



BIG DIVE

DATA SCIENCE & ANALYTICS

A SOFT INTRO TO APPLIED DATA SCIENCE

CHRISTIAN RACCA

Organized by



Designed for



In collaboration with



ISI Foundation



BEFORE STARTING



**Please,
wear your
badge**



SSID **bigdive
PWD **wifi4divers****



**09:00
17:30**

BEFORE STARTING



**Coffe
machine &
kitchen
available
in house**



**Please
notify
delay or
absence**



**We'll share
all the
teaching
materials
on GIT**

BEFORE STARTING



**Remote
Desktop
to VM**



**Slack
Channel**



**Git
Repo**

PROGRAM

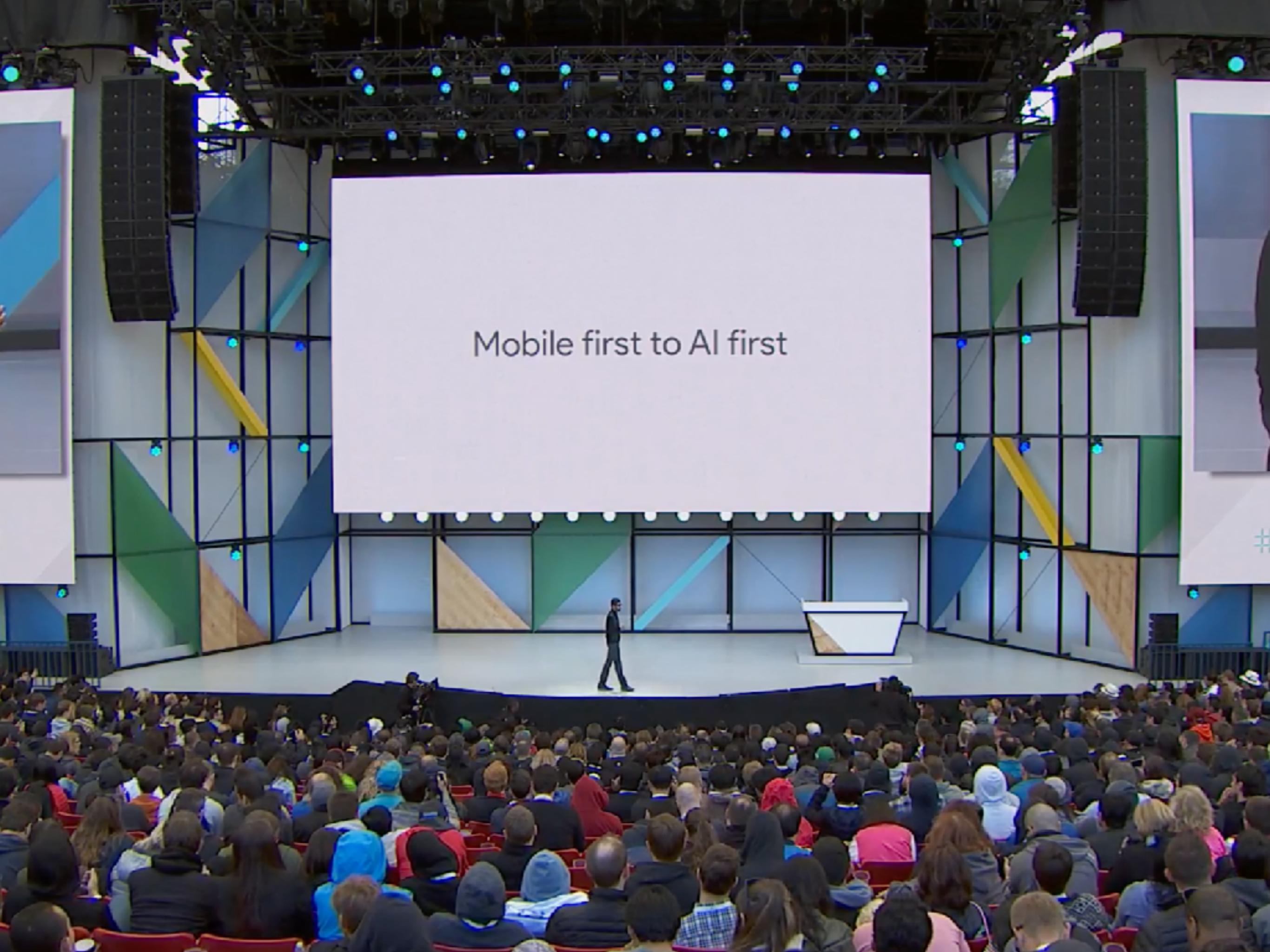
	DAY 1 (09 ott)	DAY 2 (10 ott)
09:30 - 11:00	<i>Soft intro to applied data science</i>	<i>Python basics</i>
11:00 - 11:15	<i>Break</i>	<i>Break</i>
11:15 - 12:30	<i>ISI lecture</i>	<i>Python basics</i>
12:30 - 13:30	<i>Aizoon lecture</i>	
13:30 - 14:30	<i>Lunch break</i>	<i>Lunch break</i>
14:30 - 15:30	<i>Data Viz intro</i>	<i>Python basics</i>
15:30 - 16:30	<i>Data Dev intro</i>	<i>HTML & CSS essentials</i>

(Disclaimer)
**I AM NOT A DATA
SCIENTIST**

**(BIG) DATA
OPPORTUNITIES
BEYOND THE
BUZZWORD**

THIS IS BOOKING.COM

Tipologia camera	Per	Il prezzo di oggi	Opzioni disponibili	Scegli le camere	Conferma la tua prenotazione
<p>Camera Doppia con Letti Singoli con Bagno Privato Esterno <small>Il nostro consiglio per te Super richiesta!</small></p> <p>Scelta 1 camera rimasta sul nostro sito!</p> <p>2 letti singoli </p> <ul style="list-style-type: none"> vista città TV a schermo piatto aria condizionata insonorizzazione bagno privato WIFI gratis • vista • TV • scrivania • ferro e asse da stiro • riscaldamento • disponibilità di camere comunicanti • ingresso indipendente • parquet o pavimento in legno • stand appendiabiti • doccia • asciugacapelli • accappatoio • WC • bagno privato • vasca o doccia • bidet • carta igienica • bollitore tè / macchina caffè • frigorifero • bollitore elettrico • macchina da caffè • zona pranzo all'aperto • piani superiori accessibili solo tramite scale • libri, DVD o musica per bambini <p>Include: 10 % di IVA, Colazione. Non include: 2.00 € di Tassa di soggiorno per persona a notte .</p>			<ul style="list-style-type: none"> Colazione eccezionale inclusa nel prezzo Prezzo basso - Nessun rimborso NON SERVE ALCUN PAGAMENTO ANTICIPATO - Paga in struttura 	<input style="width: 20px; height: 20px; border: 1px solid #ccc; padding: 2px; margin-right: 5px;" type="button" value="0"/>	<p>Selezione</p> <p>Conferma immediata</p> <p>I data scientist di Booking.com hanno fatto i conti: 1 altra persona sta guardando questa pagina.</p> <p>I prezzi potrebbero aumentare. Garantisce oggi la tua prenotazione.</p>
			<ul style="list-style-type: none"> Colazione eccezionale inclusa nel prezzo Prezzo basso - Nessun rimborso NON SERVE ALCUN PAGAMENTO ANTICIPATO - Paga in struttura 	<input style="width: 20px; height: 20px; border: 1px solid #ccc; padding: 2px; margin-right: 5px;" type="button" value="1"/>	



Mobile first to AI first

YES, BIG DATA IS A “BUZZWORD”

THE TERM AMBIGUITY

**quantitative
marketing
technology
research**



**qualitative
innovation
process
business**



WHAT'S “NEW” ABOUT DATA

DATA AVAILABILITY
/ Exponential growth
/ Machine VS human
/ Structured VS un-structured

INFRASTRUCTURE AS A COMMODITY
/ Cloud
/ HPC & HPN
/ Frameworks

SKILLS TO EXTRACT INFORMATION ARE NOW MORE ACCESSIBLE

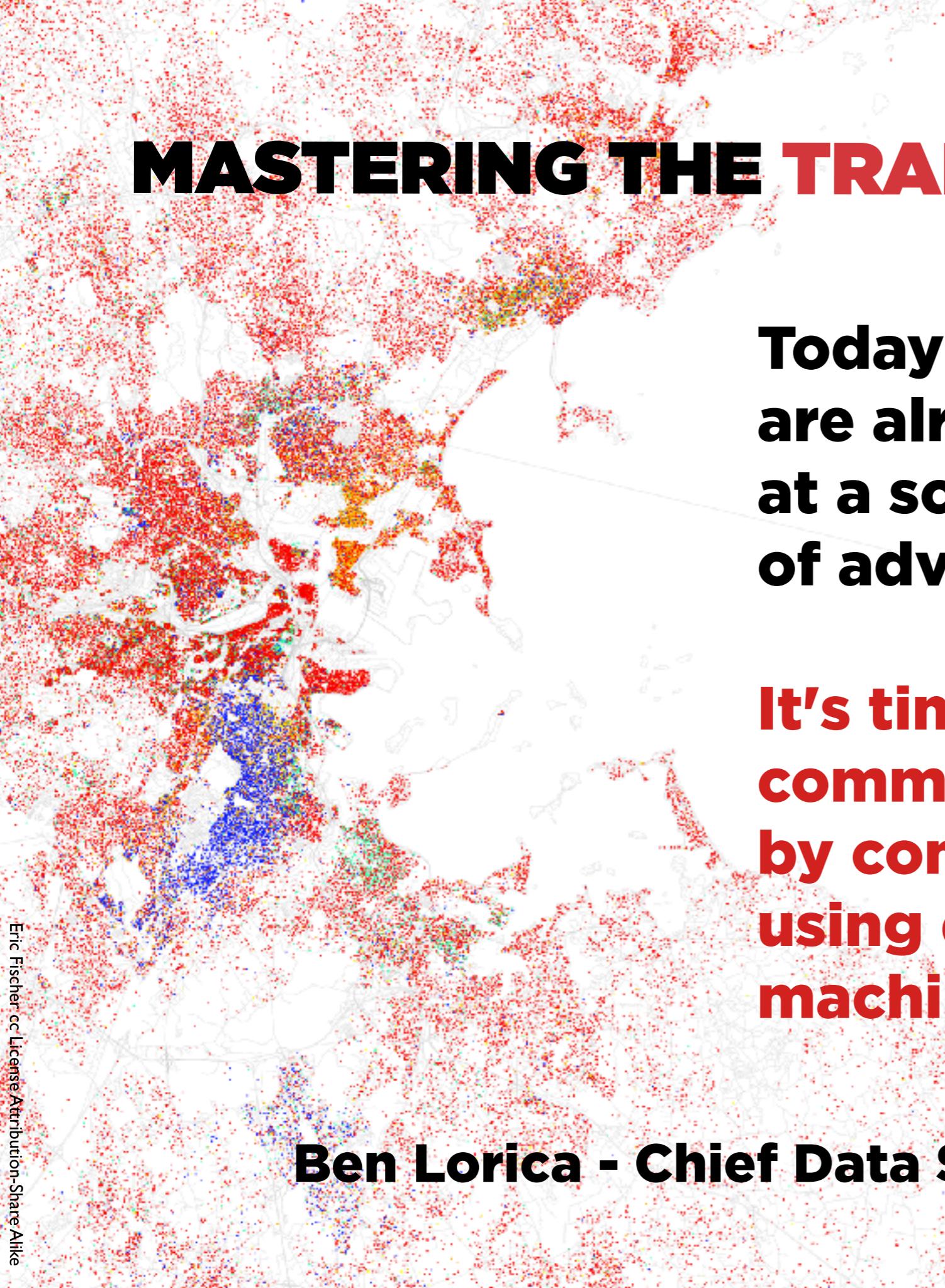
CULTURE & APPROACH
/ Complexity science
/ Network thinking
/ Open Innovation



BIG DATA + ML = The NEW STACK

Big Data technologies are used to handle core data engineering challenges, and machine learning is used to extract value from the data.

MASTERING THE TRANSITION



**Today many organizations
are already collecting data
at a scale that justifies use
of advanced analytic tools...**

**It's time to focus at the
common challenges faced
by companies interested in
using data science and
machine learning.**

Ben Lorica - Chief Data Scientist at O'Reilly Media

COMMON OPEN CHALLENGES

/ THE DATA

/ THE SKILLS

/ FROM PROTOTYPE TO...

**/ THE RESULTS INTERPRETATION AND
THE EXPLAINABILITY ISSUE**

/ “GREY ZONES” IN DATA EXPLOITATION

/ THE PURSUIT OF INNOVATION



DATA METADATA FEATURES

CHALLENGE #1

DATA REMAINS THE STARTING POINT



..... **Volume**

The effective amount of usable data. No a-priori objective parameters. On field validation is required.

DATA REMAINS THE STARTING POINT

Volume

The effective amount of usable data. No a-priori objective parameters. On field validation is required.

Metadata

“Data” that provides information about other data. {Descriptive, Structural, Administrative}

DATA REMAINS THE STARTING POINT

Volume

The effective amount of usable data. No a-priori objective parameters. On field validation is required.

Metadata

“Data” that provides information about other data. {Descriptive, Structural, Administrative}

Features Selection

Refers to the process of extracting useful information (or features) from existing data.

