



LEVERAGING MESOS AS THE ULTIMATE DISTRIBUTED DATA SCIENCE PLATFORM

(such a long title,) by @DataFellas
@Noootsab, 8th Oct. '15 @MesosCon

However, "*Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb*" is a rather long title, yet the best movie ever (IMHO)

OUTLINE

- **(LEGACY) DATA SCIENCE PIPELINE/PRODUCT**
- **WHAT CHANGED SINCE THEN**
- **DISTRIBUTED DATA SCIENCE (TODAY)**
- **LUCKILY, WE HAVE MESOS AND FRIENDS**
- **GOING BEYOND (PRODUCTIVITY)**



DATA FELLAS

5 MONTHS OLD BELGIAN STARTUP



ANDY PETRELLA



XAVIER TORDOIR

**MATHS
SCALA
APACHE SPARK**

**PHYSICS
BIOINFORMATICS**

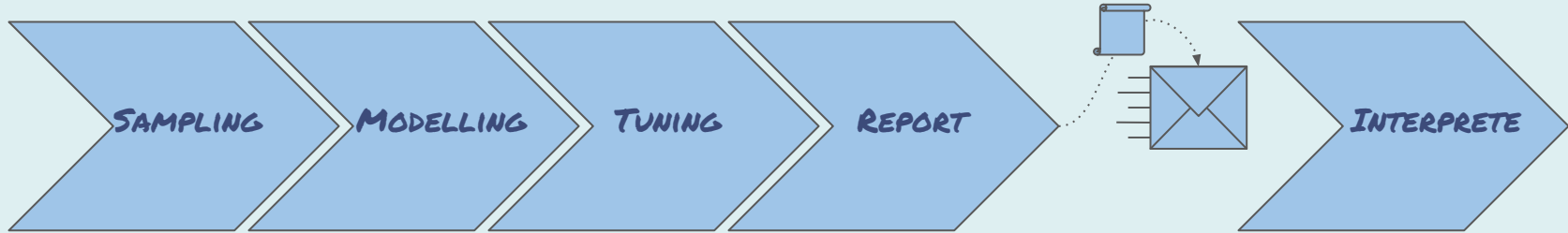
**SPARK NOTEBOOK
TRAINER
DATA BANANA**

**SCALA
SPARK**



(LEGACY) DATA SCIENCE PIPELINE

Or, so called, Data Product



STATIC RESULTS

LOT OF INFORMATION LOST IN TRANSLATION

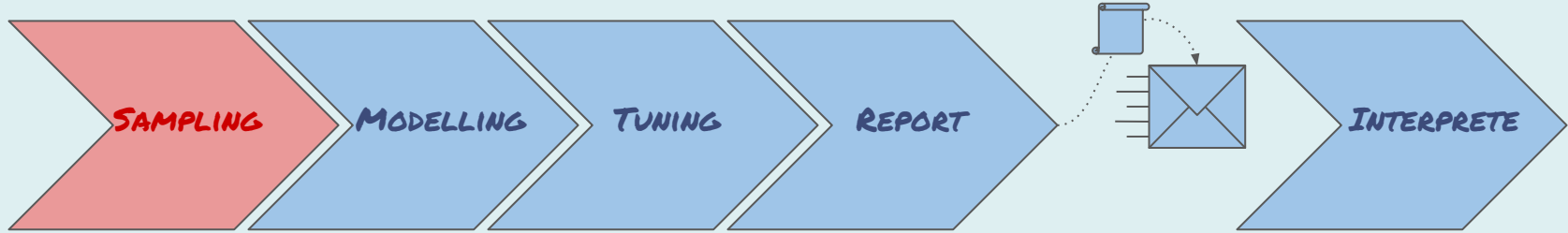
SOUNDS LIKE WATERFALL

ETL LOOK AND FEEL



(LEGACY) DATA SCIENCE PIPELINE

Or, so called, Data Product



MONO MACHINE!

