

⇒ } Imagine you were working on iPhone. Everytime users open their phones, you want to suggest one app they are most likely to open first with 90% accuracy. How would you do that? }

## ✓ Functional Requirement

1. Real time prediction — 100ms of phone unlock
2. Personalised recommendations → Moving app usages /  
Adopt with user habit every app usage
3. offline availability → offline fallback → cache prediction  
↳ 24 hrs on-device
4. Privacy → ~~GPS~~, user-id hashed

## ✓ Non-Requirements

1. latency low (Unlock → prediction) < 100ms
2. Scalability → Handle 10Mn daily active users  
(1000+ req / sec at peak)
3. Availability → 99.99%
4. Security → Encrypt }
5. Cost → ↓ 0.001 (Sagemaker, FCS, EC2) }

## Data Ingestion & Storage Layer.

## 1. User Behaviour

- Apps during diff times
  - Apps opened after unlocking
  - Time since last unlock
  - session duration
- (food → 8:17 PM }  
lunch → == }

## 27 Contextual Signals

- Time of day, Location  
Device Status (ios version, Battery, Wifi, Cellular data)  
↓  
OS  
Storage space ...

## Historical Patterns

- † Most used app last 24 hours, week ...
- † frequency of app switches (in messages, slack)
  - ↑
  - ↑

Storage (Ans)

Real-time data  $\rightarrow$  Kafka / Amazon Kinesis Data Streams

Raw data archives , processed S3

(partition vsig date/  
hour)

↓  
Give Data  
Catalog.

Real-time

[ User\_123 → 8:30 AM ]

Batch

↳ PySpark SOL → Ranking

time spent on the app  
over user id, hour

S3 → store batch features

store real-time info (key-value)

+ user 8:00 AM  
+ user 9:00 AM

(DynamoDB)

<https://aws.amazon.com/dynamodb/>

## Machine Learning Pipeline

90%

1. Data Versioning → Track datasets S3 object versioning

2. Training → Sagemaker → ml. c5.2xlarge (spot instance)  
+ Hyperparameter tuning (Dept... '08')

3. Validation

+ Holdout validation (30 days time-series split)

+ AUC-ROC

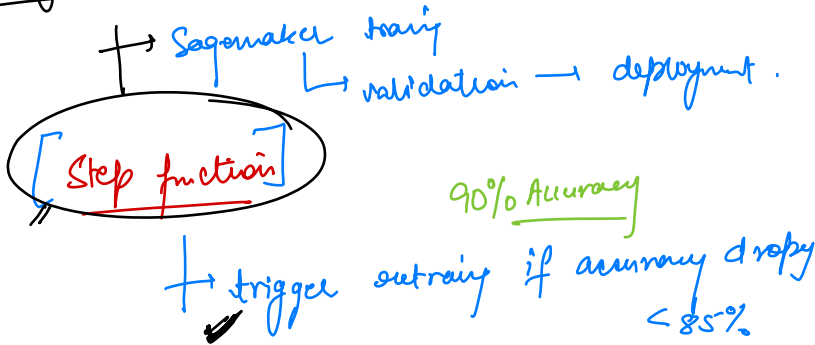


4. Model Registry

+ Version model in Sagemaker Model Registry

Production / Staging

## 5. Deployment



<https://aws.amazon.com/step-functions/>

Model → ML model selection

① LightGBM

- < 10ms
- interpretable
- Stop values

② XGBOOST

- < 10ms
- interpretable

→ sequential data (Struggles) → app usage history.

③ Transformer Models

→ app sequences (calendar → [ ] → Zeroth  
↓  
Divert)

→ GPU inference  
→ latency could be higher

+ Complex patterns

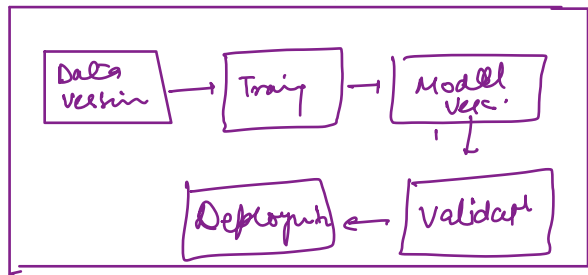
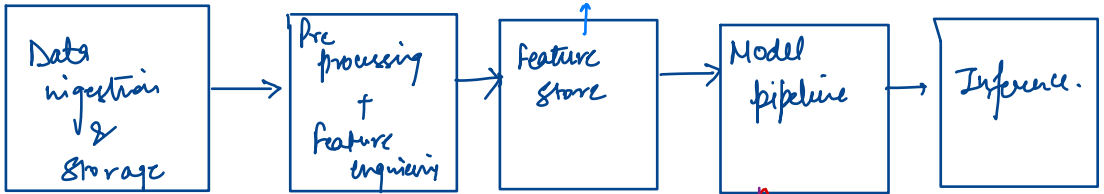
## Optional Extension

A/B testing → Sagemaker AB Testing  
↳ 5% traffic to challenger model

<https://aws.amazon.com/blogs/machine-learning/a-b-testing-ml-models-in-production-using-amazon-sagemaker/>

Canary deployment  
↳ 1% users in specific region

Federated Learning →





Sample req → { user id: 123  
features : { hour: 8, location: new delhi,  
last-app-used: Calendar } }

{ battery: 75% , network: - }

{ slack: 10 }

H/W ↴

An e-commerce company is trying to minimize the time it takes customers to purchase their selected items. As a machine learning engineer, what can you do to help them?