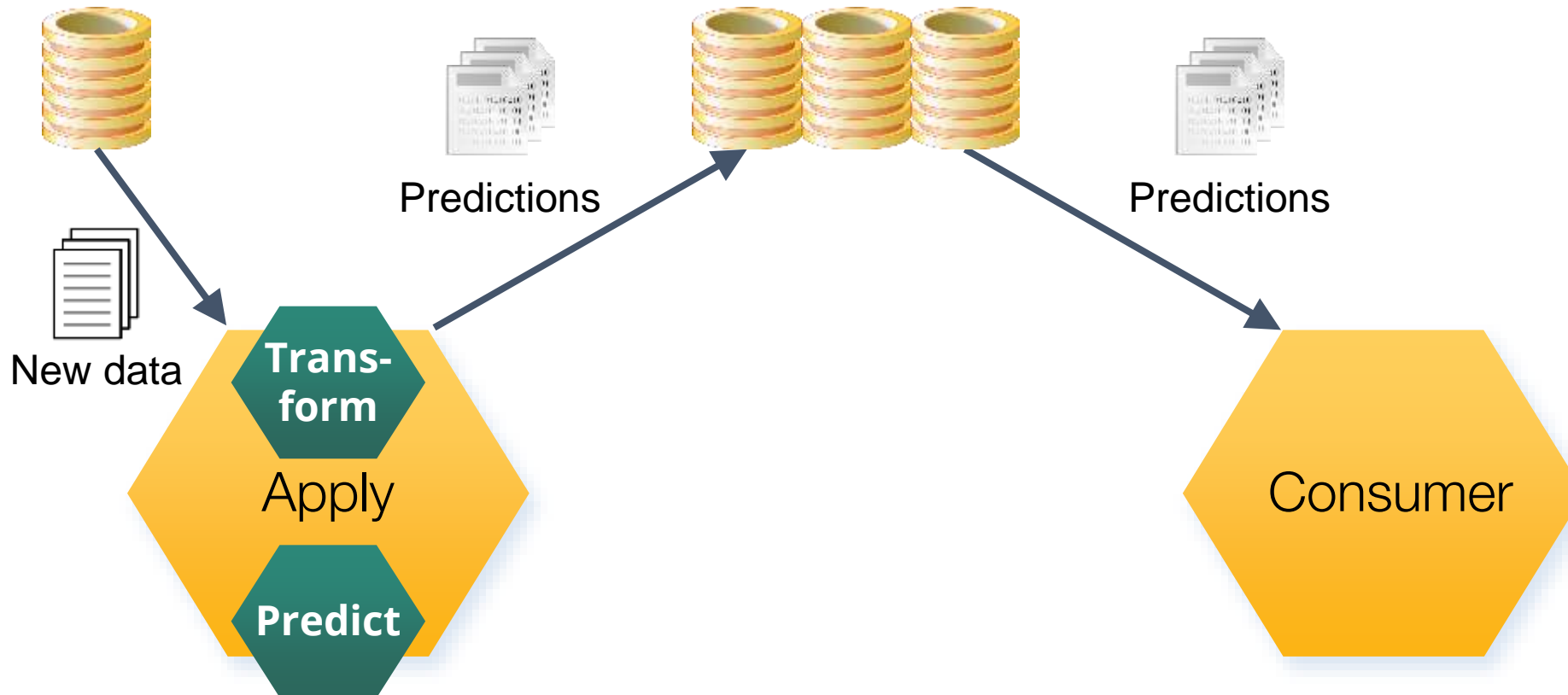
Interface between building and scoring

- In-memory - model is never persisted, train then score
 - single application, also streaming
- Data only - linear coefficients, PMML, etc
 - code is independent
- Serialized objects - Pickle, R, Spark, custom
 - reuse code
- Code + Data - e.g., H2O's POJO
 - code is generated