ABOUT FEATURES...

FROM SOURCE DATA

Features “reduction”

Noisy or redundant data makes it more difficult to discover meaningful patterns.

High-dimensional dataset requires more complex models/algorithms and more computational power.

Data augmentation

Enriching existing data with open data or through third-party data providers.



TO RELEVANT DATA

DATA REMAINS THE STARTING POINT

Volume

The effective amount of usable data. No a-priori objective parameters. On field validation is required.

Metadata

“Data” that provides information about other data. {Descriptive, Structural, Administrative}

Features Selection

Refers to the process of extracting useful information (or features) from existing data.

DATA REMAINS THE STARTING POINT

Volume

The effective amount of usable data. No a-priori objective parameters. On field validation is required.

Metadata

“Data” that provides information about other data. {Descriptive, Structural, Administrative}

Features Selection

Refers to the process of extracting useful information (or features) from existing data.

Data Quality

Traceability, expiration, completeness, currentness compliance, understandability, accuracy.

DATA QUALITY

Characteristic	Metric
Traceability	Track of creation
	Track of updates
Currentness	Percentage of current rows
	Delay in publication
Expiration	Delay after expiration
Completeness	Percentage of complete cells
	Percentage of complete rows
Compliance	Percentage of standardized columns
	eGMS Compliance
	Five star Open Data
Understandability	Percentage of columns with metadata
	Percentage of columns in comprehensible format
Accuracy	Percentage of accurate cells
	Accuracy in aggregation

DATA QUALITY

Name	Gender	Street	House #	Zip code	City	State	D.O.B
John Doe	Male	60th street	45		New York	New York	08/12/64
Jane Doe	Female	Jonathan	36	10023	Poughkeepsy	NY	21-dec-1954

Name	Gender	Street	House #	Zip code	City	State	D.O.B
John Doe	Male	E 60 th St	45	10022	New York	NY	08/12/64
Jane Doe	Female	Jonathan St	36	10023	Poughkeepsie	NY	12/21/54

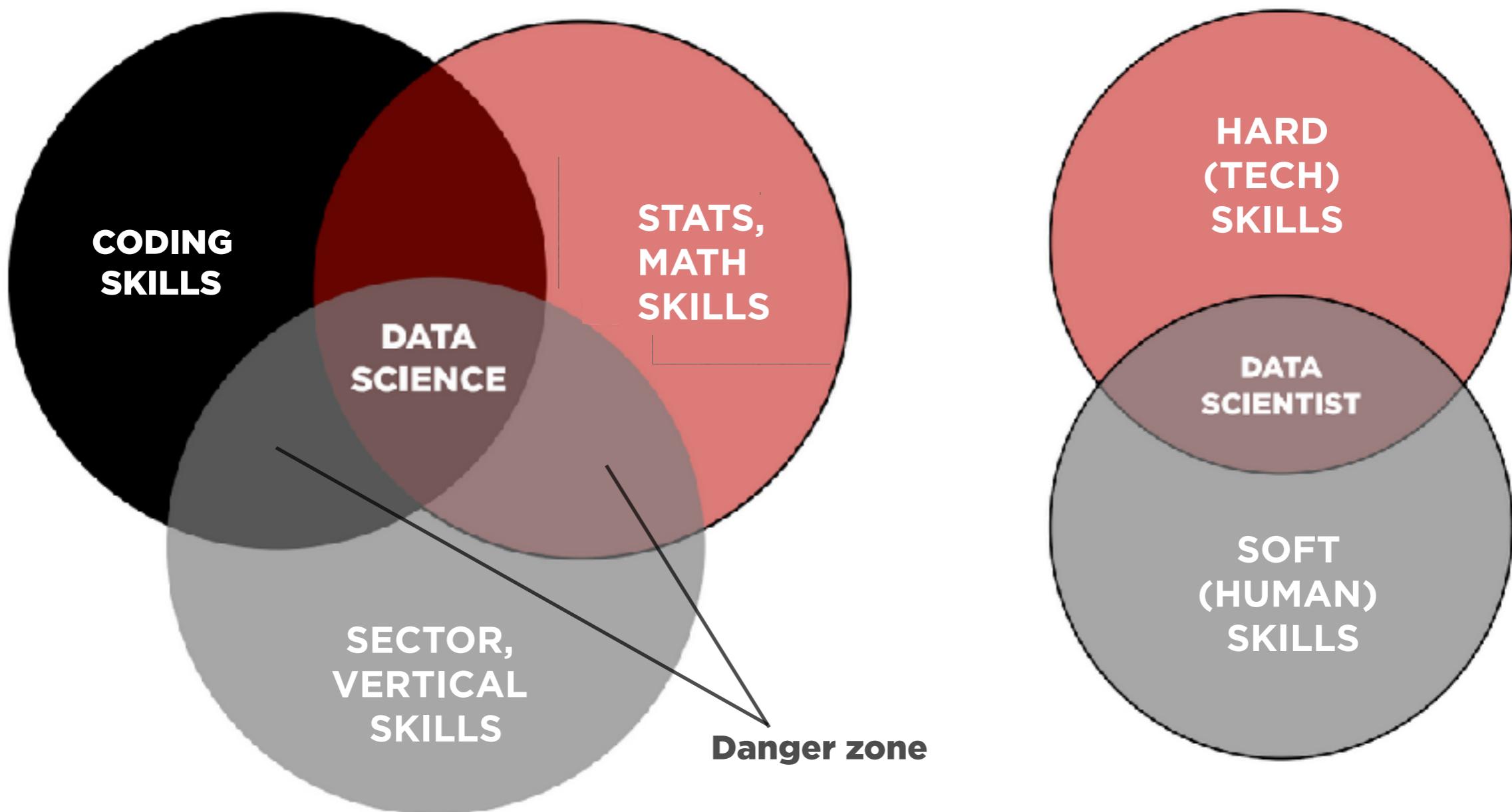
- Completeness
- Accuracy
- Conformity
- Consistency
- Uniqueness

THE DATA TEAM

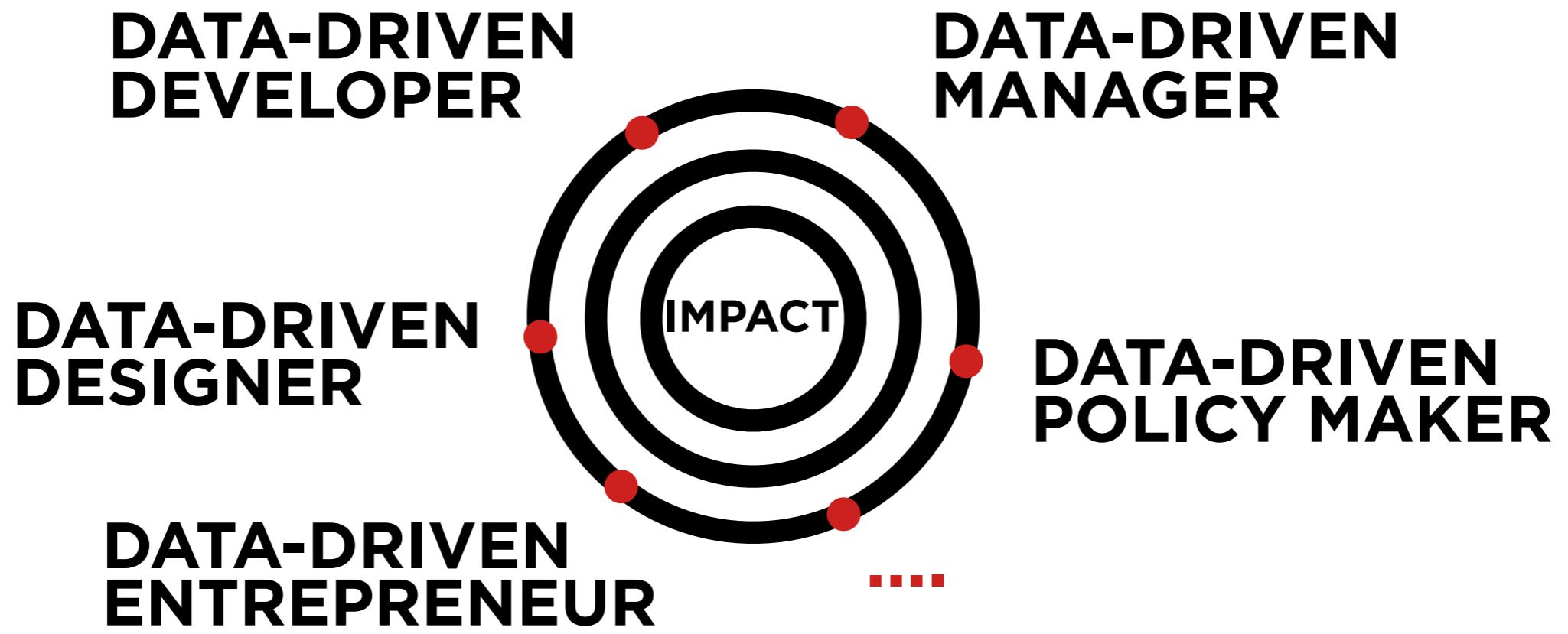
CHALLENGE #2

LOOKING FOR UNICORNS

Re-arrangement of the THE DATA SCIENCE VENN DIAGRAM
by Drew Conway



THE DATA TEAM



SINGLE VS TEAM

FROM PROTOTYPE TO PRODUCTION

CHALLENGE #3



COMPARISON

— Jobs in Data Science —



Data Scientist



vs Data Engineer



vs Statistician

These people use their analytical and technical capabilities to extract meaning insights from data.

These people ensure uninterrupted flow of data between servers and applications. They are responsible for data architecture.

These people understand statistics theoretically and apply them to real life problems.

Responsibilities

Develop and plan required analytic projects in response to business needs.

Contribute to data mining architectures, modeling standards, reporting, and data analysis methodologies.

Collaborate with stakeholders to integrate data mining results with existing systems.

Monitor data mining system performance and implement efficiency improvements.

Design, construct, install, test and maintain highly scalable data management systems

Improve data foundational procedures, guidelines and standards

Integrate new data management technologies and software engineering tools into existing structures

Create custom software components (e.g. specialized UDFs) and analytics applications

Apply statistical theories and methods to solve practical problems of various industries

Determine methods for finding or collecting data

Design surveys or experiments or opinion polls to collect data

Analyze, interpret & undertake data analysis

Report conclusions from their analyses

Skills

Programming, Mathematics, Business Understanding, Statistics, Data Visualization, Machine Learning, Attention to detail

Database design, Production coding, Data collection, data warehousing, Data transformation, Work diligently with data

Technical and Analytics Skills, Mathematics, Operational Research, Writing skills, Ability to Analyze, Model and interpret data, Flair of explaining difficult concepts in simple manner

A BABEL OF (CODING) LANGUAGES

PRODUCTION

JAVA, C, C++, ...

**DATA-DRIVEN
PROTOTYPE**

PYTHON, R, D3.JS

REFACTORING

**DATA
ENGINEERING**
[SCALA, ...]

RESULTS INTERPRETATION

&

EXPLAINABILITY

CHALLENGE #4

THE EXPLAINABILITY ISSUE

LOW EXPLAINABILITY

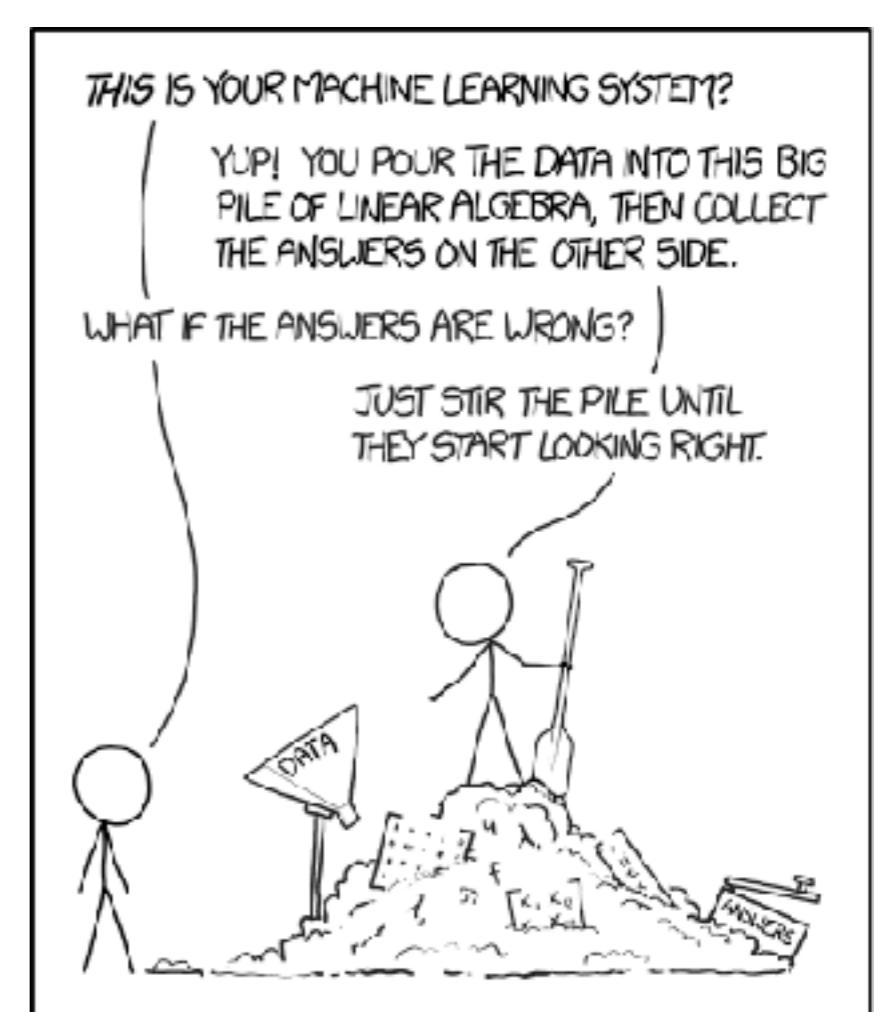


Deep learning

Machine learning

Inferential statistics

HIGH EXPLAINABILITY



MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:
/ Errors in the training set

