

André Simões
andre.simoes@xpand-it.com
@ITXpander



GENERAL ISSUES

How busy is my cluster?

Which tools does each department use?

Who are the main cluster users/departments?

Do I need to plan on an upgrade?

How much is process X costing me?

Are there available time slots?

(SOME OF) THE PROBLEMS

ORCHESTRATION IS HARD

MONITORING IS HARDER

NOTIFICATIONS... WELL...



CONTROLLING THE ETL LAUNCH

CENTRALIZED ORCHESTRATION

POOL MANAGEMENT

EXECUTION DELEGATION

ORCHESTRATION METADATA

REMOTE AGENTS

CLUSTER DATA

DEPENDENCIES

SCHEDULING



ETL METADATA

POOL METADATA

USER MAPPING

REMOTE AGENTS

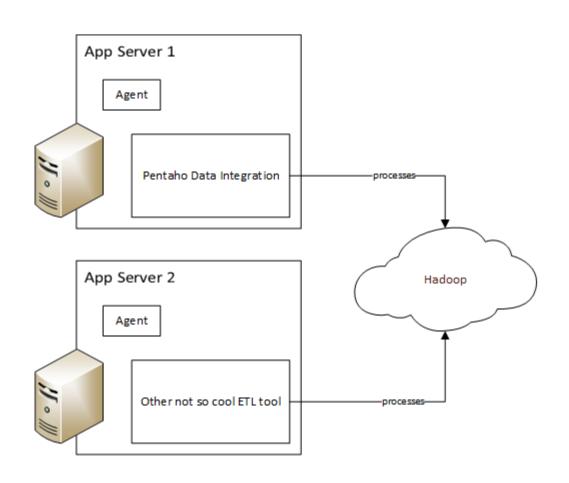
FINE GRAINED CONTROL

ETL TOOL SPECIFIC

REAL TIME LOGGING

ERROR RECOVERY

ASYNC EXECUTION



GATHERING EXECUTION DATA

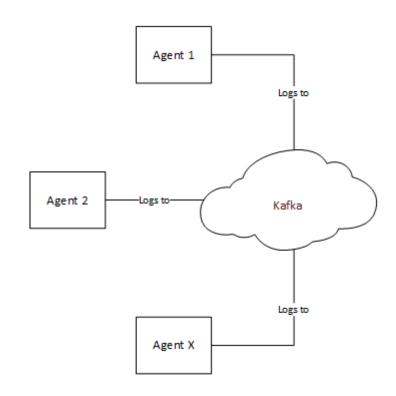
USE KAFKA AS A LOG SINK

FAULT TOLERANT

REAL TIME

CONSISTENT

CONTROLLED BY THE AGENT



PDI EXTENSION POINTS

CAPTURE LOG START

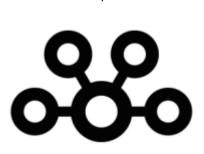
CAPTURE LOG END

CAPTURE CONNECTION TYPE

CAPTURE STEP LINEAGE DETAIL

GENERATES NOTIFICATIONS





COLLECT LOG DATA IN (AS) REALTIME (AS POSSIBLE)