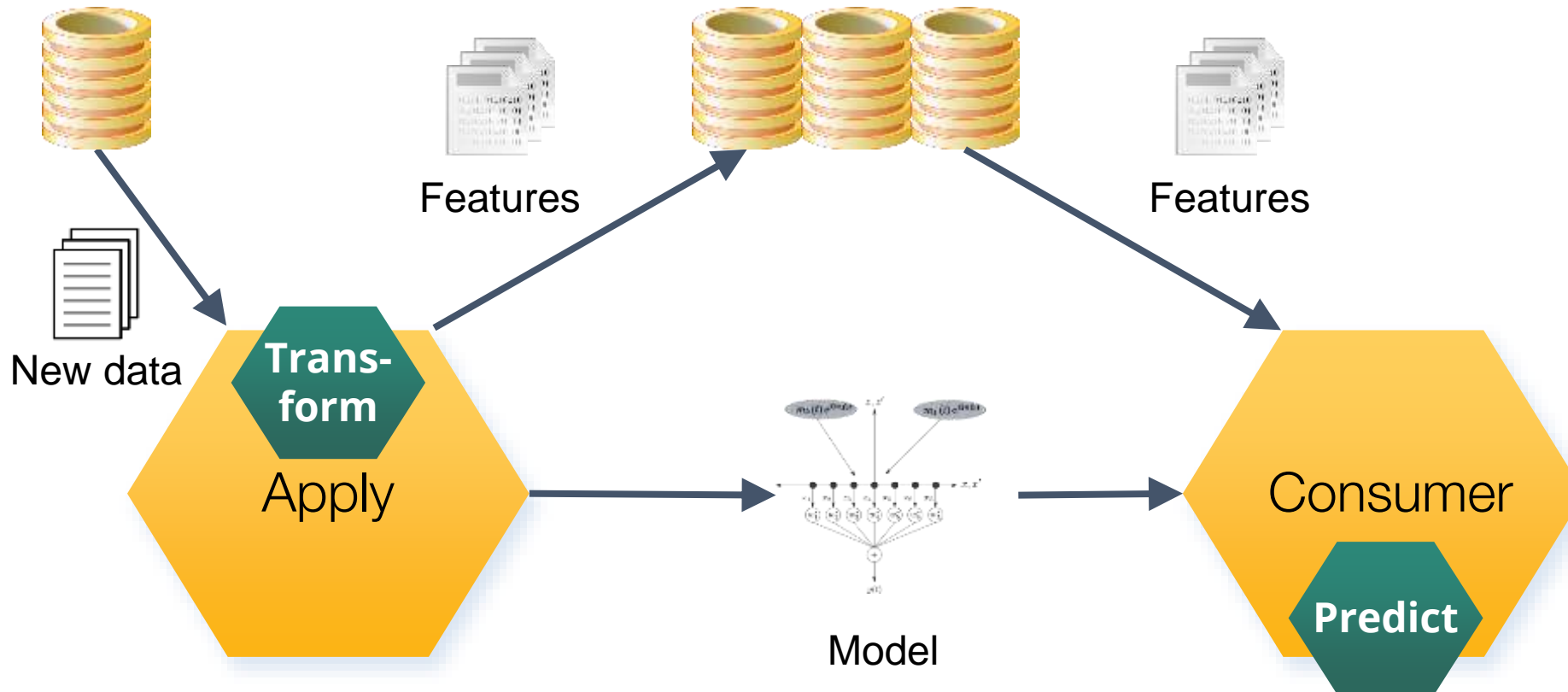
Scoring systems



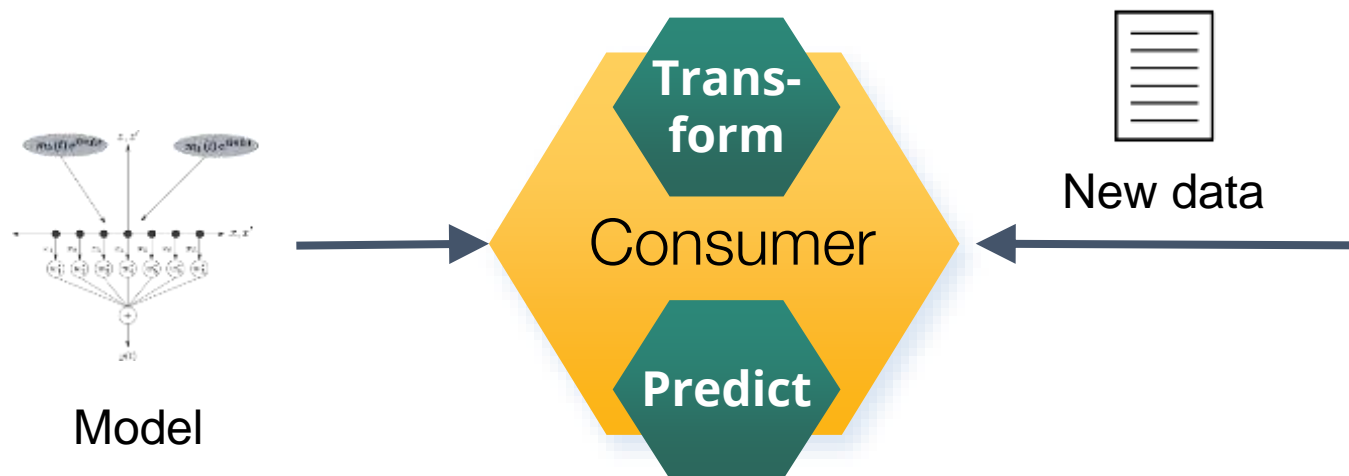
Batch processing



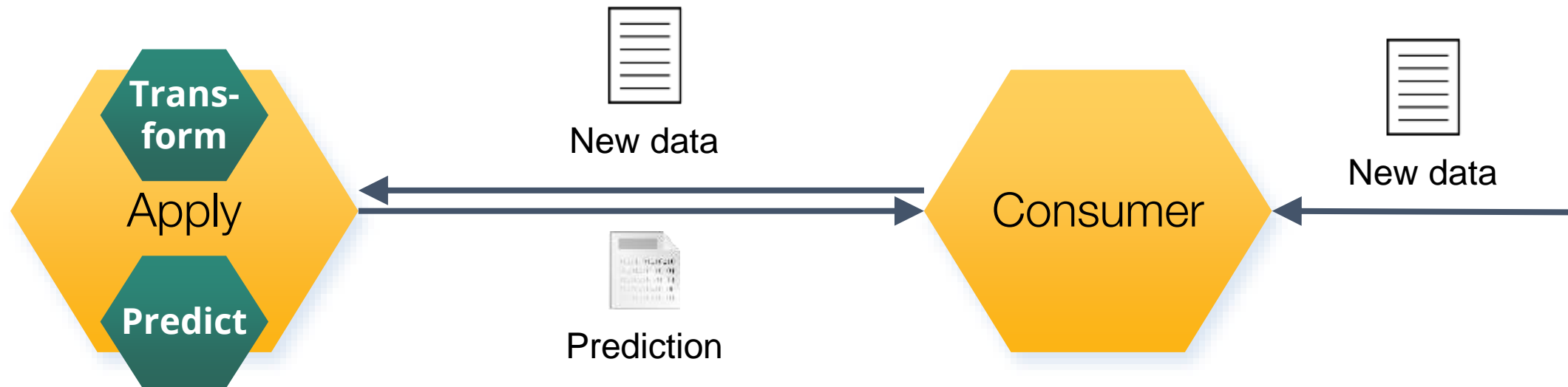
Batch features; consumer predicts



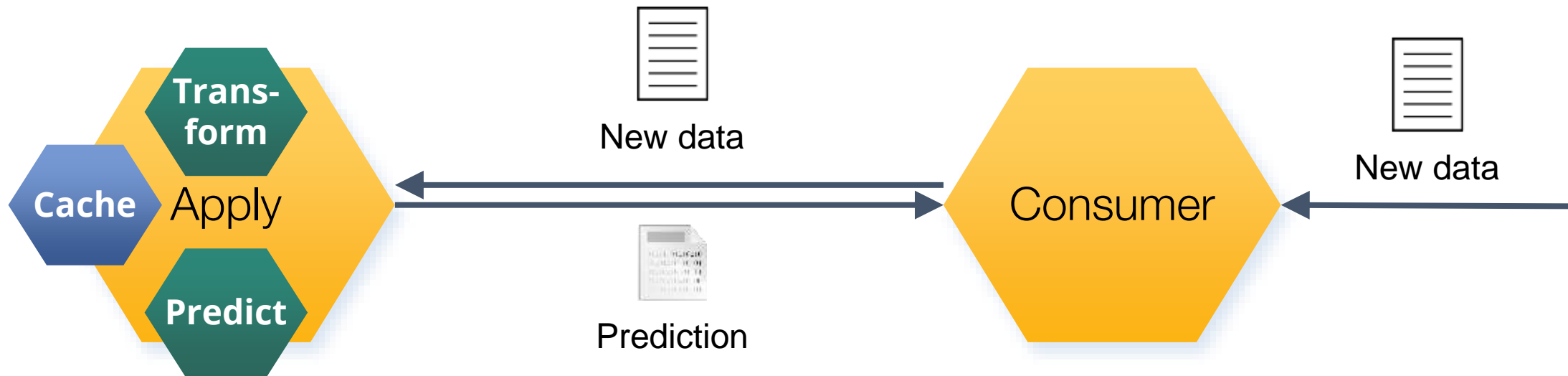
Single-row predictions



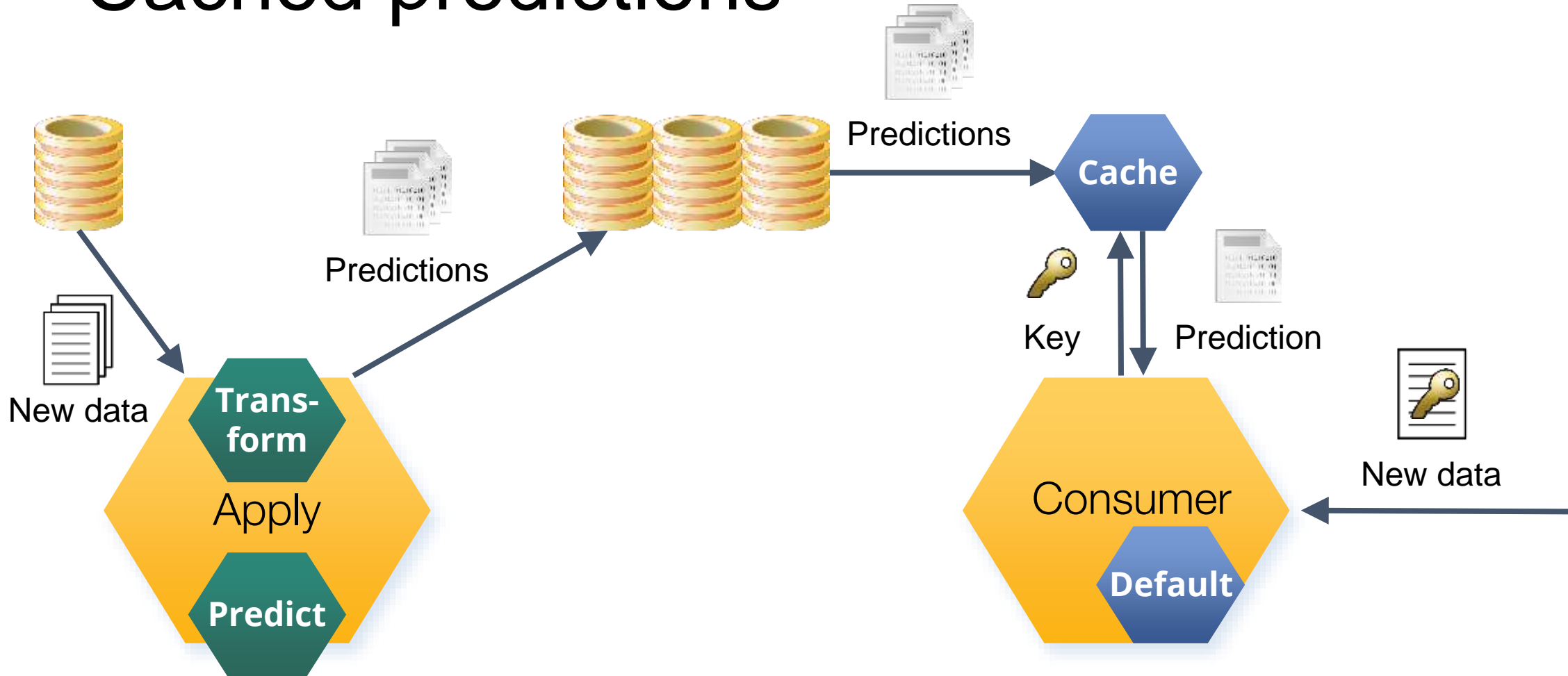
A service for a single row consumer