GARBAGE-IN GARBAGE-OUT

**MACHINE learns ONLY through the training data
(no prejudices, no additional elaboration)**

Training set:

$$2+2 = 5$$

$$2+2 = 5$$

$$2+2 = 5$$

$$2+2 = 5$$

$$2+2 = 5$$

$$2+2 = 5$$

$$2+2 = 4$$

>>> MACHINE SAYS THAT $2+2 = 5$

MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:
/ Errors in the training set

MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:

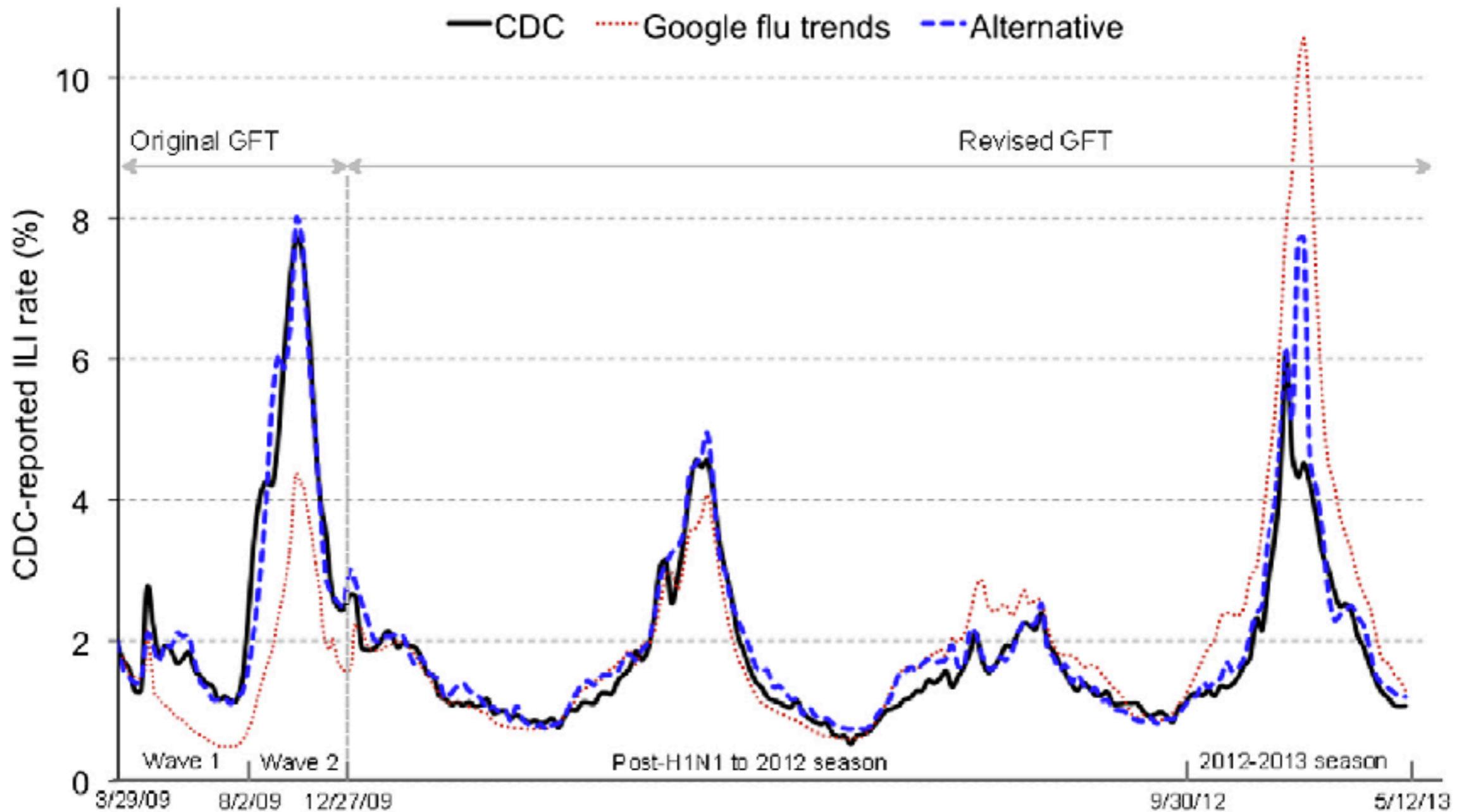
- / Errors in the training set**
- / Errors in generalization**

MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:

- / Errors in the training set**
- / Errors in generalization**
- / Missing update in the model**

ABOUT CORRELATION: GOOGLE FLU



MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:

- / Errors in the training set**
- / Errors in generalization**
- / Missing update in the model**

MACHINE LEARNING IS NOT PERFECT

BAD LEARNING CAUSES:

- / Errors in the training set**
- / Errors in generalization**
- / Missing update in the model**
- / Human errors**

MACHINE LEARNING IS NOT PERFECT

**One of the biggest error in creating a training dataset
is assuming people think and behave in a common-
similar way.**

Timer: 00:00:00 of 7 days Want to work on this HIT? Want to see other HITs?

Total Earned: \$0.00 Total HITs Submitted: 0

Accept HIT Skip HIT

Describe a picture in your own words

Requester: Scott Lobdell

Qualifications Required: None

Reward: \$0.05 per HIT HITs Available: 29 Duration: 7 days

Using your own words, describe this man however you want.

Submit



A screenshot of a Mechanical Turk task interface. The top bar shows a timer at 00:00:00 of 7 days, and buttons for accepting or skipping the HIT. It also displays total earnings of \$0.00 and 0 submitted HITs. Below the bar, requester information (Scott Lobdell) and qualification requirements (None) are listed, along with reward details (\$0.05 per HIT, 29 HITs available, 7-day duration). The main task area contains instructions to "Describe a picture in your own words" and a text input field with a "Submit" button. Below the input field is a photograph of a shirtless man standing on a beach, leaning against a palm tree trunk, holding a surfboard.

MACHINE LEARNING IS NOT PERFECT

**One of the biggest error in creating a training dataset
is assuming people think and behave in a common-
similar way.**

Timer: 00:00:00 of 7 days Want to work on this HIT? Want to see other HITs?

Total Earned: \$0.00 Total HITs Submitted: 0

Accept HIT Skip HIT

Describe a picture in your own words

Requester: Scott Lobdell Qualifications Required: None

Reward: \$0.05 per HIT HITs Available: 29 Duration: 7 days

Using your own words, describe this man however you want.



MACHINE LEARNING IS NOT PERFECT

**One of the biggest error in creating a training dataset
is assuming people think and behave in a common-
similar way.**

Timer: 00:00:00 of 7 days Want to work on this HIT? Want to see other HITs?

Total Earned: \$0.00 Total HITs Submitted: 0

Accept HIT Skip HIT

Describe a picture in your own words

Requester: Scott Lobdell Qualifications Required: None

Reward: \$0.05 per HIT HITs Available: 29 Duration: 7 days

Using your own words, describe this man however you want.