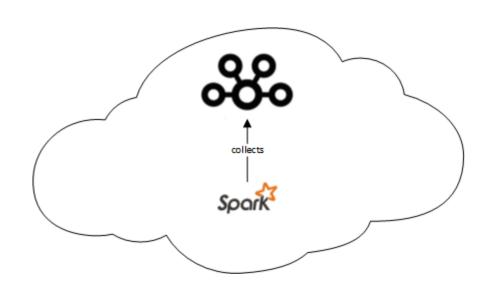
SPARK AS KAFKA COLLECTOR

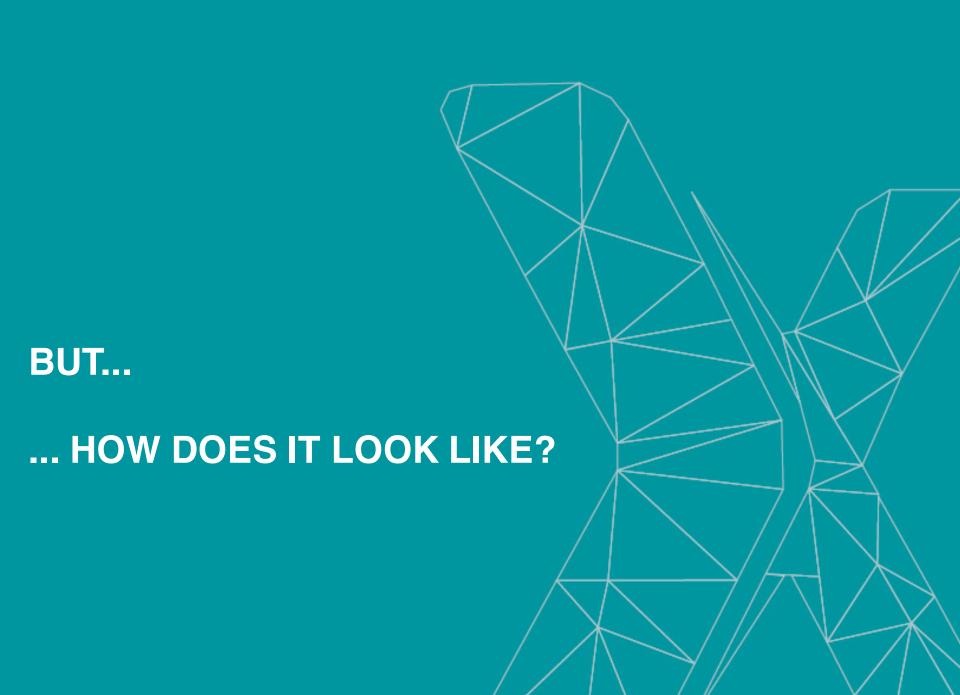
REAL TIME LOG PARSING

ETL TOOLADAPTABLE

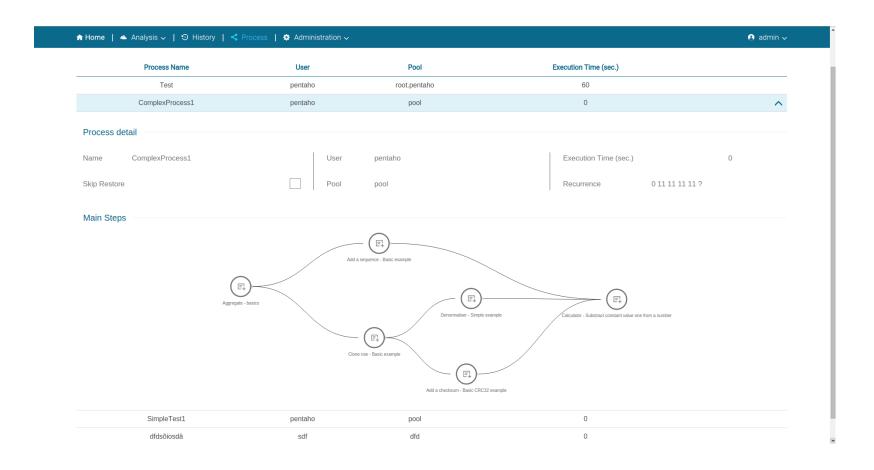
DATA DUMPS IN IMPALA AND HBASE

GENERATES NOTIFICATIONS

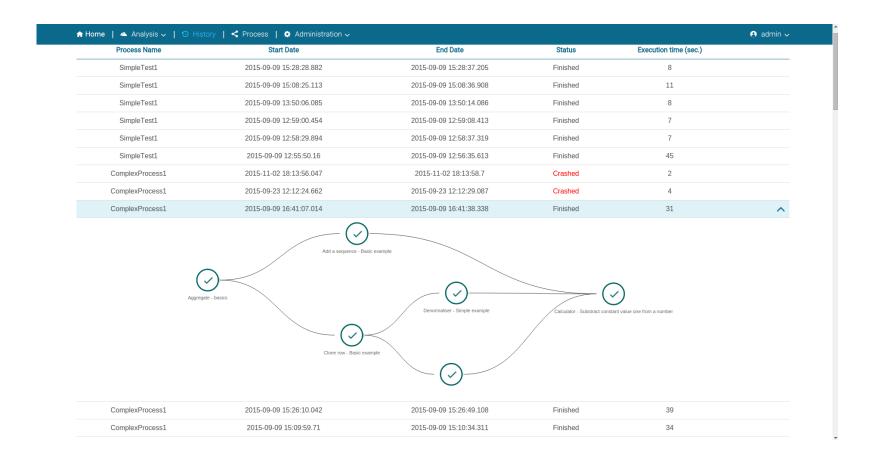




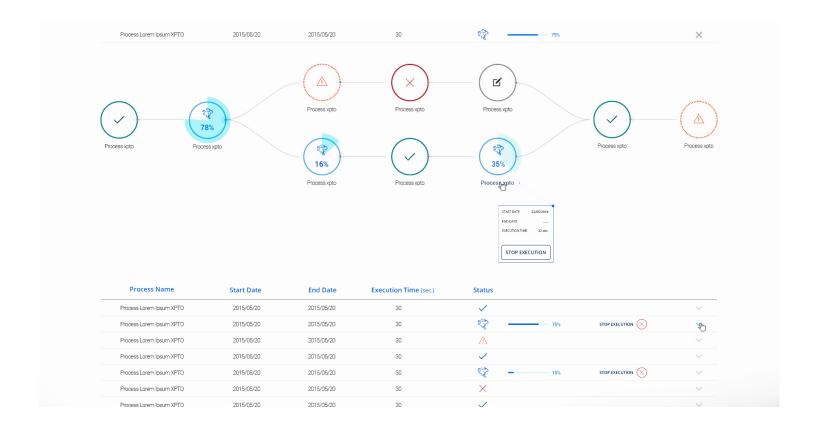
CONFIGURATION