CPU BOUNDS

MEMORY BOUNDS



OUR WORLD TODAY

No, it wasn't better before

FACTS

DATA GETS BIGGER OR, PRECISELY, THE AMOUNT OF AVAILABLE SOURCE EXPLODES

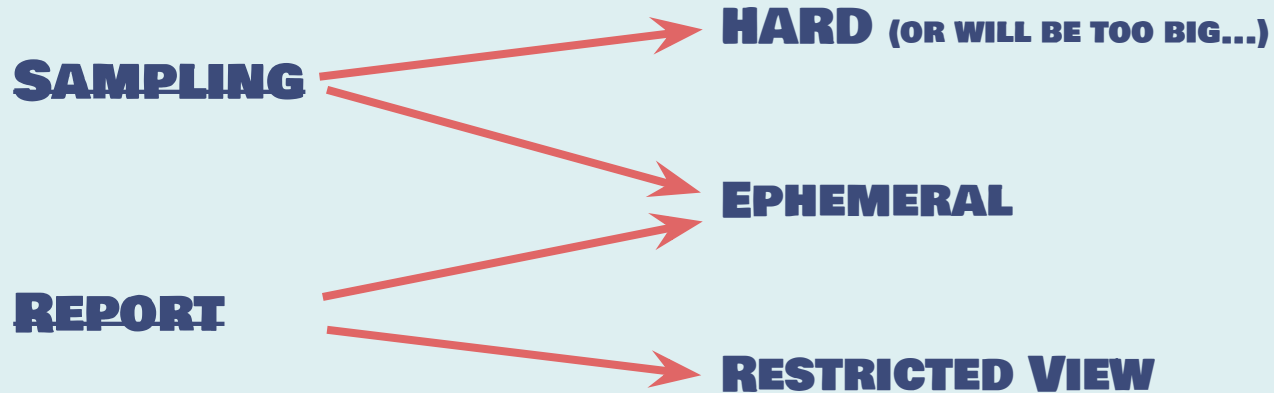
**DATA GETS FASTER (AND FASTER), ONLY EVEN CONSIDER:
WATCHING NETFLIX OVER 4G 🙄**



OUR WORLD TODAY

No, it wasn't better before

CONSEQUENCES



OUR WORLD TODAY

No, it wasn't better before

CONSEQUENCES

INTERPRETATION

⇒ **TOO SLOW TO GET REAL ROI OUT OF THE OVERALL SYSTEM**

HOW TO WORK THAT AROUND?



OUR WORLD TODAY

No, it wasn't better before

NEEDS

ALERTING SYSTEM OVER DESCRIPTIVE CHARTS

MORE ACCURATE RESULTS

MORE OR HARDER MODELS (E.G. DEEP LEARNING)

MORE DATA

CONSTANT DATA FLOW

ONLINE INTERACTIONS UNDER CONTROL (E.G. DIRECT FEEDBACK)



OUR WORLD TODAY

No, it wasn't better before

NEEDS

DISTRIBUTED SYSTEMS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

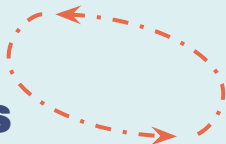
FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES



WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTENT, CONTENT QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS

YO!

**AREN'T WE TALKING ABOUT
"BIG" DATA ?**

FAST DATA ?

**SO COULD REALLY (ALL) RESULTS BEING
NEITHER BIG NOR FAST?**

**ACTUALLY, RESULTS ARE BECOMING
THEMSELVES**

**"BIG" DATA !
FAST DATA !**



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

"CREATE" CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY **HOW DO WE ACCESS DATA SINCE 90'S? REMEMBER SOA?**
→ **SERVICES!**

TUNE PERFORMANCES
NOWADAYS, WE'RE TALKING ABOUT MICRO SERVICES.

WRITE RESULTS TO SINKS
HERE WE ARE, ONE SERVICE FOR ONE RESULT.

ACCESS LAYER

USER ACCESS



DISTRIBUTED DATA SCIENCE

System/Platform/SDK/Pipeline/Product/... whatever you call it

“CREATE” CLUSTER

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

ACCESS LAYER

USER ACCESS

**C’MON, CHARTS/TABLES CANNOT ONLY BE THE
ONLY VIEWS OFFERED TO CUSTOMERS/CLIENTS
RIGHT?**

**WE NEED TO OPEN THE CAPABILITIES TO UI
(DASHBOARD), CONNECTORS (THIRD PARTIES),
OTHER SERVICES (“SOA”) ...**

...

OTHER PIPELINES !!!



WHERE IS MESOS?

(Almost) EVERYWHERE!