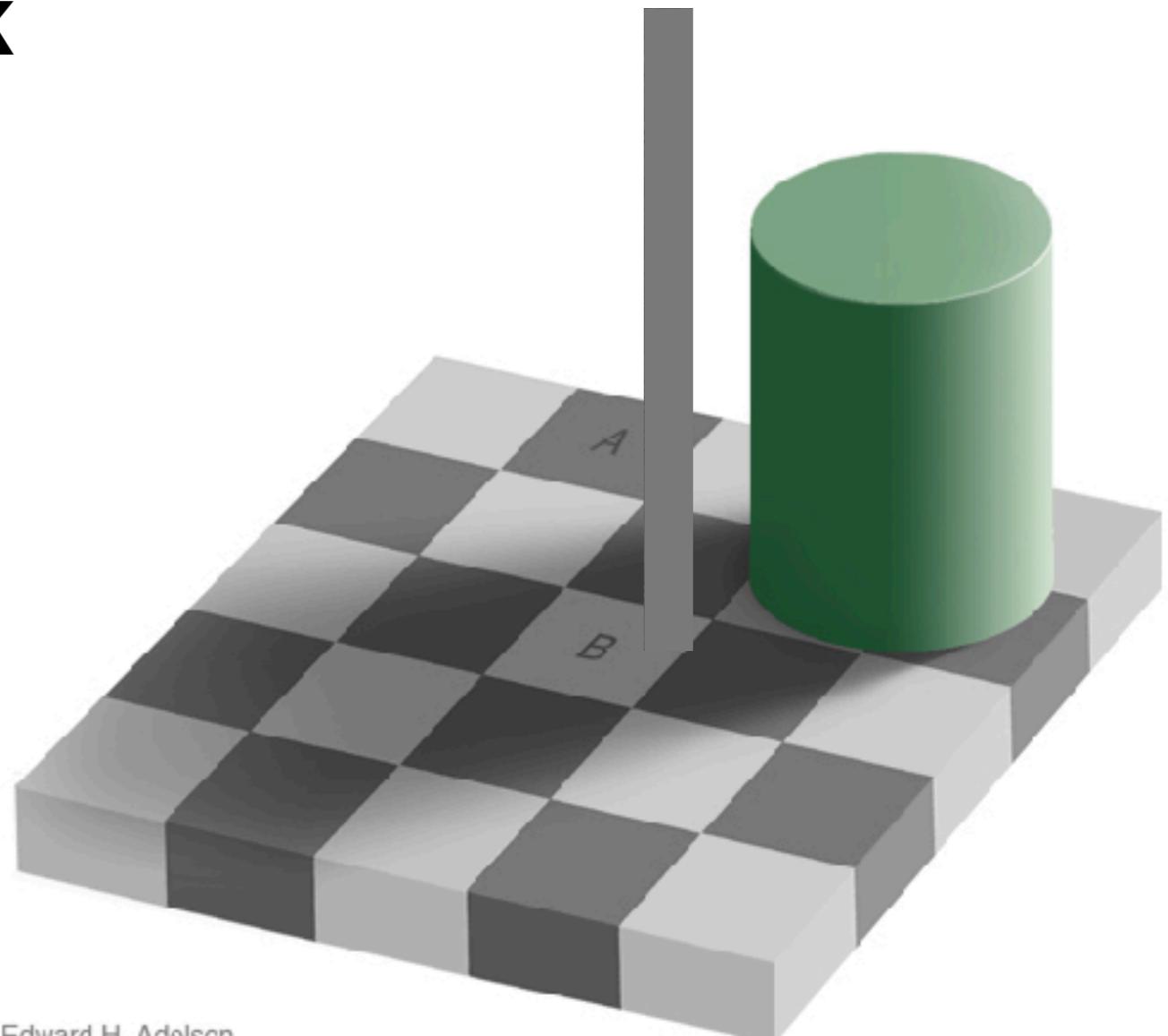
Submit



MACHINE LEARNING IS NOT PERFECT

OPTICAL ILLUSIONS TRICK

HUMAN MIND: do A & B
**squares have the same
color ?**



Edward H. Adelson

**GREY ZONES
IN
DATA EXPLOITATION**

CHALLENGE #5

THE DARK SIDE OF BIG DATA

/ DATA DEMOCRACY

**/ CORRELATION DOES NOT MEAN
CAUSATION**

/ BUBBLE FILTERS

/ BIAS

/ PRIVACY & LEGAL ISSUES

/ HUMAN RIGHTS

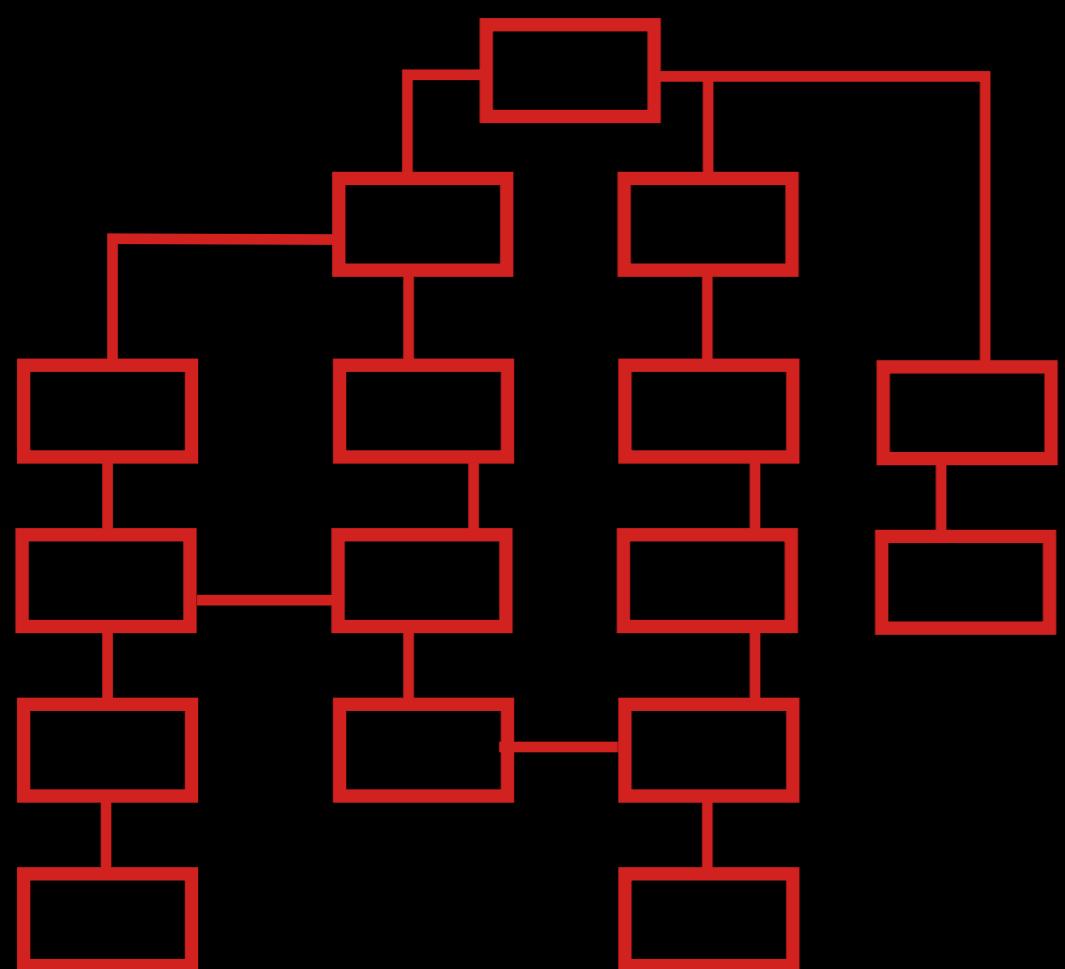
/ ...

GRID NATIVES
VS
COMPLEX NATIVES

CHALLENGE #6

HOW TO BECOME A DATA DRIVEN COMPANY ?

THE GRID VS THE COMPLEXITY



GRID “NATIVES”

**HAVE LOW
INTERCONNECTION**

**DON’T LIKE CHANGES
& “CHAOS”**

**GET OLD AND NEED
CONTINUOUS FIXING**

**ARE HIGHLY EXPOSED TO
BIG “CRISIS” (black swans)**

**YOU NEED A MAP TO
UNDERSTAND THEM**

COMPLEX “NATIVES”

***Have high
interconnection***

***Change a lot & like
chaos (edge)***

***Get old but also evolve/
adapt (auto-repair)***

***Prefer failing (small
mistakes) fast & soon***

***You need dynamic data
to understand them***

GRID “NATIVES”

SOFTWARE VENDORS

BANKS

HARDWARE VENDORS

PORTALS

...

COMPLEX “NATIVES”

service (experience) vendors

mobile payment platforms

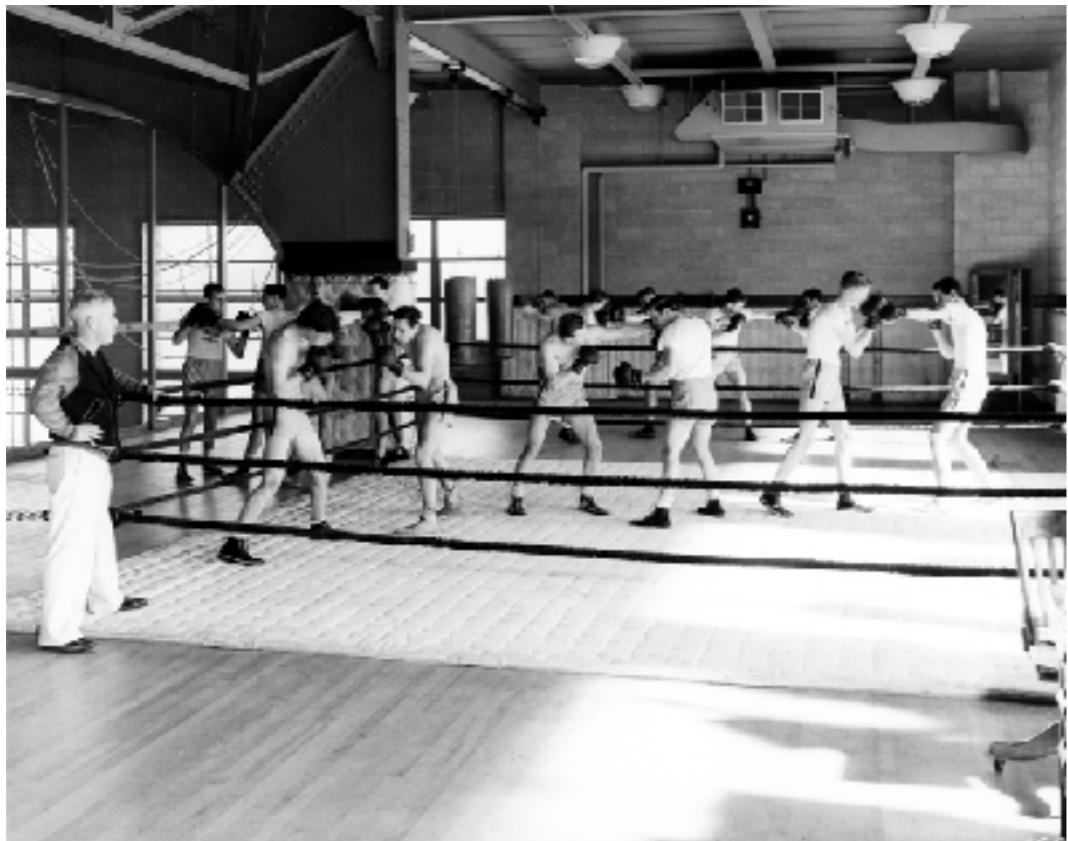
cloud vendors

social networks

...