EXECUTION HISTORY



REALTIME (AS POSSIBLE) DATA



ANALYTICS

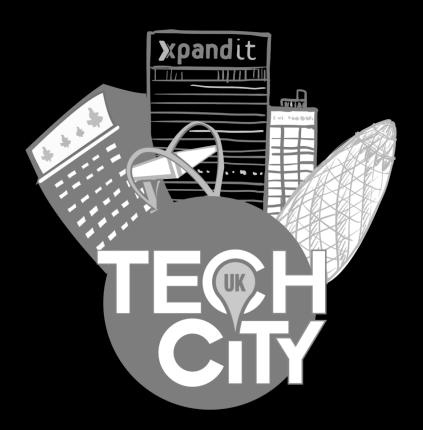
Intentionally left blank.

THE BAD AND THE UGLY

LIMITED POOL MANAGEMENT

MANAGING KERBERIZED CLUSTERS

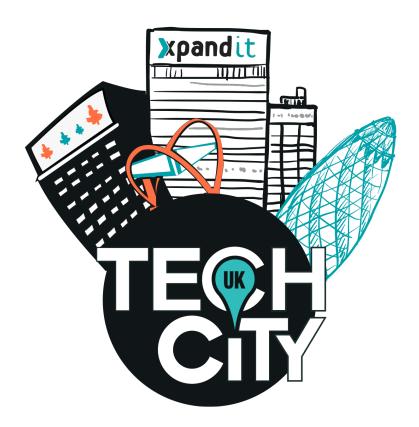
MULTI-POINT DEPLOYMENT



NEW LONDON OFFICE!

THANK YOU!

NEWSFLASH



We are now in London!