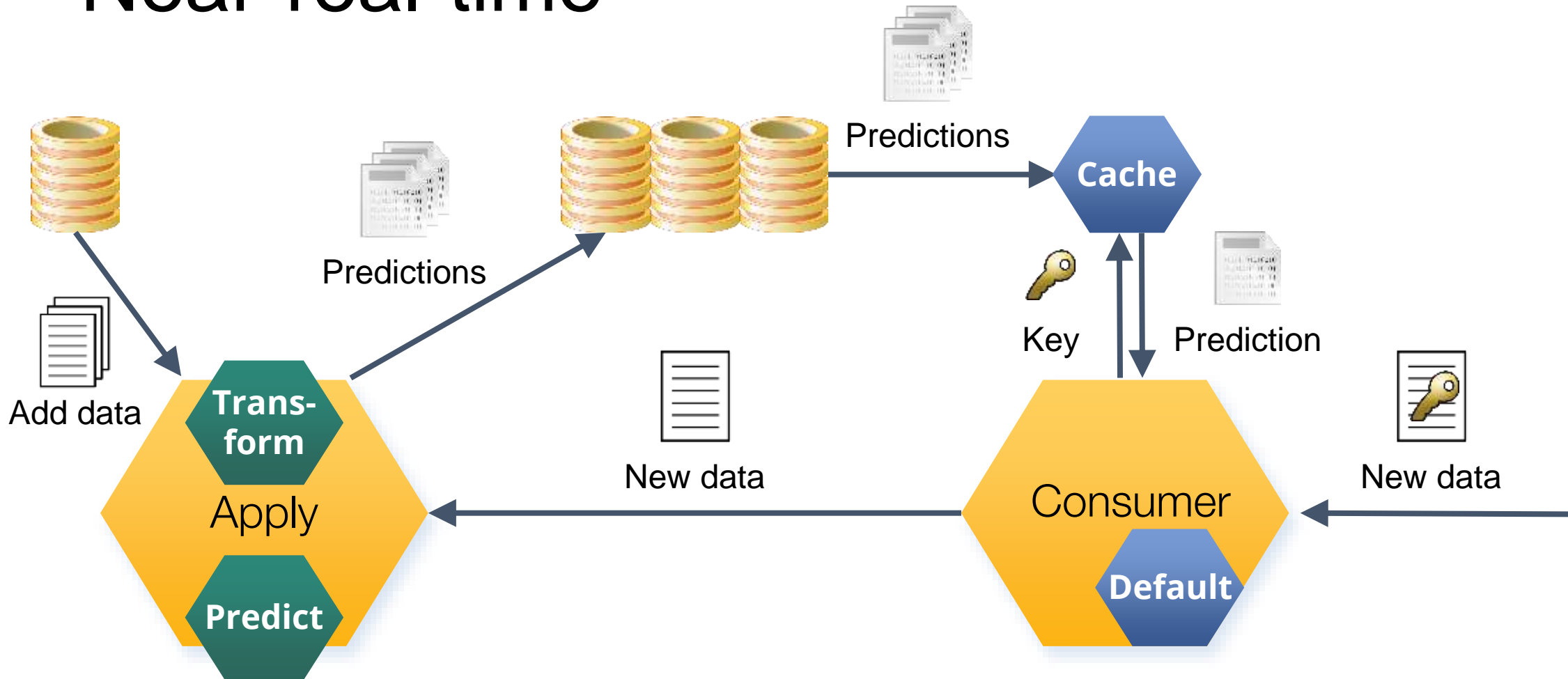
A service for a single row consumer



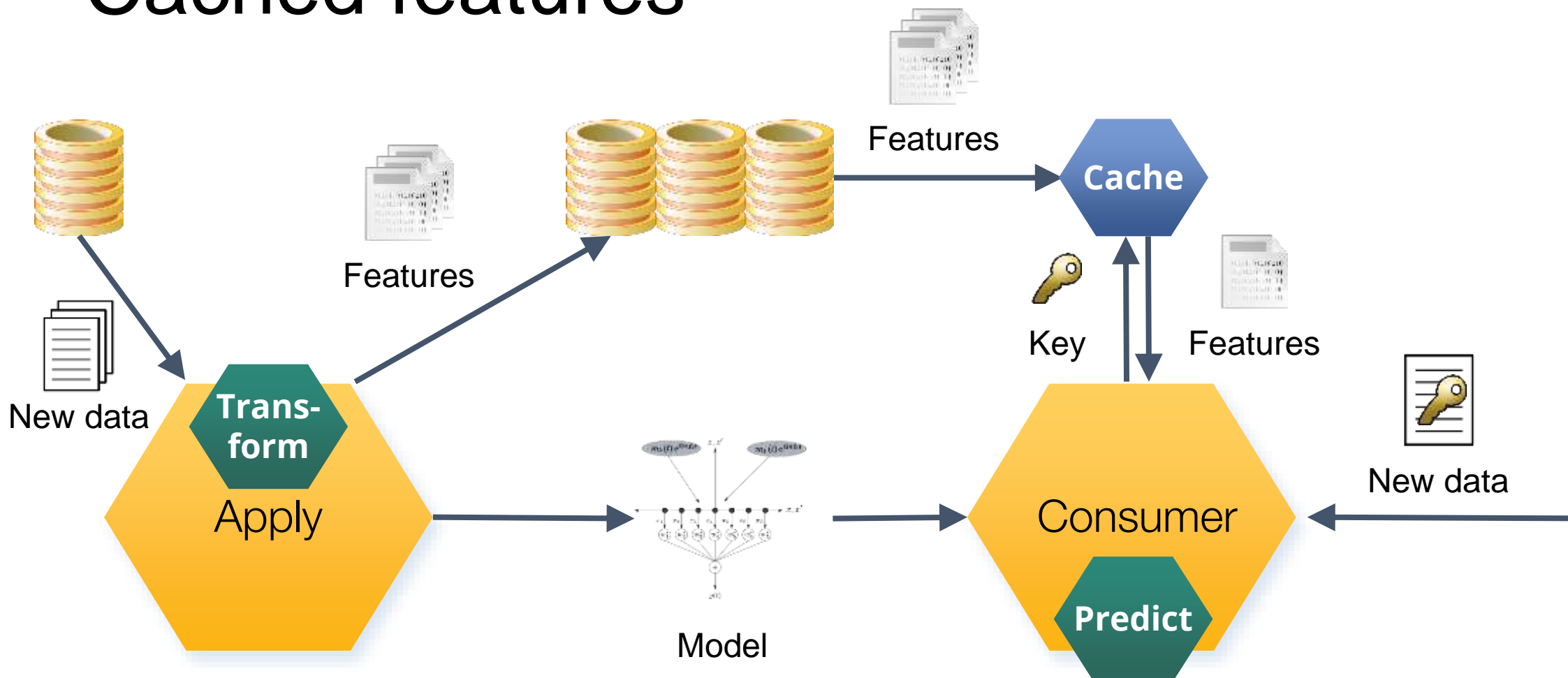
Cached predictions



Near-real-time

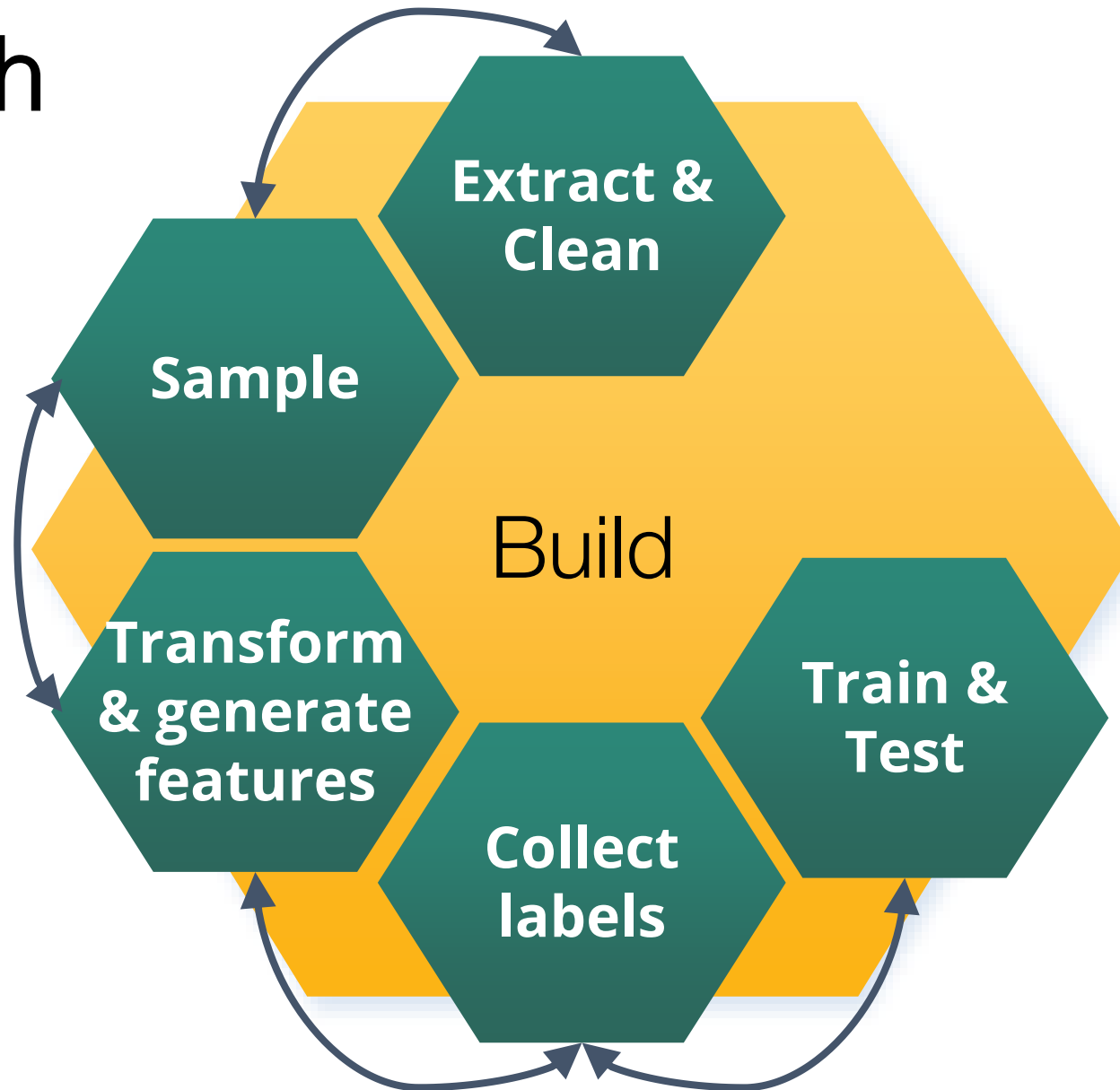


Cached features

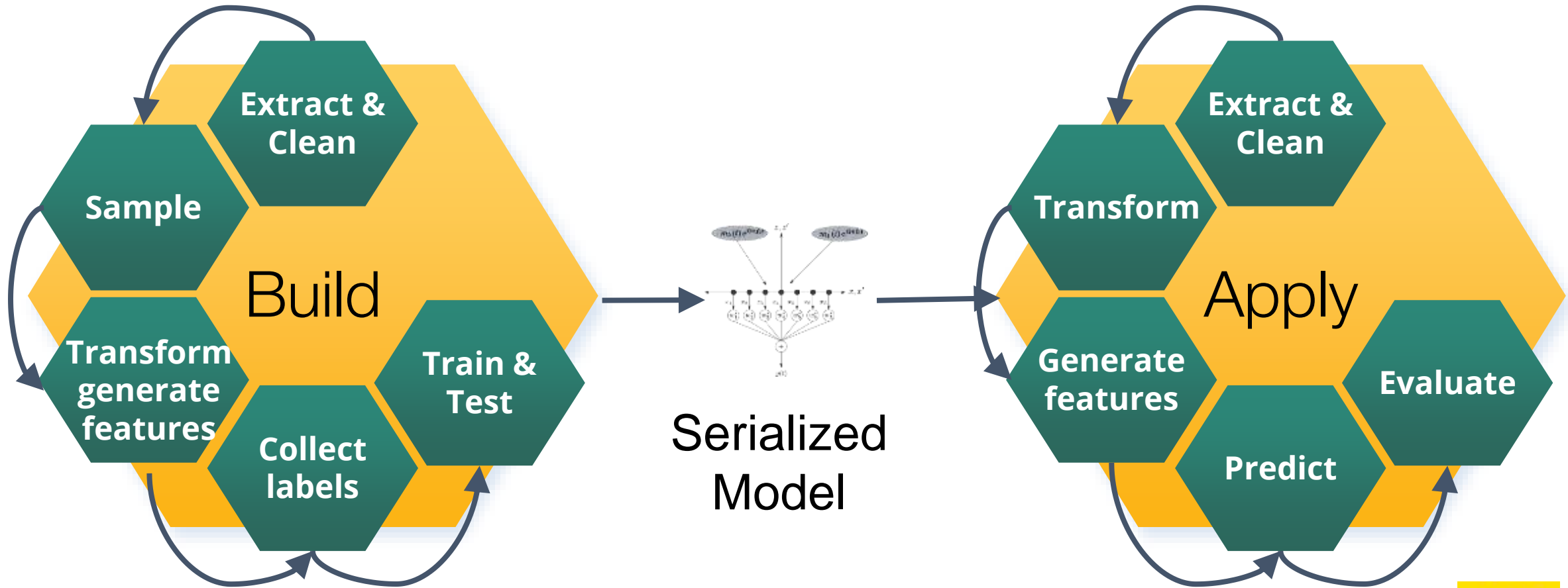


Evolution of ML systems

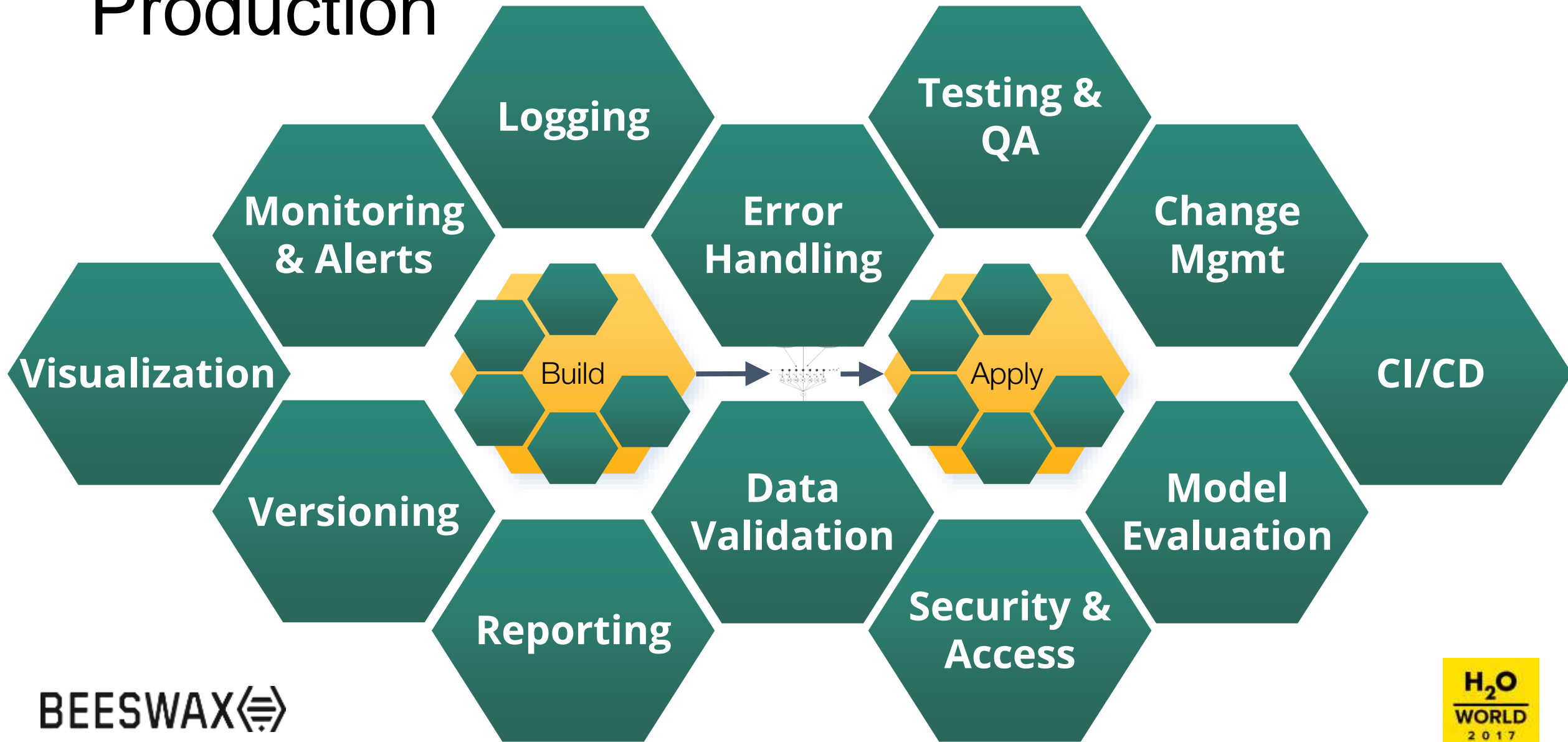
Research



Prototype



Production



Production Fault-tolerant systems

- We should expect problems
 - models don't converge; input data changes; bugs
- Produce acceptable results - even when something fails
- Handle predictable error conditions automatically
- Minimal human intervention
- Easy diagnostics and recovery in case of fatal errors

Conway's law:

"Any piece of software reflects the organizational structure that produced it."

People Questions

- Who is developing training?
- Who is developing scoring?
- Who is responsible for training in production?
- Who is responsible for scoring in production?
- Who deploys new models and model updates?
- Who is responsible for quality control?

