

# THE BIG DATA REVOLUTION



VIKTOR MAYER-SCHÖNBERGER

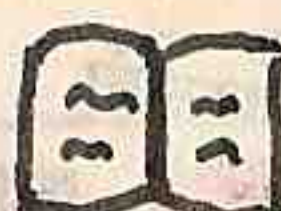
@ Imperial College London



Introduction by Aija Leiponen

**MANY** data providers  
relatively **UNKNOWN**

- acxiom ← **1BN\$**
- CoreLogic
- Datalogix
- e Bureau
- id:analytics
- INTELIUS



**DATA BROKERS**

Report by FTC 2/14



by Viktor

**BIG DATA: A REVOLUTION**

**DELETE: THE VIRTUE OF FORGETTING IN THE DIGITAL AGE**

THE



Prediction case

→ TRADITIONAL APPROACH ←

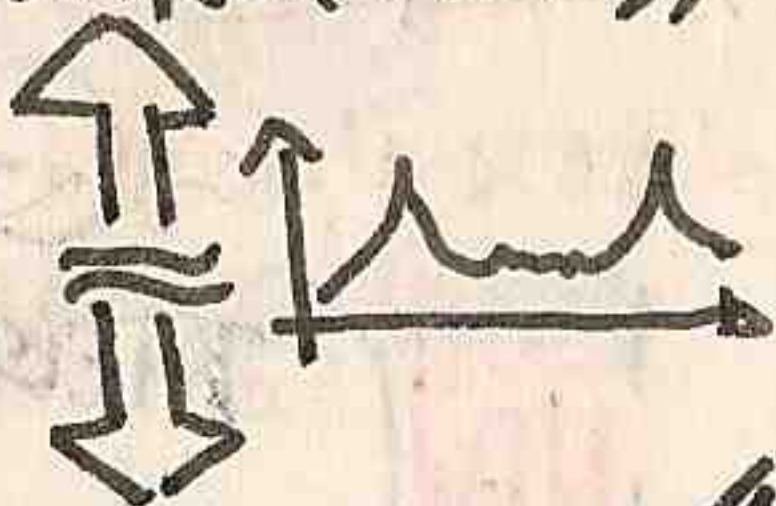


2 WEEKS

VS

→ GOOGLE APPROACH (FOLLOW SEARCH TERMS) ←

**REAL TIME**

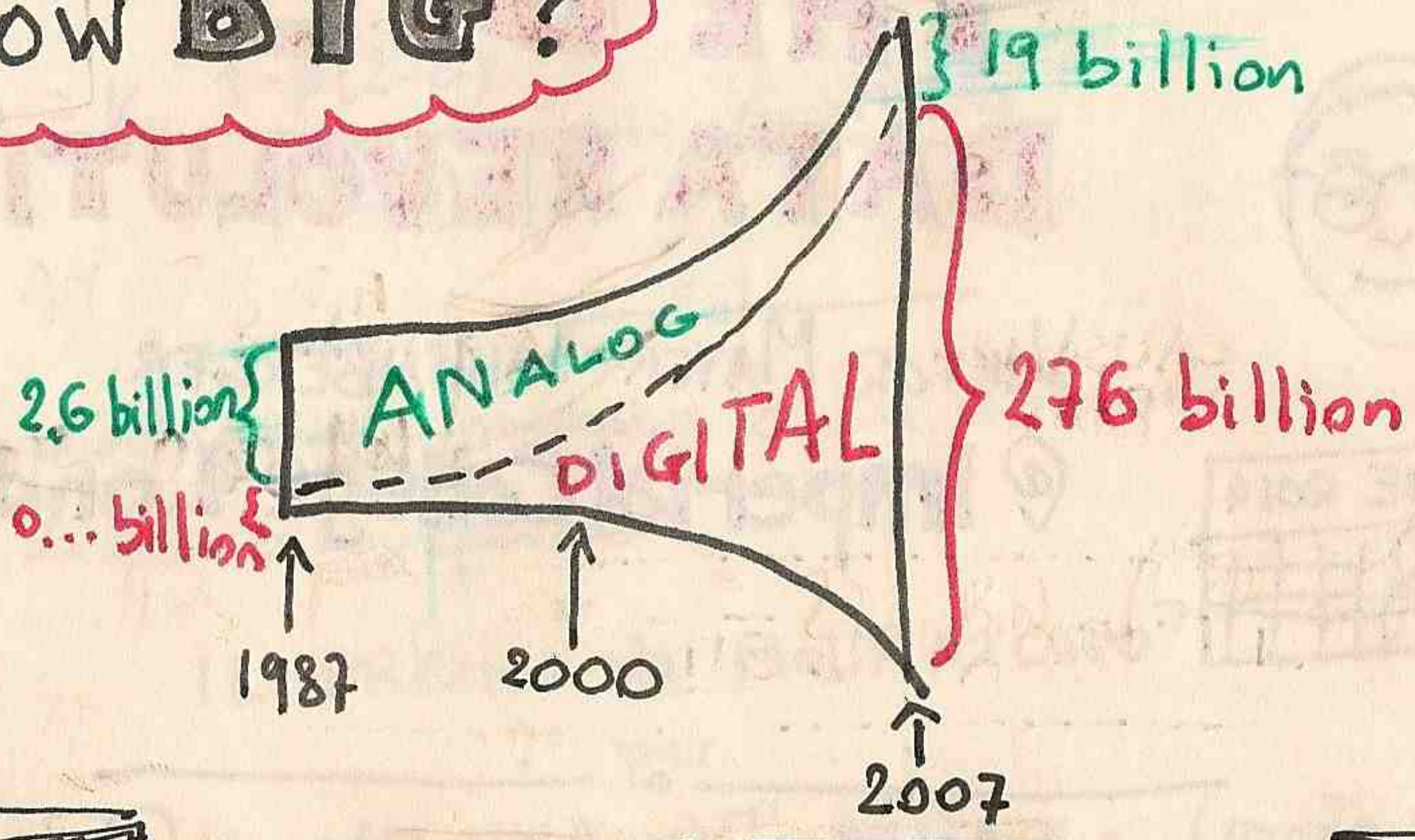


CDC\* data over 10 years

\* CDC = Center for Disease Control



# How BIG?



## THEN and NOW

DATA COLLECTION WAS EXPENSIVE \$\$\$

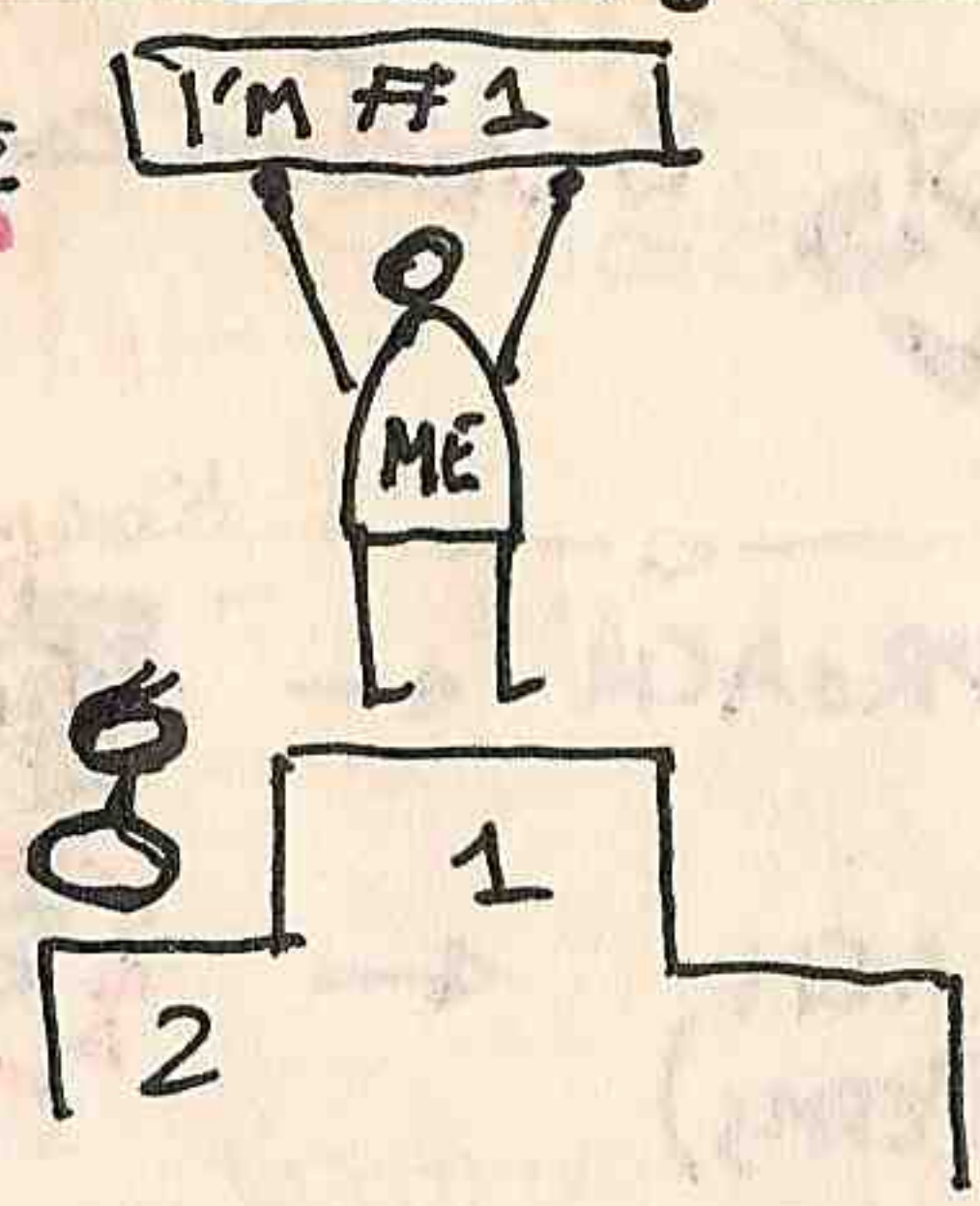
WHAT IF DATA COLLECTION IS CHEAP?