THE BIG DIVE CREDO

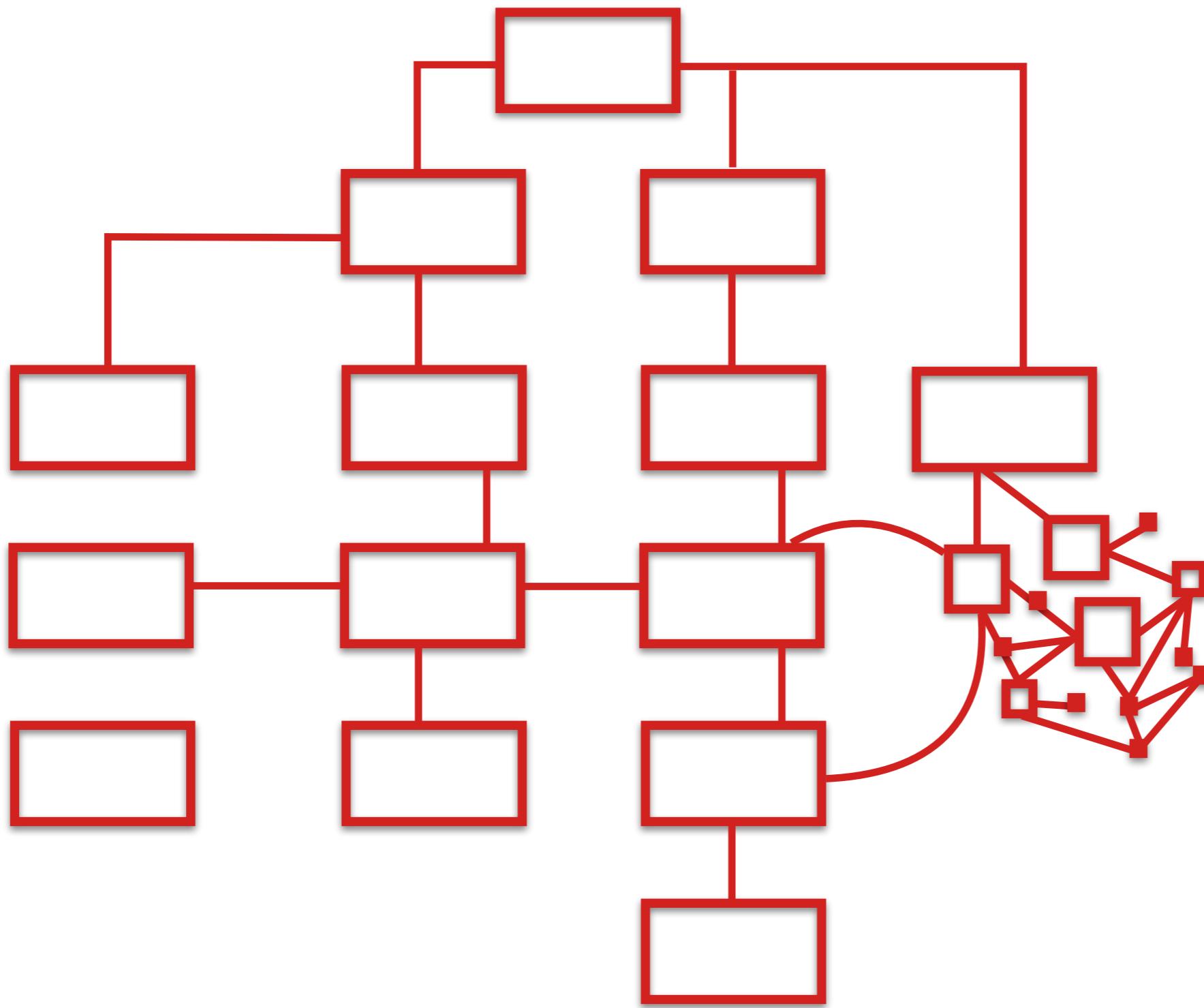
**WE DON'T NEED A
UNIVERSITY**



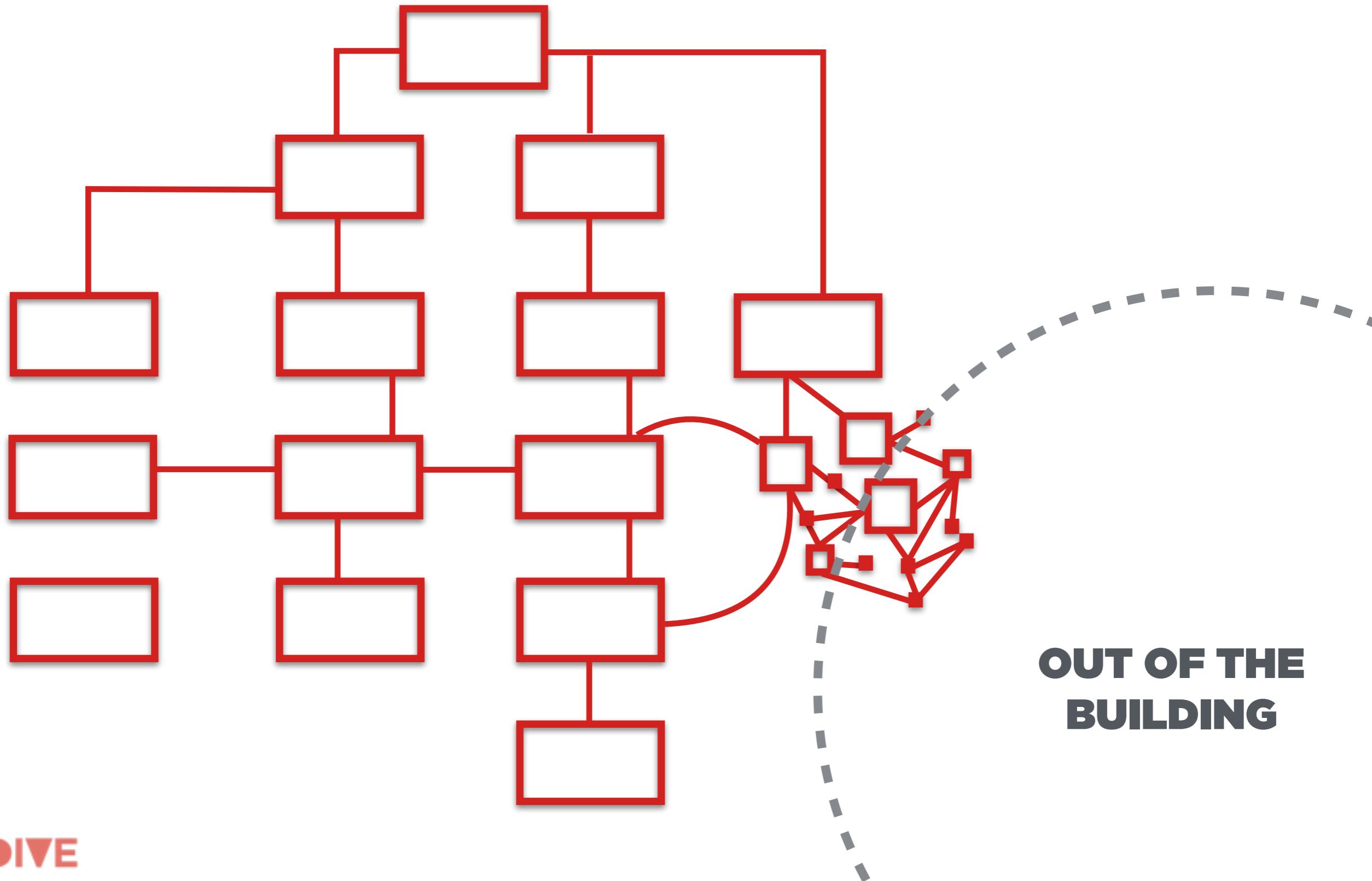
**WE NEED
STREET FIGHTING**



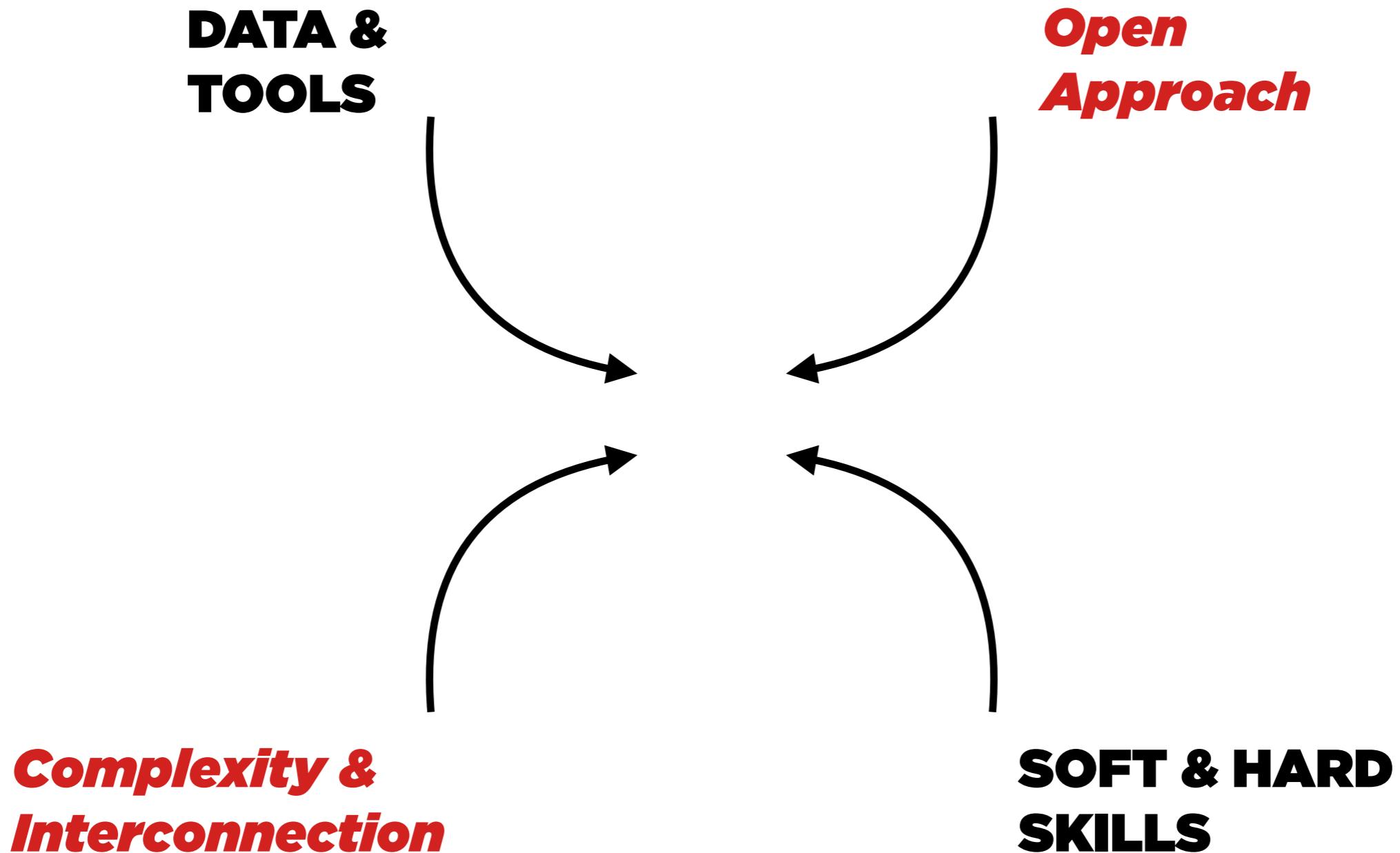
MAKING THE GRID *MORE COMPLEX*



IN VS OUT



MASTERING THE TRANSFORMATION



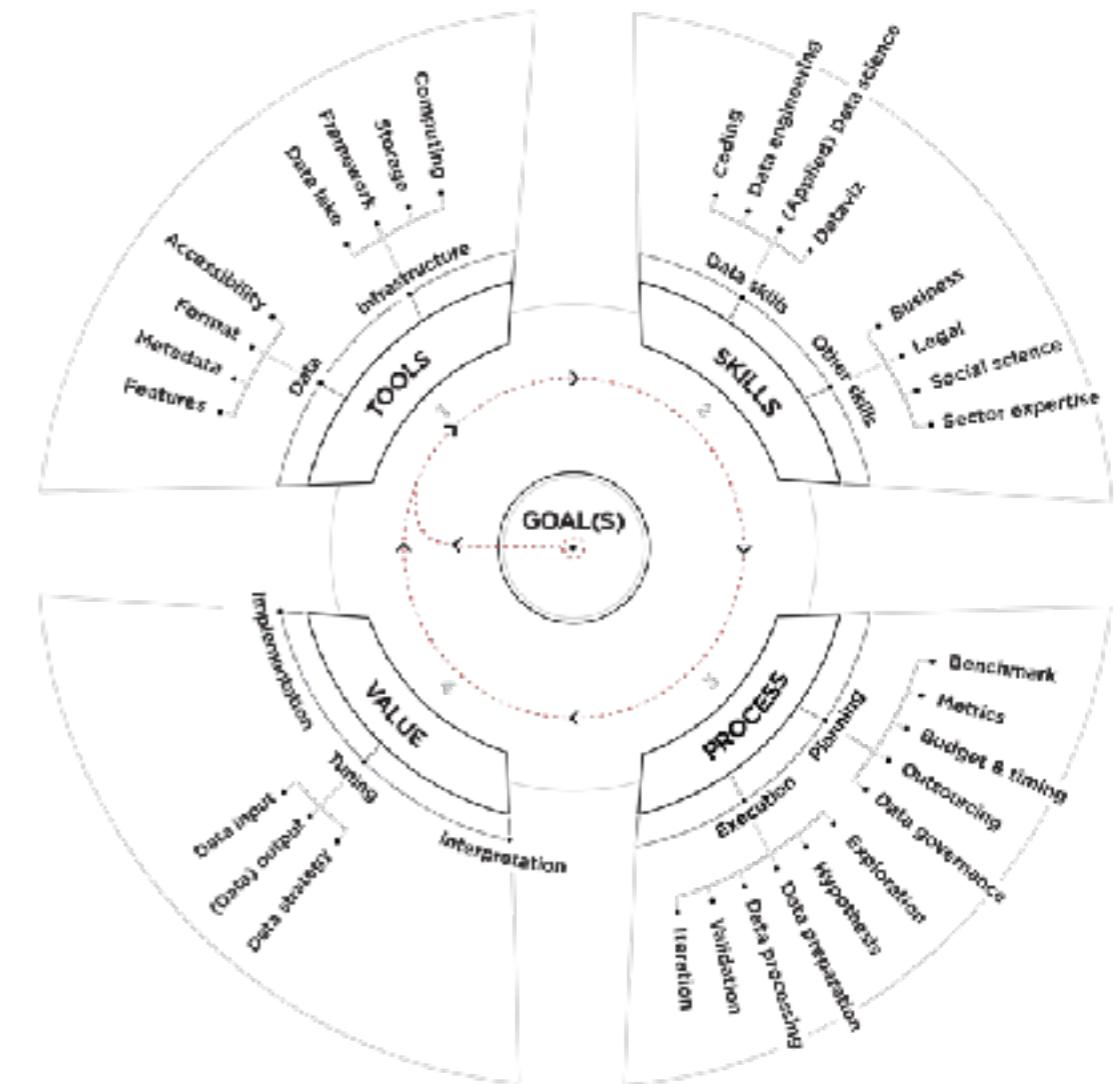
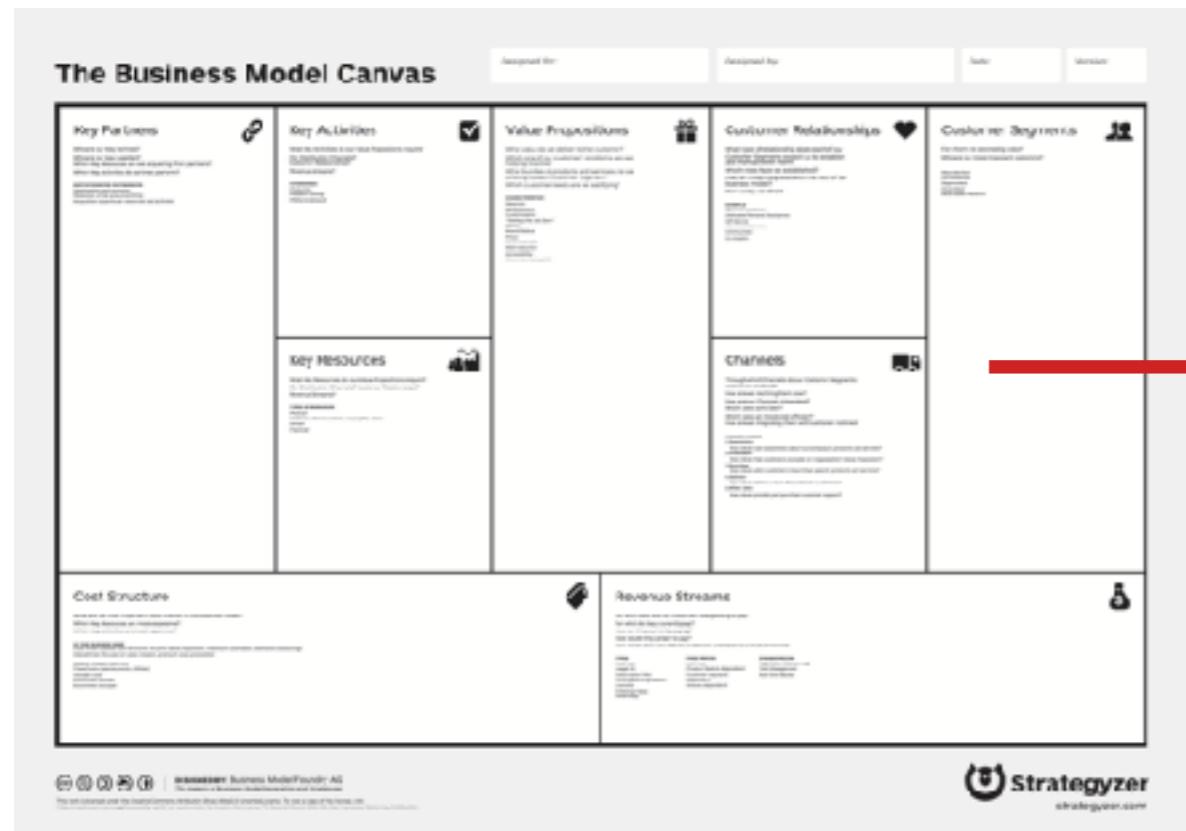
LESSONS LEARNED ABOUT DATA

- 1. DATA is a core asset that it is advisable to develop internally (at the companies).**
- 2. DON'T FOCUS ON TOOLS, FOCUS ON SCIENCE & INNOVATION PROCESSES.**
- 3. GO OPEN (source)! Open Community is the new R&D department.**
- 4. Approach must be SYSTEMIC, LEARN-BY-DOING, DYNAMIC, VALUE-BASED.**
- 5. DECENTRALIZATION is the new mantra.**
- 6. NEVER SEPARATE DATA FROM CONTEXT & SOURCES.**

**LEVERAGING (BIG) DATA
OPPORTUNITIES
REQUIRES
A PROPER METHOD**

THE DATA RING

THE CANVAS APPROACH



The inspiring precursor

The Data Ring

WHY A DATA CANVAS

- / It forces the project owner to state crystal clear the **value proposition** of the project.
- / It is an analytical tool, devoted to **self-diagnosis** and to define and respect an internal strategy.
- / It provides for a complete representation of the process that can be explained to **third parties** too.
- / It is not a “static shot” but it **evolves through time** according to project evolution.