People's Functions

Product management

Data science

Data engineering

Application engineering (RT
server-side applications, client-
side applications)

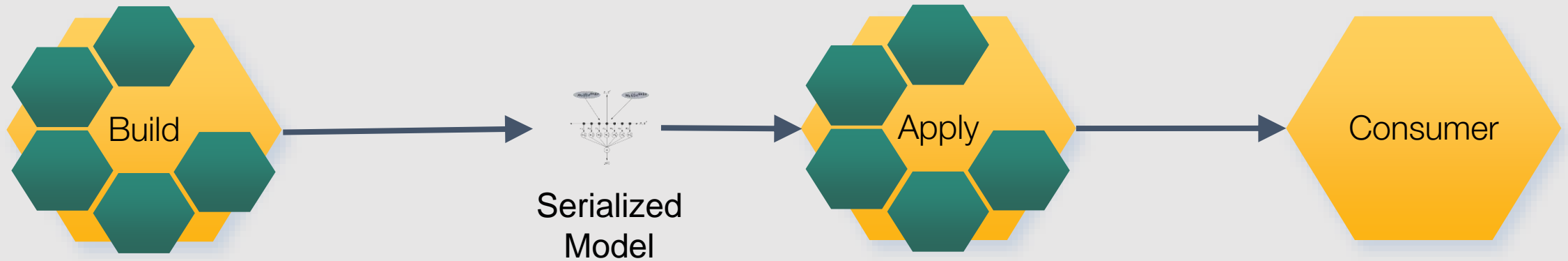
UX

BEESWAX $\langle \equiv \rangle$

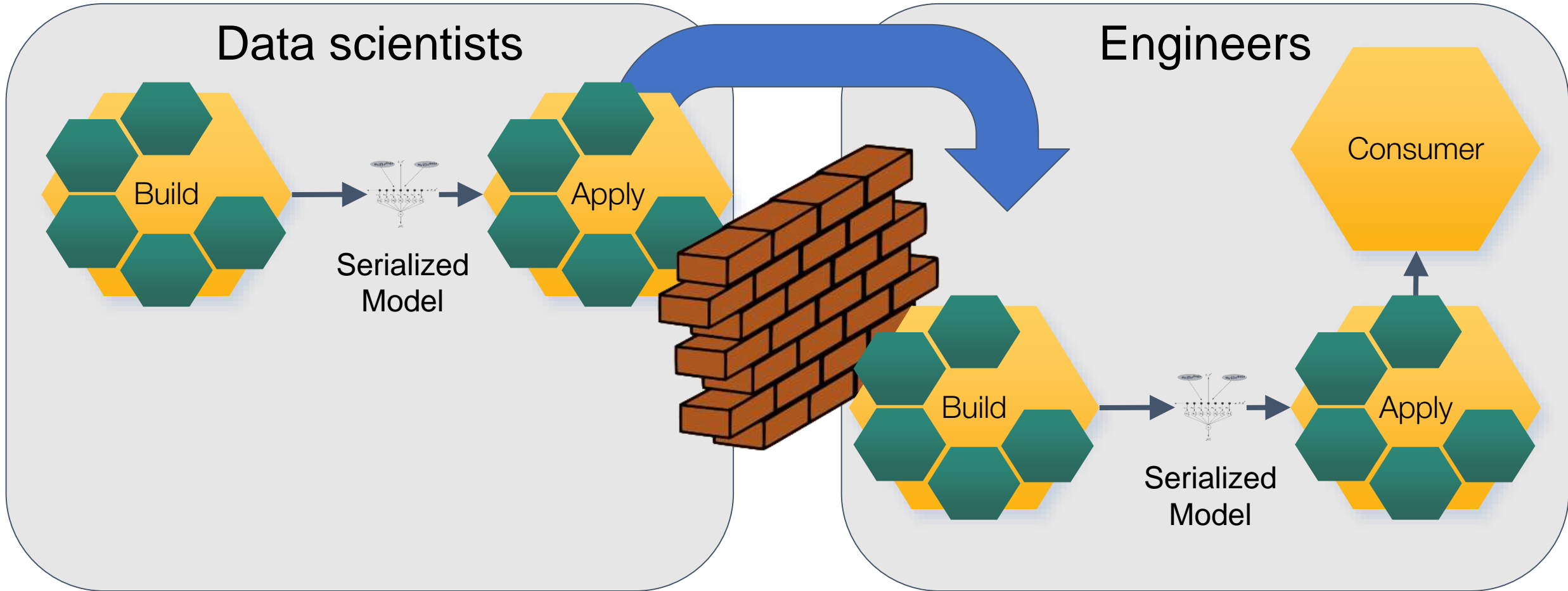
- Front-end development
- Data collection (e.g., logging)
- Code deployment
- Testing and QA
- Infrastructure provisioning
- Support

Data scientists as consumers

Data scientists

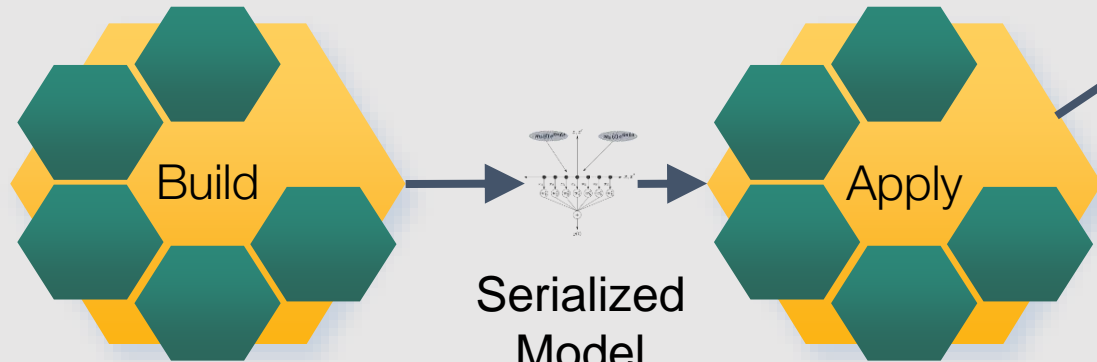


Over the wall



Deliver predictions only (aka “black box”)