COLLECT DATA JUST AS NEEDED 01001

00110010110  
1010001100  
00100...  
...  
COLLECT ALL YOU CAN

MAYBE RETHINKING STATISTICS

### BACK TO THE FLU CASE

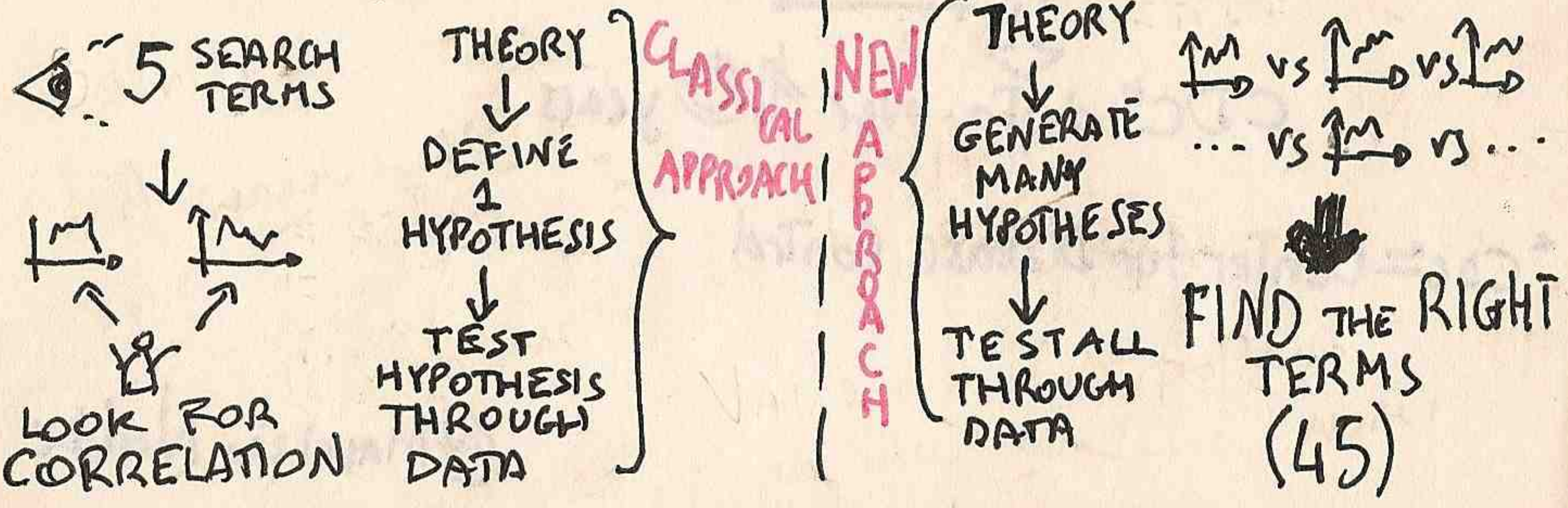


A CANADIAN RESEARCHER CLAIMS TO HAVE BEEN FIRST

BUT... THE APPROACH WAS DIFFERENT!