DATA RING CANVAS



GOAL(S) SETTING



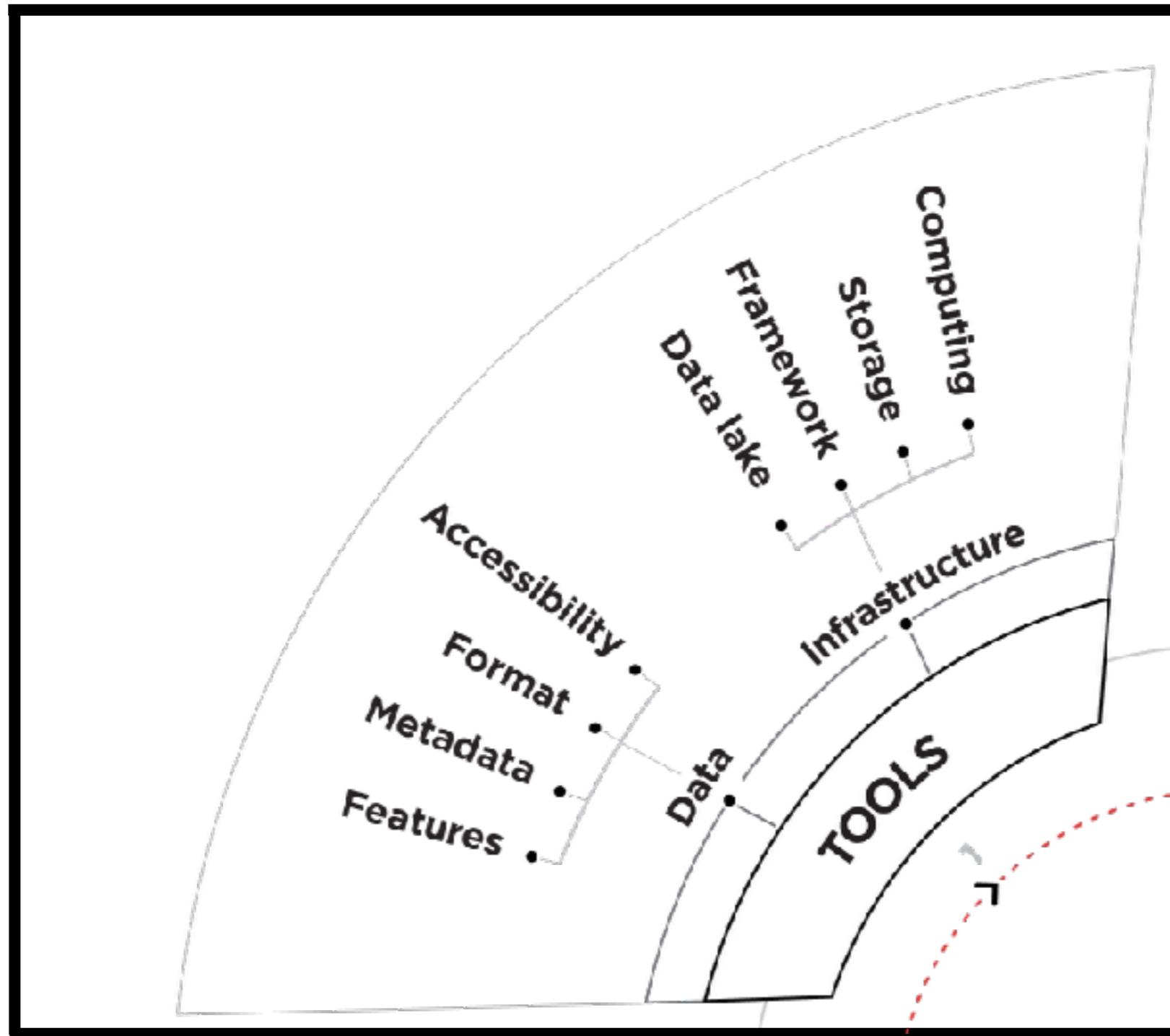
BEST PRACTICES:

- 1) Start from the problem-setting
- 2) Be specific!
- 3) Try to be quantitative (not only qualitative)

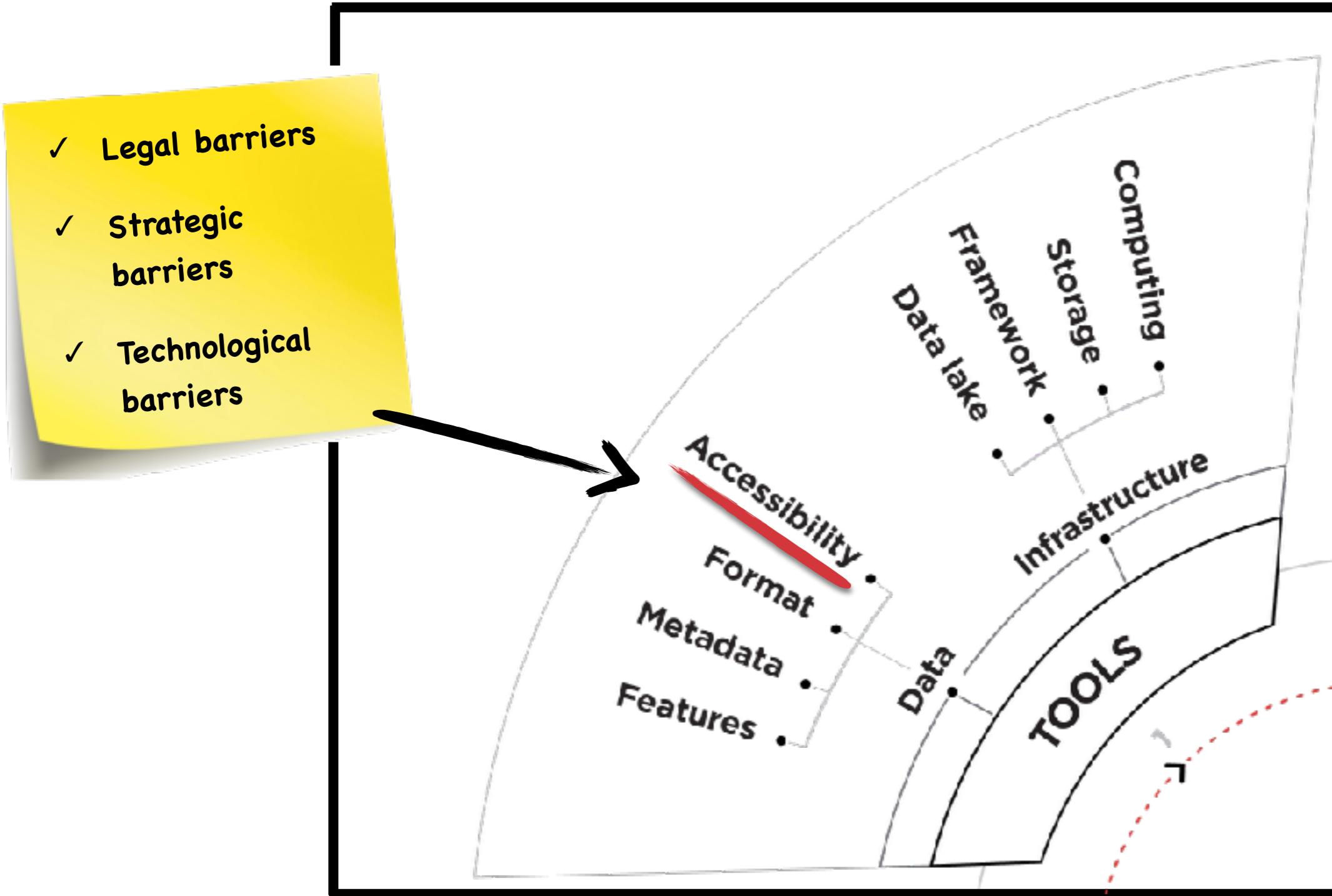
GOAL(S) SETTING



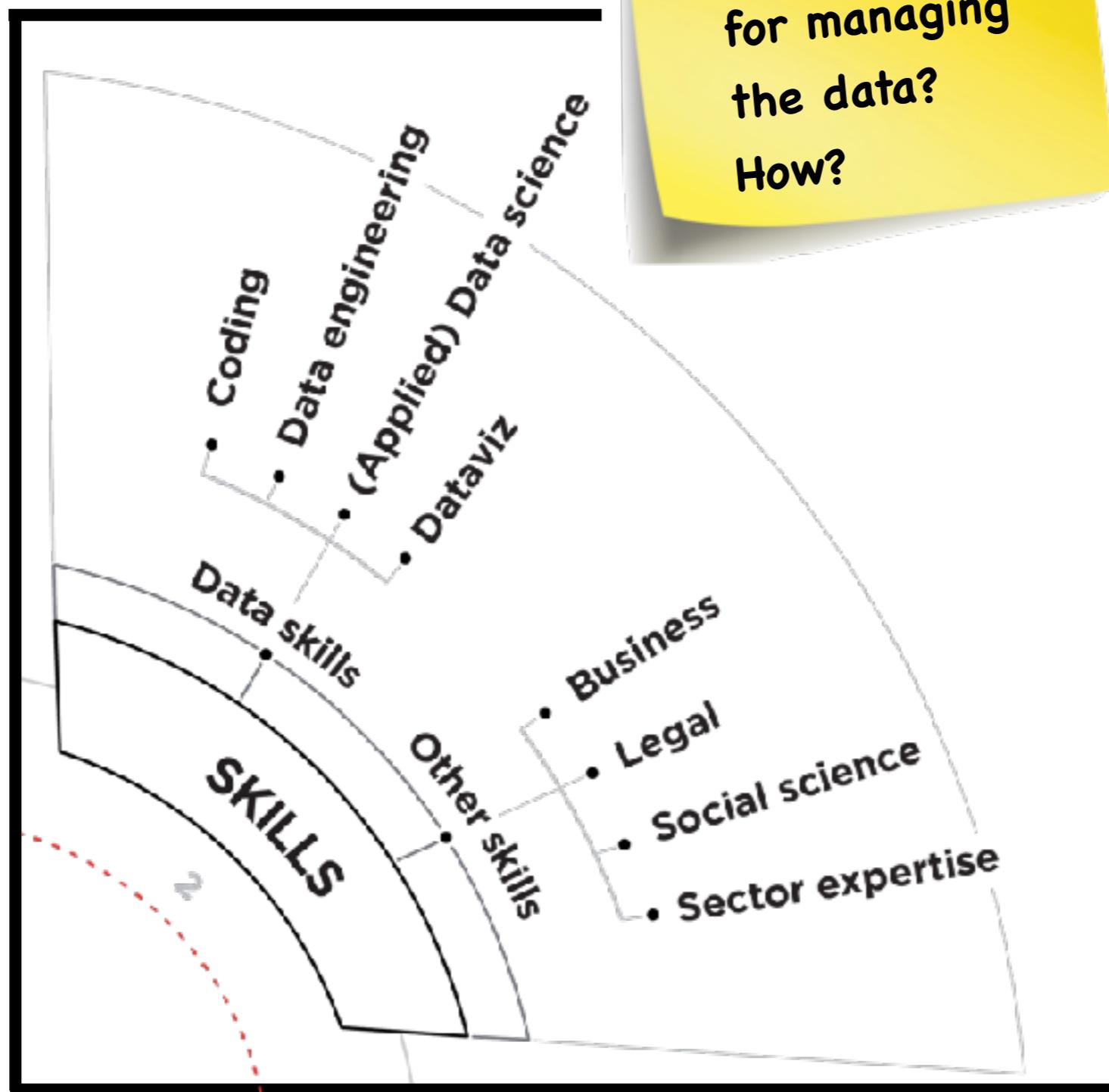
TOOLS



TOOLS



SKILLS



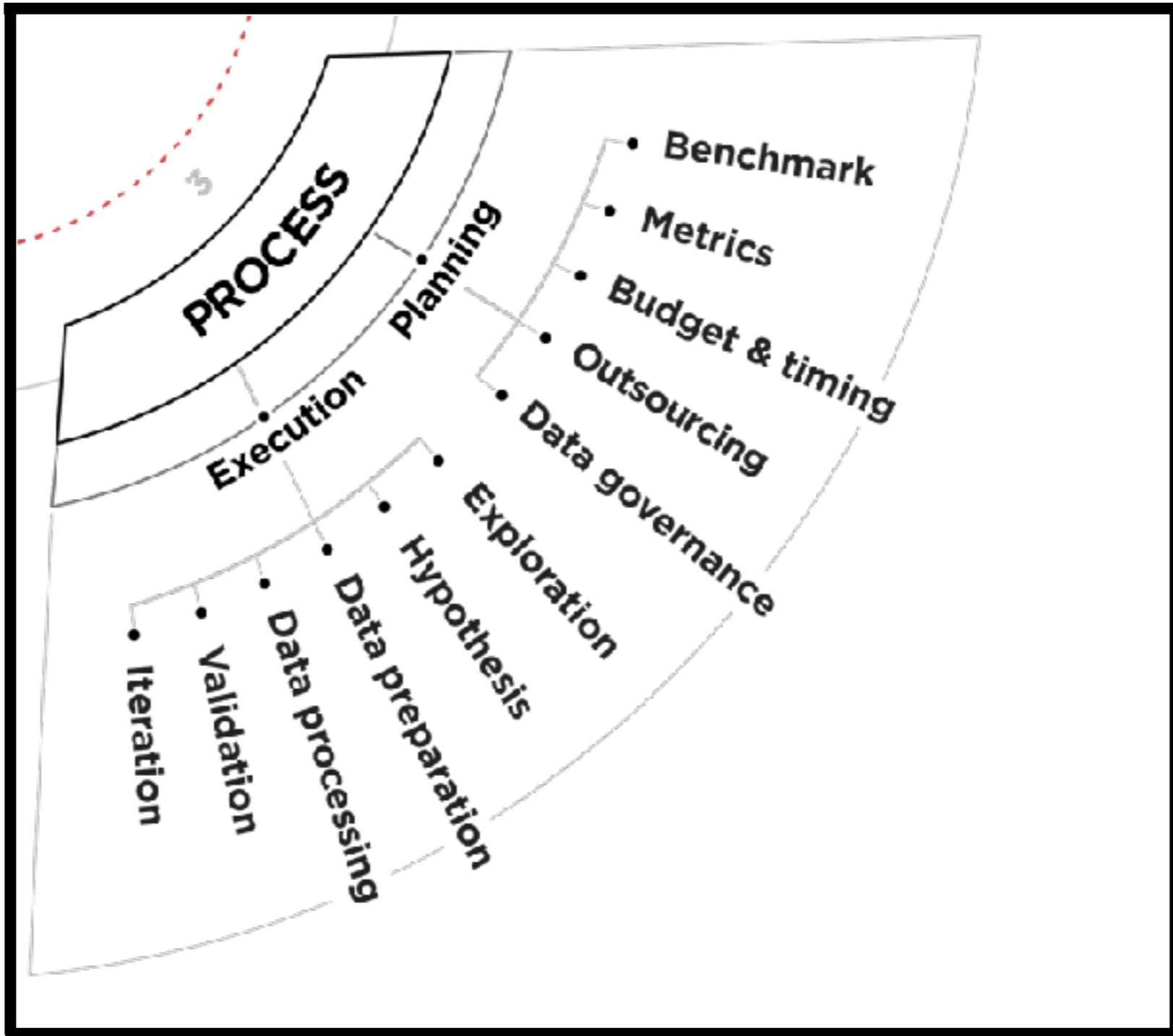
Who is responsible for managing the data?
How?