"CREATE" CLUSTER

IMPLIES ALLOCATION

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

CREATE DISTRIBUTED DATA PIPELINE/MODEL

TUNE ACCURACY

IMPLIES SCALABILITY

TUNE PERFORMANCES

WRITE RESULTS TO SINKS

IMPLIES DEPLOYMENT

ACCESS LAYER

IMPLIES DEPLOYMENT

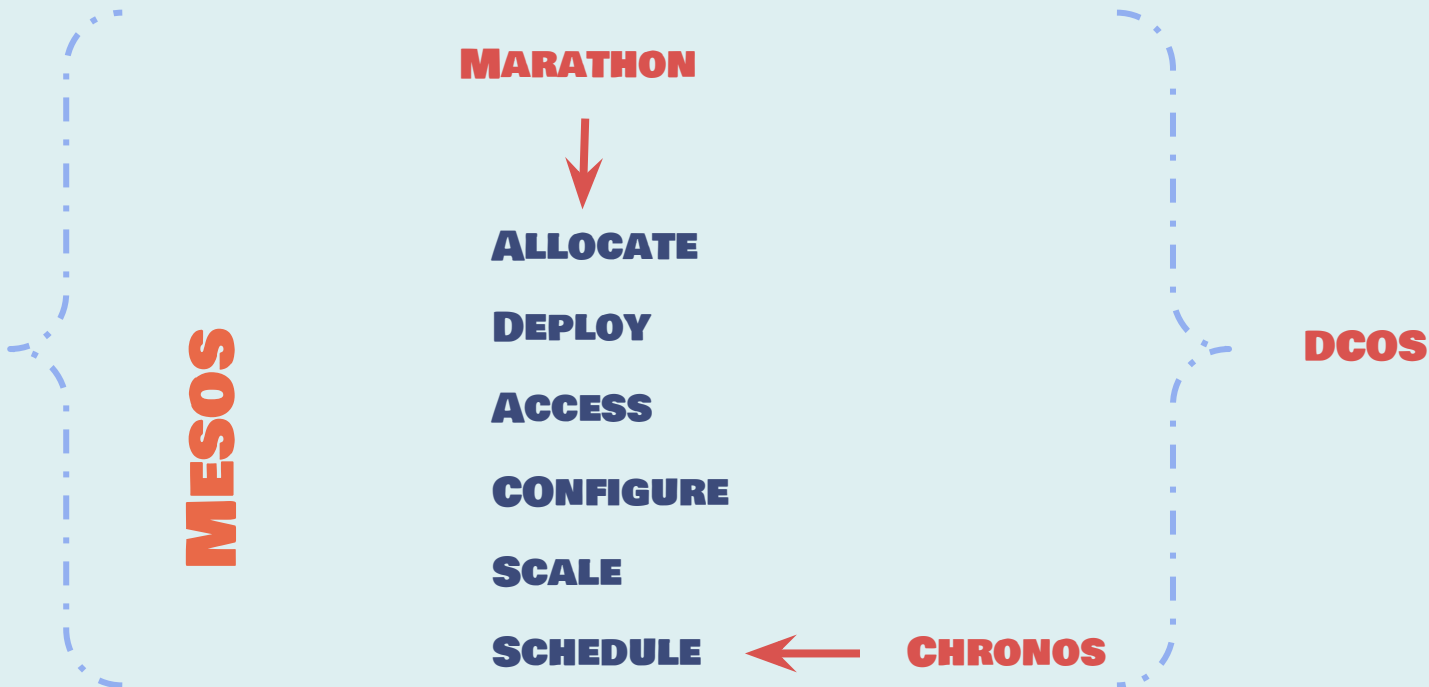
USER ACCESS

IMPLIES SCALABILITY



WHY MESOS?

Because it can... (and even more)



WHAT ABOUT PRODUCTIVITY?

Streamlining development lifecycle most welcome

"CREATE" CLUSTER

OPS

FIND AVAILABLE SOURCES (CONTEXT, CONTENT, QUALITY, SEMANTIC, ...)

DATA

CONNECT TO SOURCES (STRUCTURE, SCHEMA/TYPES, ...)

OPS

DATA

CREATE DISTRIBUTED DATA PIPELINE/MODEL

SCI

TUNE ACCURACY

SCI

TUNE PERFORMANCES

SCI

OPS

WRITE RESULTS TO SINKS

OPS

DATA

ACCESS LAYER

WEB

OPS

DATA

USER ACCESS

WEB

OPS

DATA

SCI



WHAT ABOUT PRODUCTIVITY?

Streamlining development lifecycle most welcome

- **LONGER PRODUCTION LINE**
- **MORE CONSTRAINTS (RESOURCES SHARING, TIME, ...)**
- **MORE PEOPLE**
- **MORE SKILLS**

OVERLOOKING THESE POINTS AND YOU'LL BE

KICKED

SOON OR SOONER

SO, HOW TO HAVE:

- **RESULTS COMING FAST ENOUGH WHILST KEEPING ACCURACY LEVEL HIGH?**
- **RESPONSIVITY TO EXTERNAL/UNPREDICTABLE EVENTS?**



WHAT ABOUT PRODUCTIVITY?

Streamlining development lifecycle most welcome

AT DATA FELLAS, WE THINK THAT WE NEED INTERACTIVITY AND REACTIVITY TO TIGHTEN THE FRONTIERS (WITHIN TEAM AND IN TIME).

HENCE, DATA FELLAS

- **EXTENDS THE SPARK NOTEBOOK (INTERACTIVITY)**
- **IN THE SHAR3 PRODUCT (INTEGRATED REACTIVITY)**



THAT'S ALL FOLKS

Thanks for listening/staying

POKE US ON

@DATAFELLAS

@SHAR3_FELLAS

@SPARKNOTEBOOK

@XTORDOIR & @NOOOTSAB

ANALYSIS
SPARK NOTEBOOK

NOW @TYPESAFE: [HTTP://T.CO/01BT6DQTGH](http://t.co/01BT6DQTGH)

FOLLOW UP SOON ON [HTTP://NOETL.ORG](http://noetl.org)

(HI5 TO @CHIEFSCIENTIST FOR THAT)