Data scientists



Serialized
Model



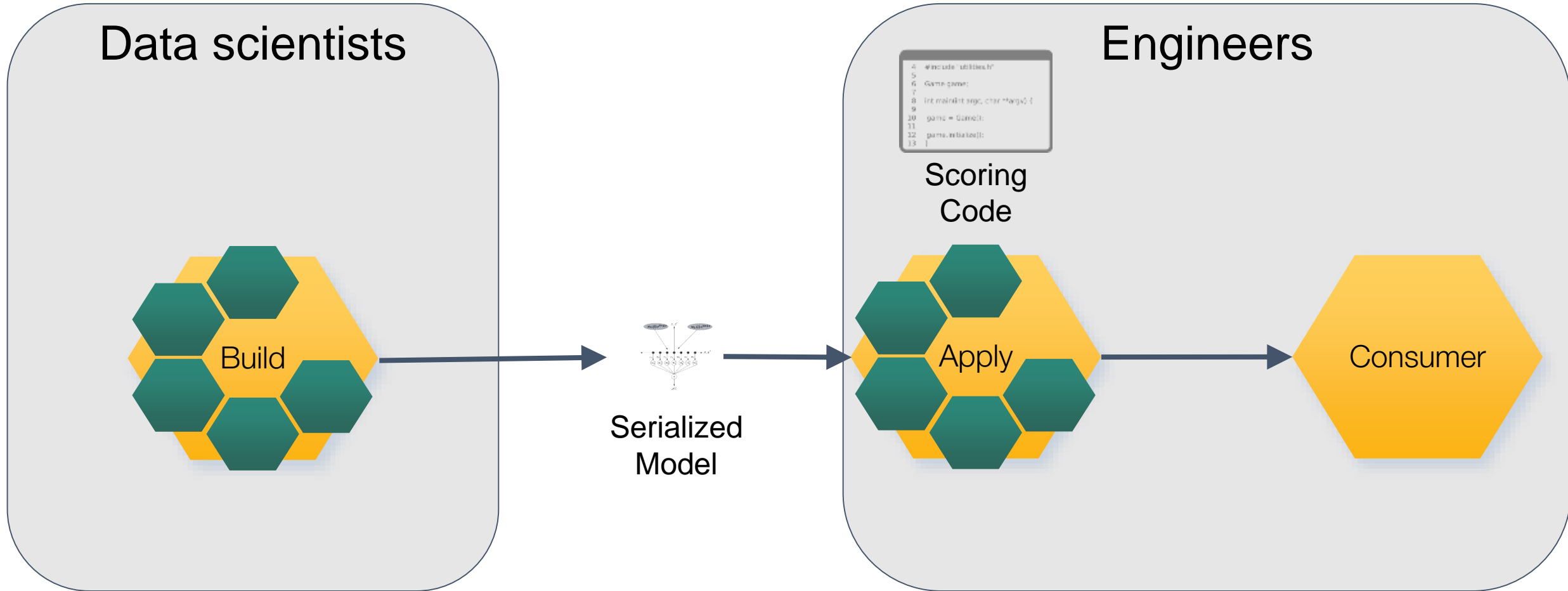
Predictions



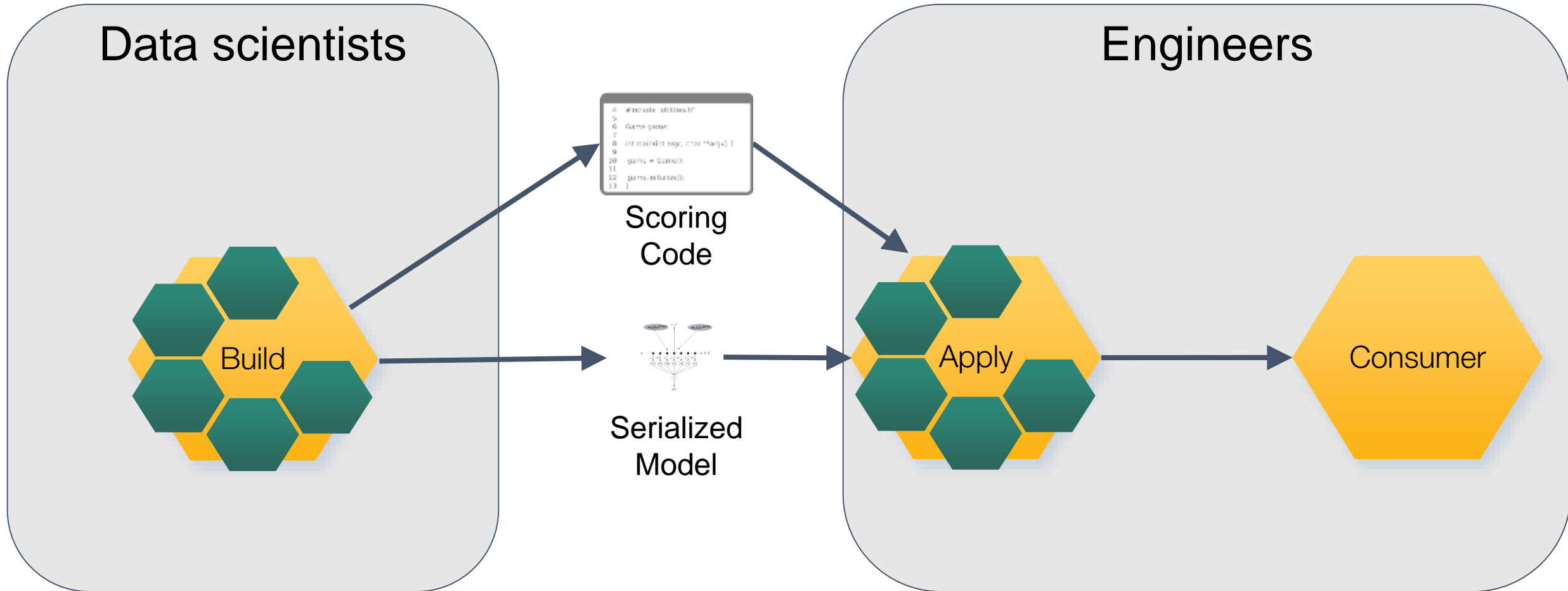
Engineers

Consumer

Deliver data (serialized model)

