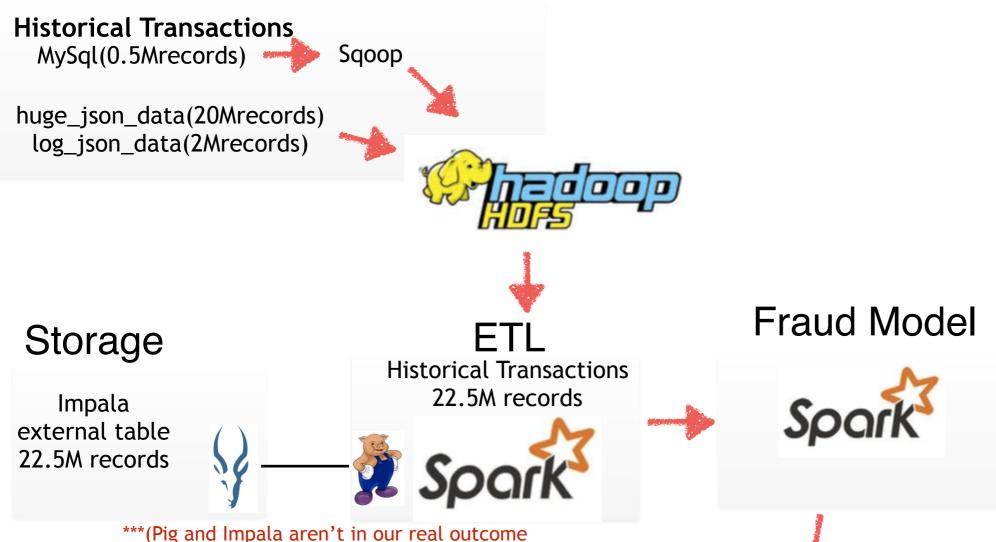# Near Real-time Fraud Detection Framework

Sirapat Na Ranong 57070503438
Chaleamkwan Bunyarattanamongkon 57070503444
Saket Khandelwal 57070503483

# Diagram/Workflow

## Ingest to HDFS

**Historical Transactions**
MySql(0.5Mrecords) → Sqoop

huge_json_data(20Mrecords)
log_json_data(2Mrecords)

## Storage

Impala
external table
22.5M records

## ETL
Historical Transactions
22.5M records

***(Pig and Impala aren't in our real outcome
but we have used it for study cases)***

## Fraud Model

## Alert System
(Pushbullet&Gmail)

## Real-time Transactions

Stream.py

## Fraud Detector

# Running Process Examples (Spark[ETL&Model],Pig,Impala)

# Result/Demo

| Friends | Me | Following |

FRAUD DETECTED!!!! Your account ID539 has withdraw 199348 baht (It is above 107420.943035 )

**Fraud Detected**
FRAUD DETECTED!!!! Your account ID185 has withdraw 159545 baht (It is above 107117.141358 )

**Fraud Detected**
FRAUD DETECTED!!!! Your account ID652 has withdraw 188350 baht (It is above 107721.530259 )

**Fraud Detected**
FRAUD DETECTED!!!! Your account ID487 has withdraw 151408 baht (It is above 107817.834739 )

21:50

To   ♾ All Devices

📷 Send a message                Send

Pushing     Channels     Settings

Gmail ▾

COMPOSE

Inbox (300)
Starred
Sent Mail
Drafts

(no subject)  Inbox  x

dummyparallel7235@gmail.com                9:50 PM (1
to bcc: me ▾

FRAUD DETECTED!!!! Your account ID173 has withdraw 164948 baht (It is above 107866.402112 )