How do you ensure scientific validation?

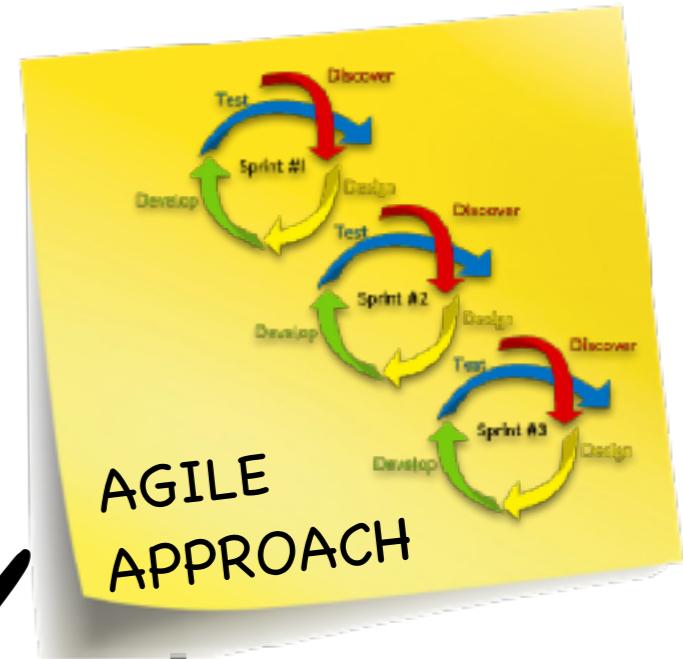
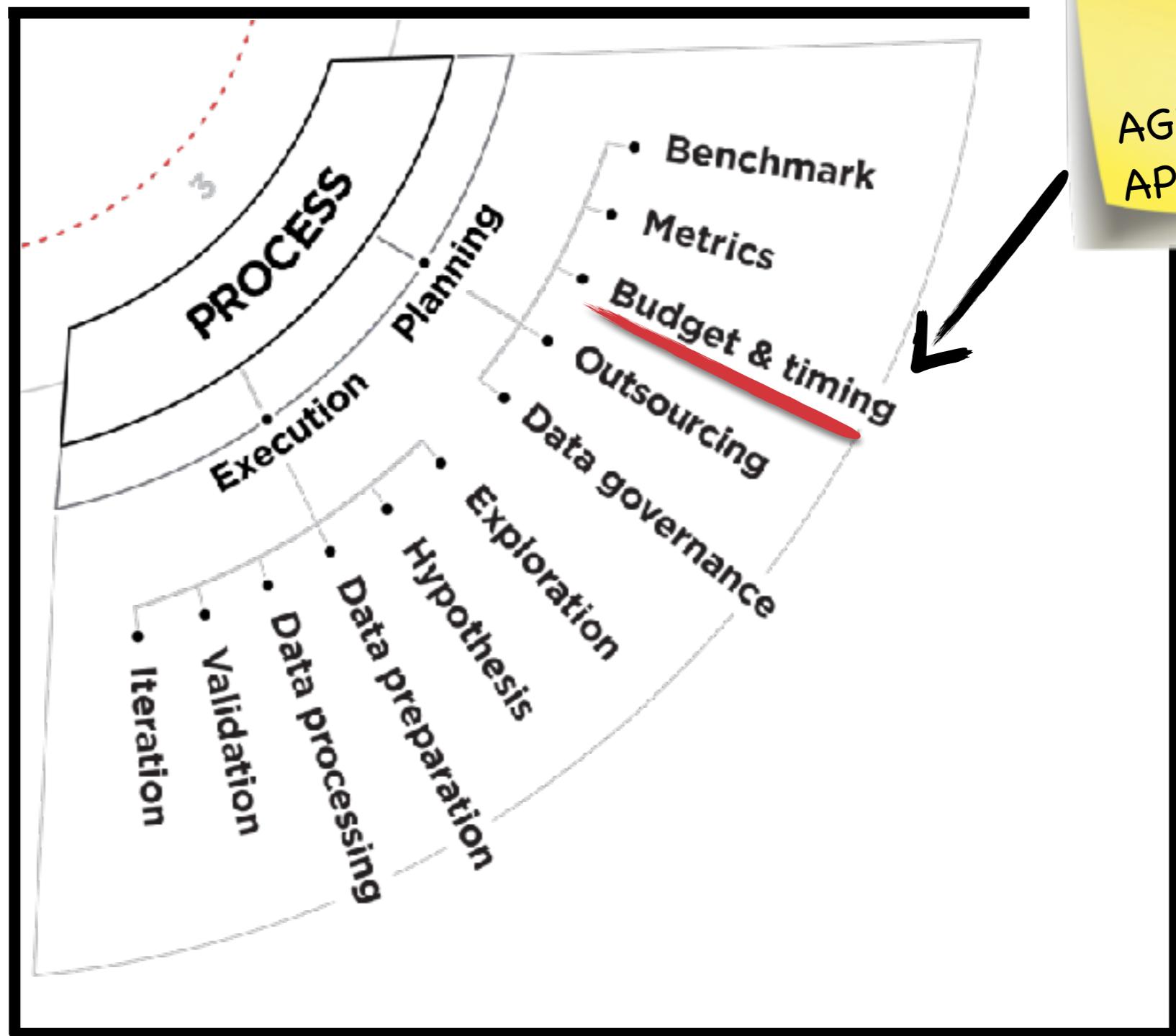
Which are your recruiting channels?

Is there a collaboration between the data team and other business units?