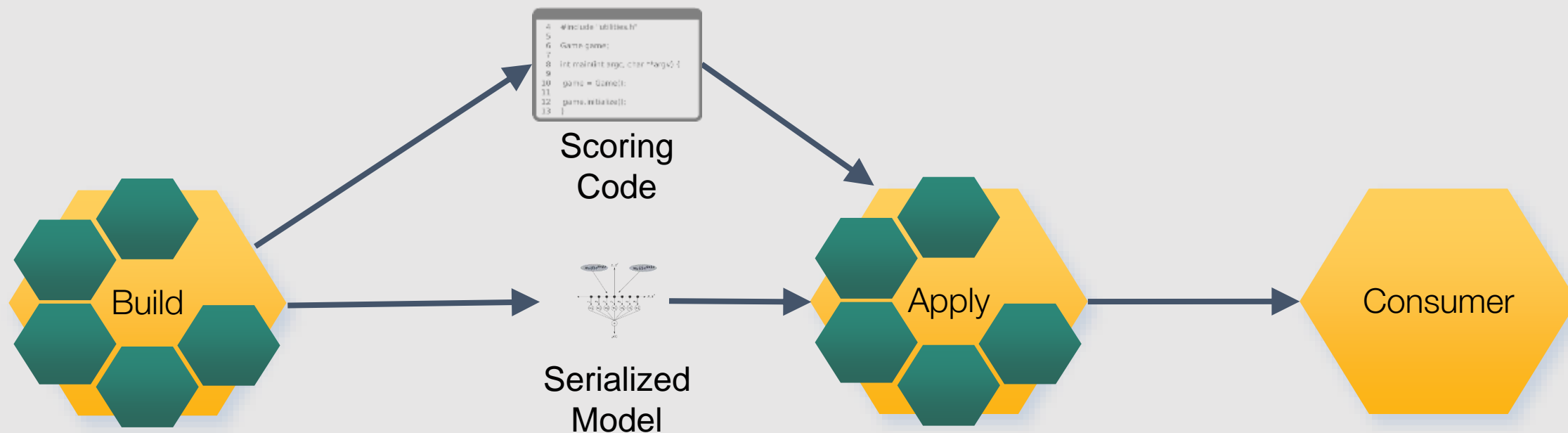


Deliver code and data

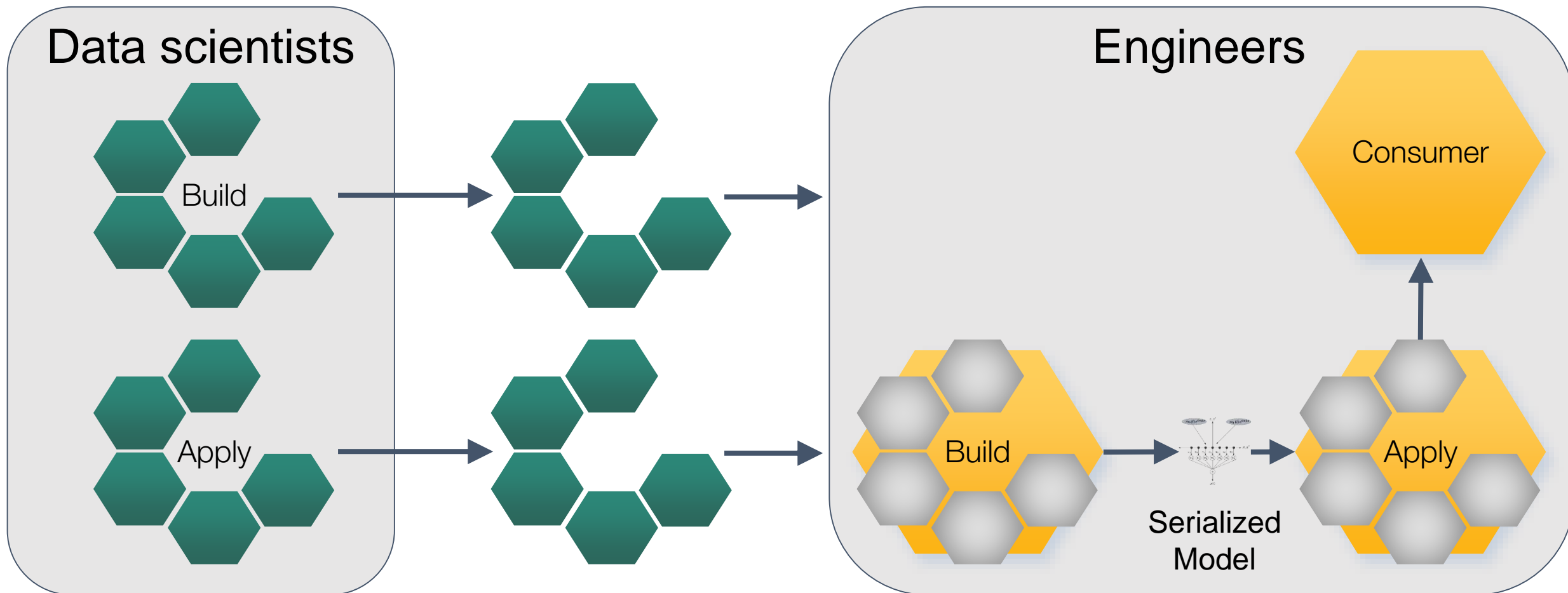


Cross-functional team

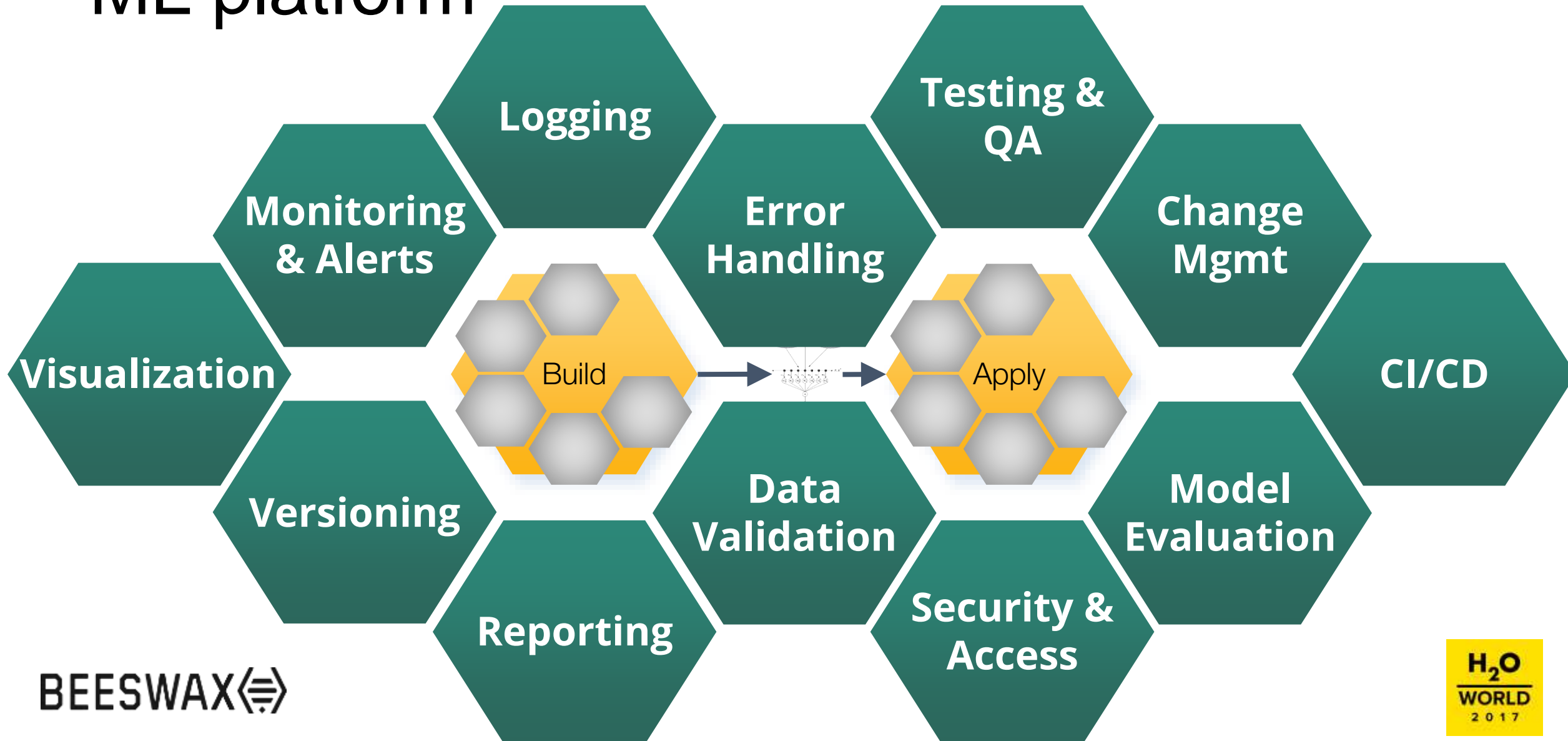
Team: Data Scientists + Engineers



ML platform



ML platform



Conclusion

- Find the right problem
- Define constraints
- Design components and interfaces
- Take into account organizational constraints
- Production can't be an afterthought
- The process is a lot of work, but it's not rocket science

Questions?

Yes, we are hiring...

sergei@beeswax.com