### Google





SMALL DATA



CAUSALITY BIAS

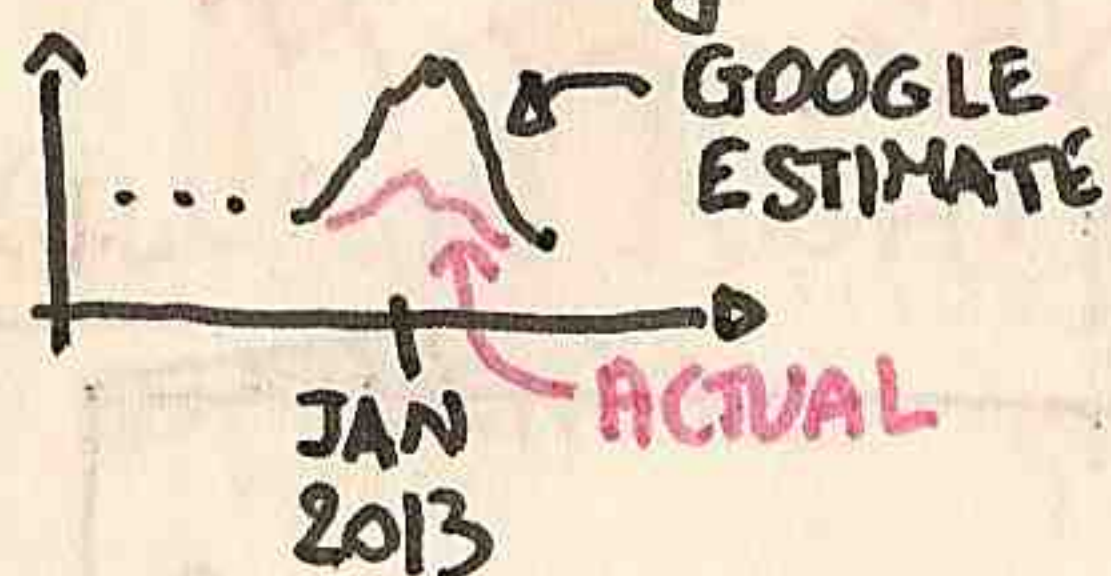
BIG DATA



SPURIOUS SIGNALS

THEORY IS STILL ALIVE!

BUT Google was WRONG in winter 2012-2013



WHY?

MODEL BUILDING

MODEL APPLICATION

2009 ← 4 yrs → 2013

...NO UPDATE...

UPDATE THE MODEL!  
(BAYES)