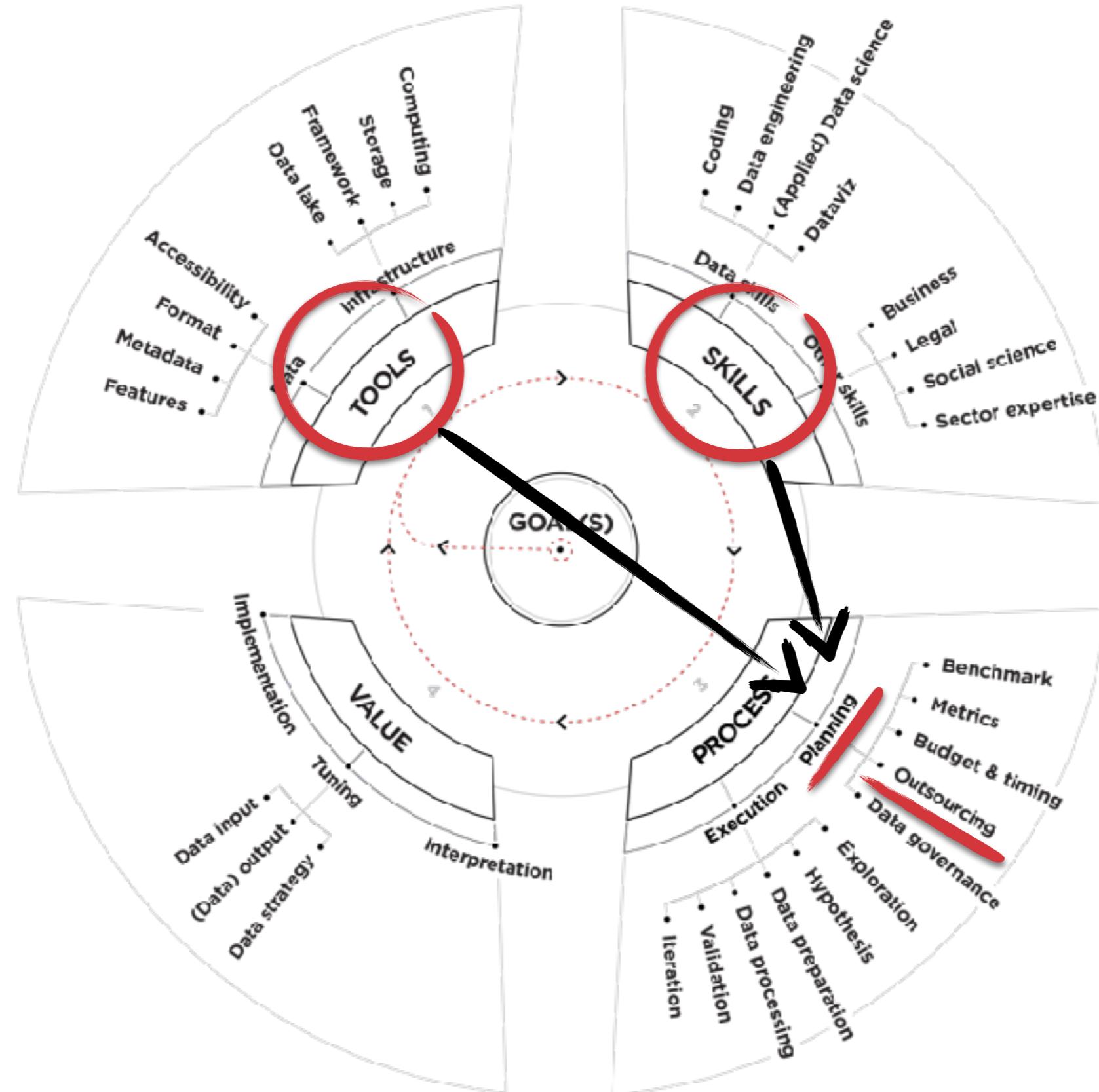
PROCESS



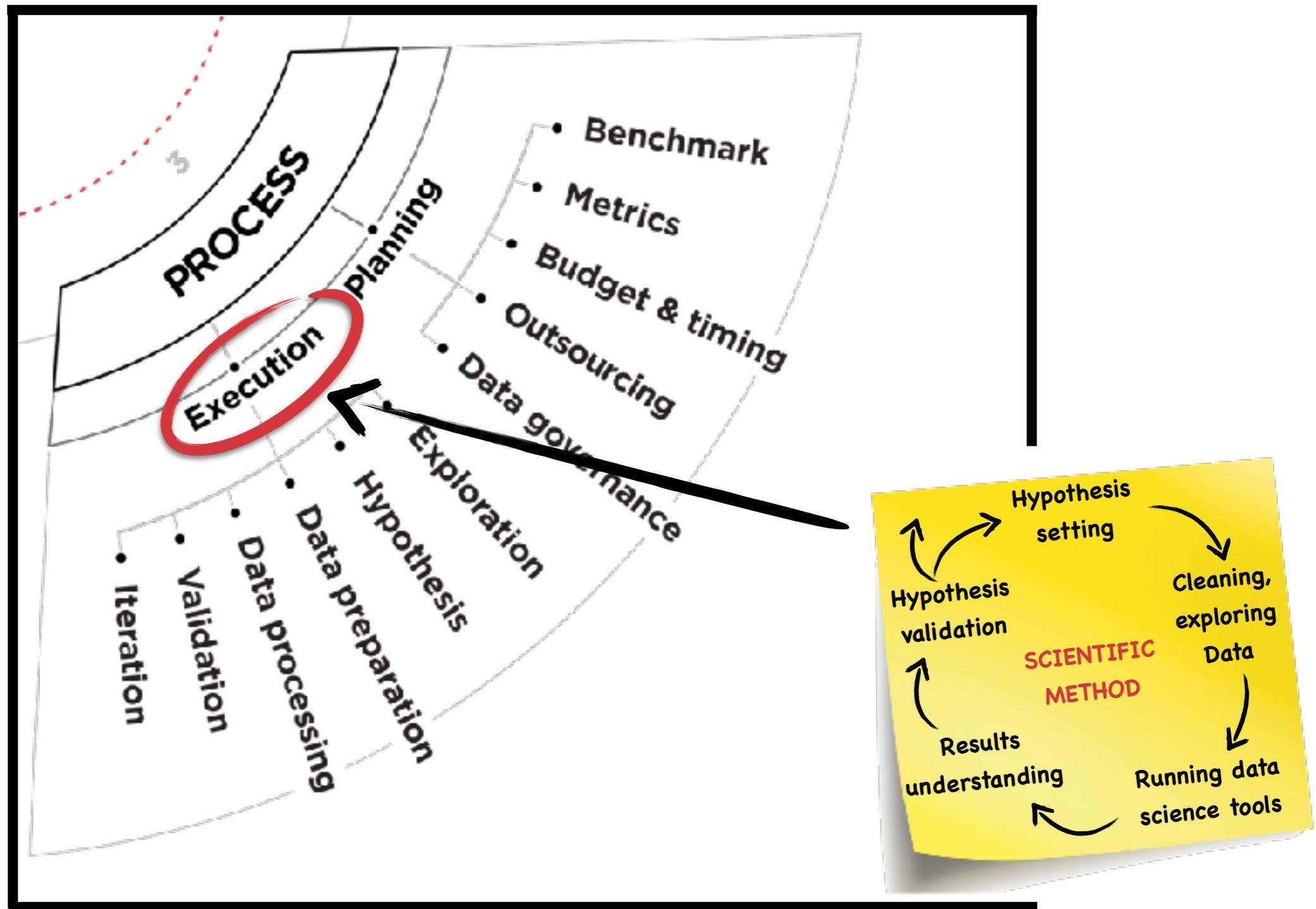
PROCESS



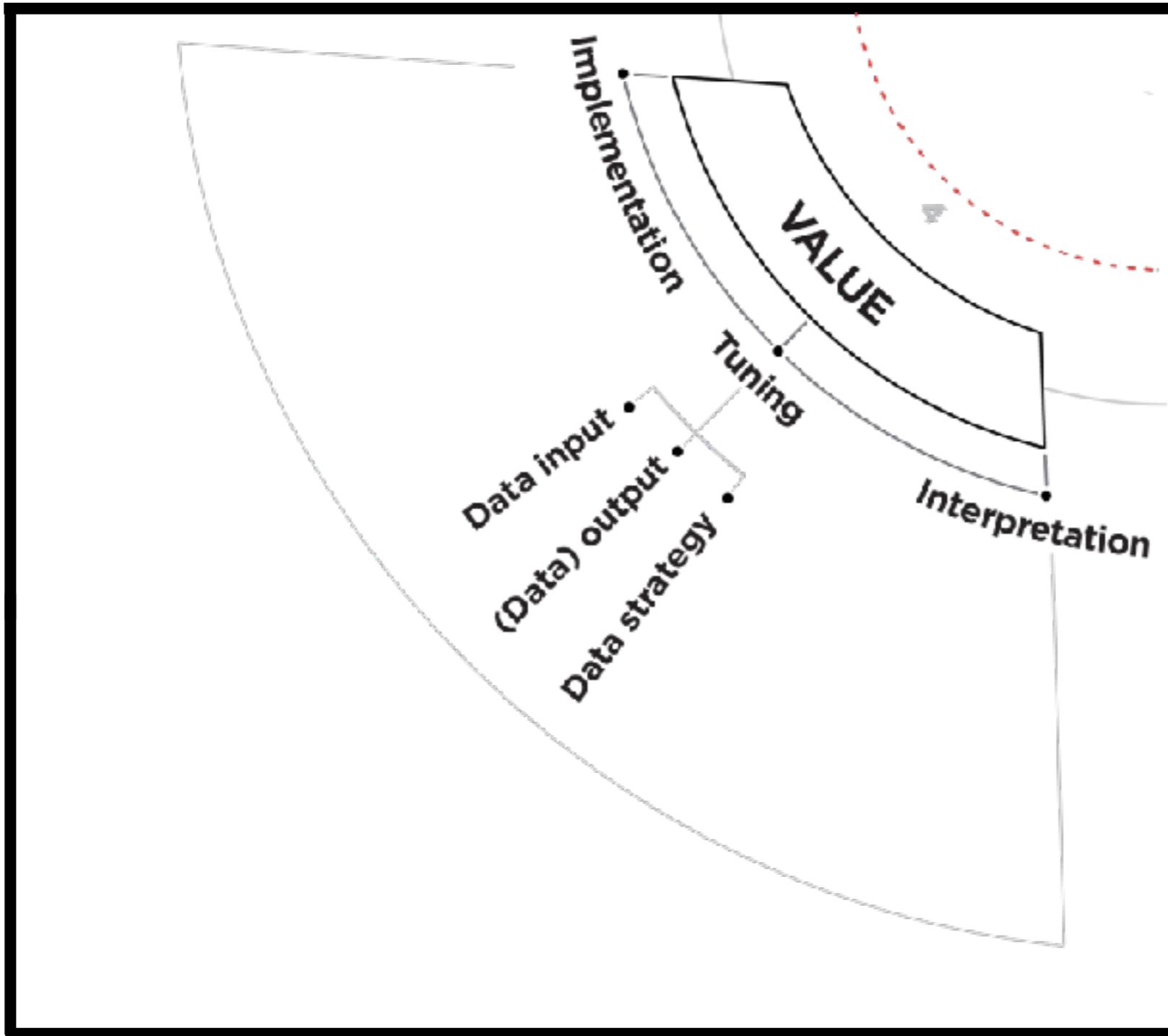
PROCESS



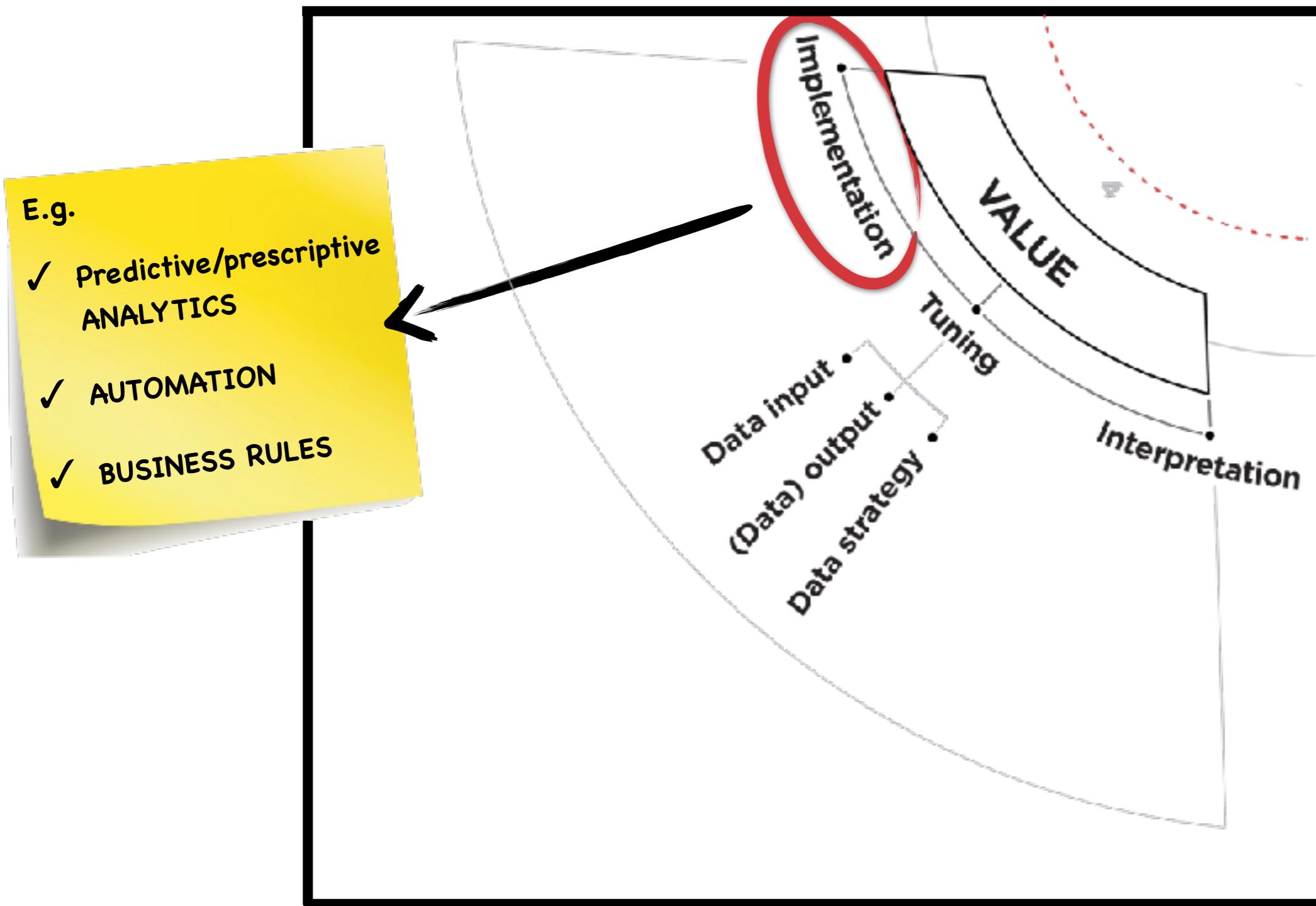
PROCESS



VALUE



VALUE



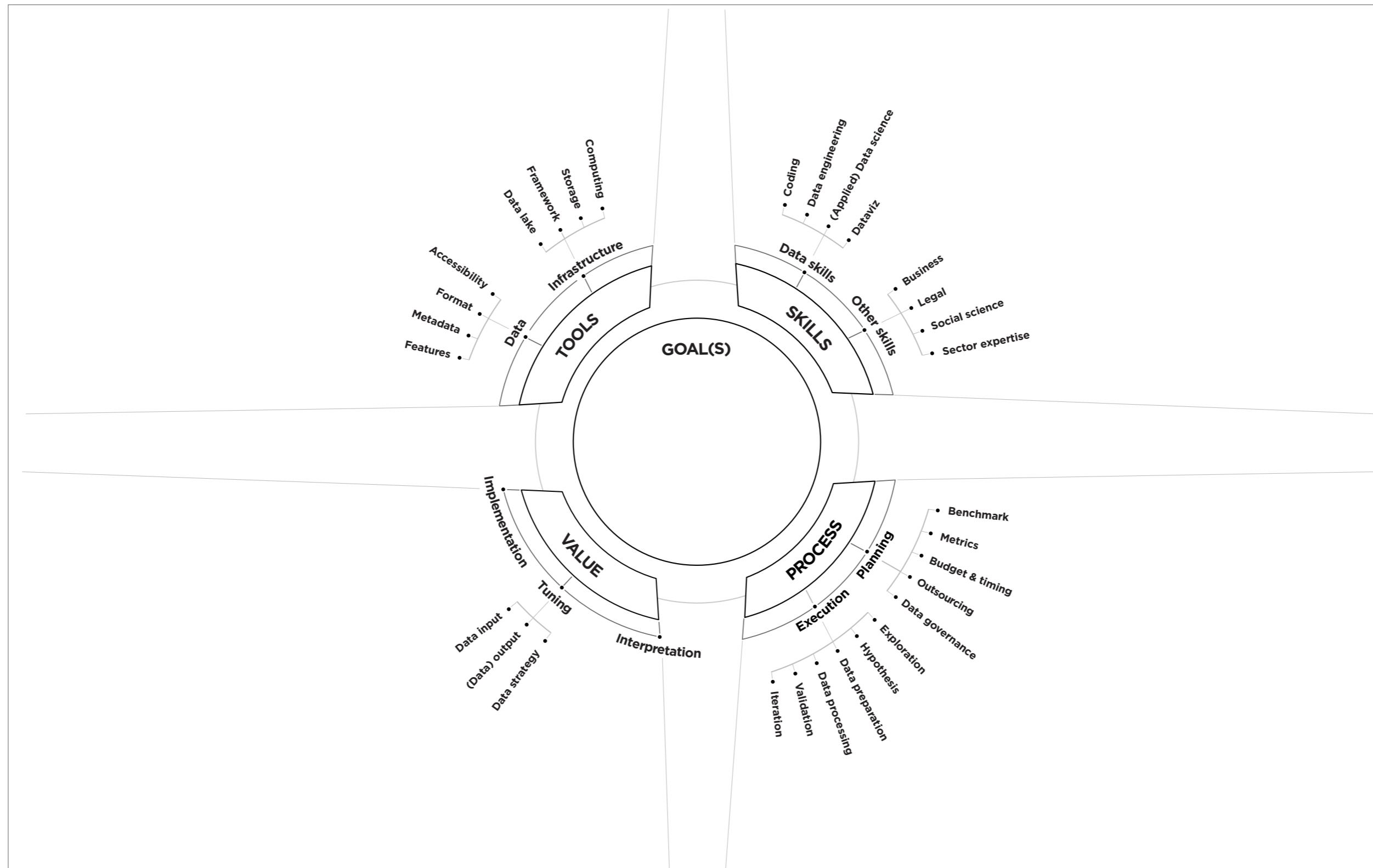
The Data Ring Canvas

Project name:

Designed by:

Date:

Version:



**DOWNLOAD, USE,
COMMENT,
REFINE IT (CC LICENSE)**
- www.dataring.eu-

The Data Ring Canvas

Project name:

Designed by:

Date:

Version:



Q & A

christian.racca@top-ix.org

www.top-ix.org

www.bigdive.eu

@top_ix

@bigdive_eu

THANKS!