SMALL DATA

DATA COLLECTION was TARGETED



UNUSABLE AFTERWARDS

BIG DATA

REUSE DATA

DATA ≠ CONTENT

REUSE MUST BE INCENTIVIZED

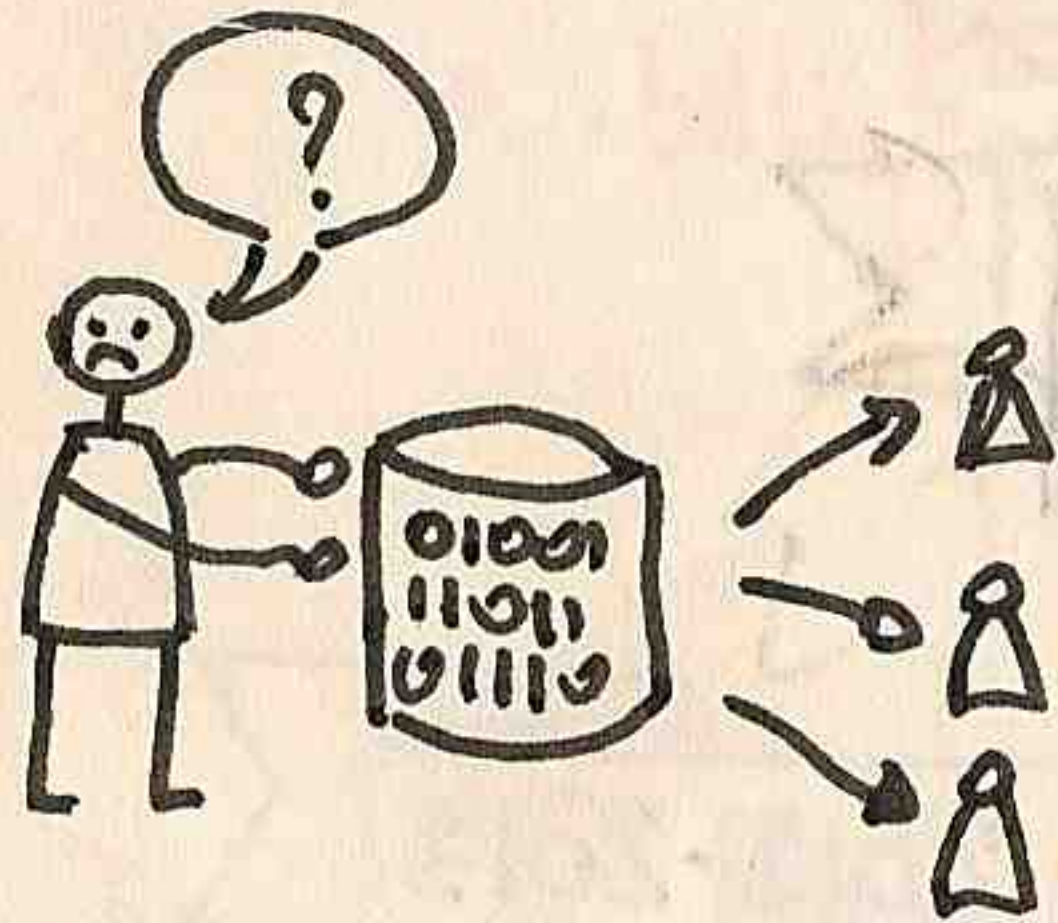
REUSE INCENTIVIZED by INTELLECTUAL PROPERTY LAW



PROTECTION  $\equiv$  EXCLUSION  $\Rightarrow$

AR  
\$\$\$  
E

NO EXCLUSION  $\Rightarrow$  NO INCENTIVE TO PRODUCE



THE DATA HOLDER'S DILEMMA

TO SHARE  
OR  
NOT TO SHARE?

HOW TO INCENTIVIZE REUSE/SHARING?

LOW TRANSACTION COSTS



TRUST BY DATA HOLDERS

CREATE DATA USER OPPORTUNITY

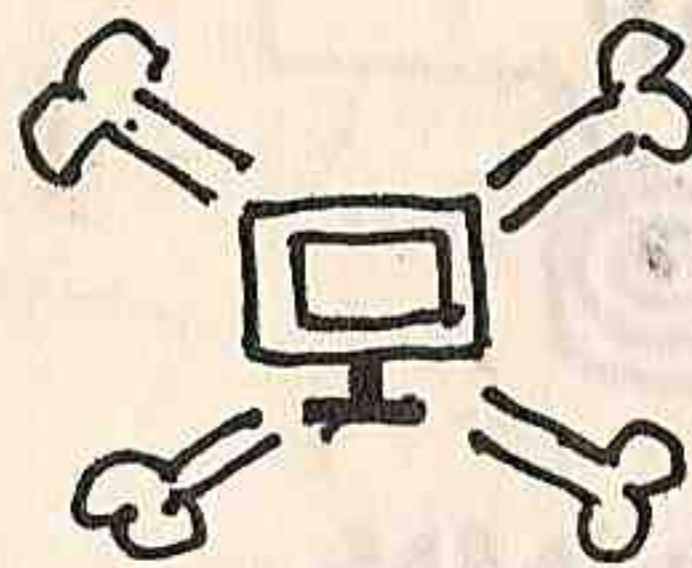
3 POSSIBILITIES

CONTRACTUAL ARRANGEMENTS

PROPERTIZATION

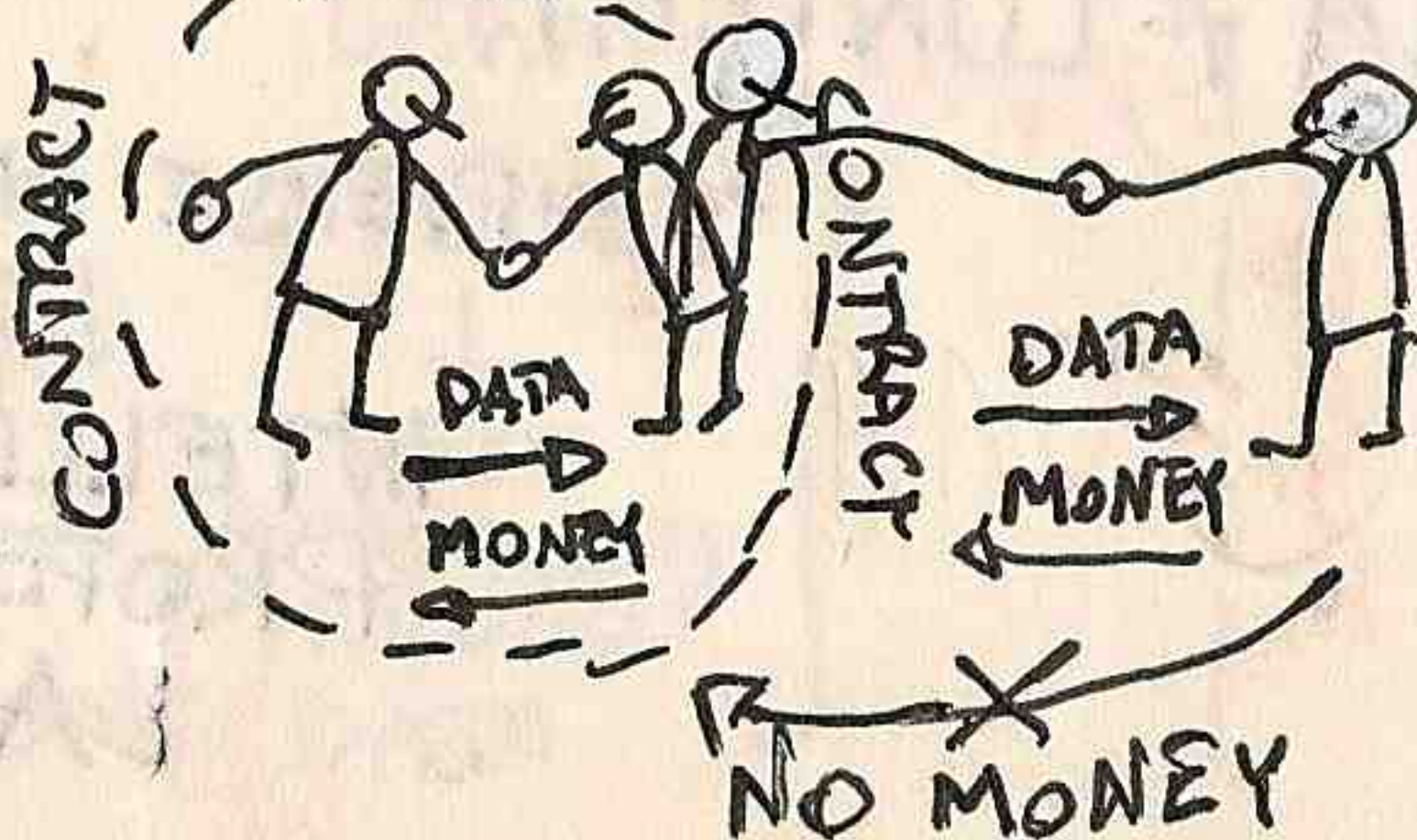
EXCLUSION RIGHT SUI GENERIS

CONTRACTUAL ARRANGEMENTS



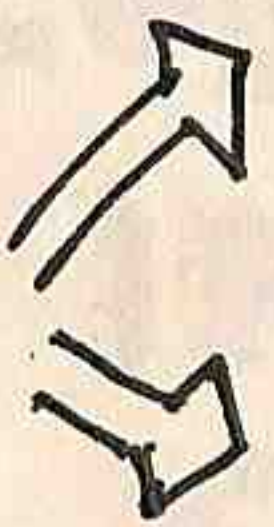
DATA PIRACY IS DIFFICULT TO DETECT

LIMITED TO CONTRACT PARTIES

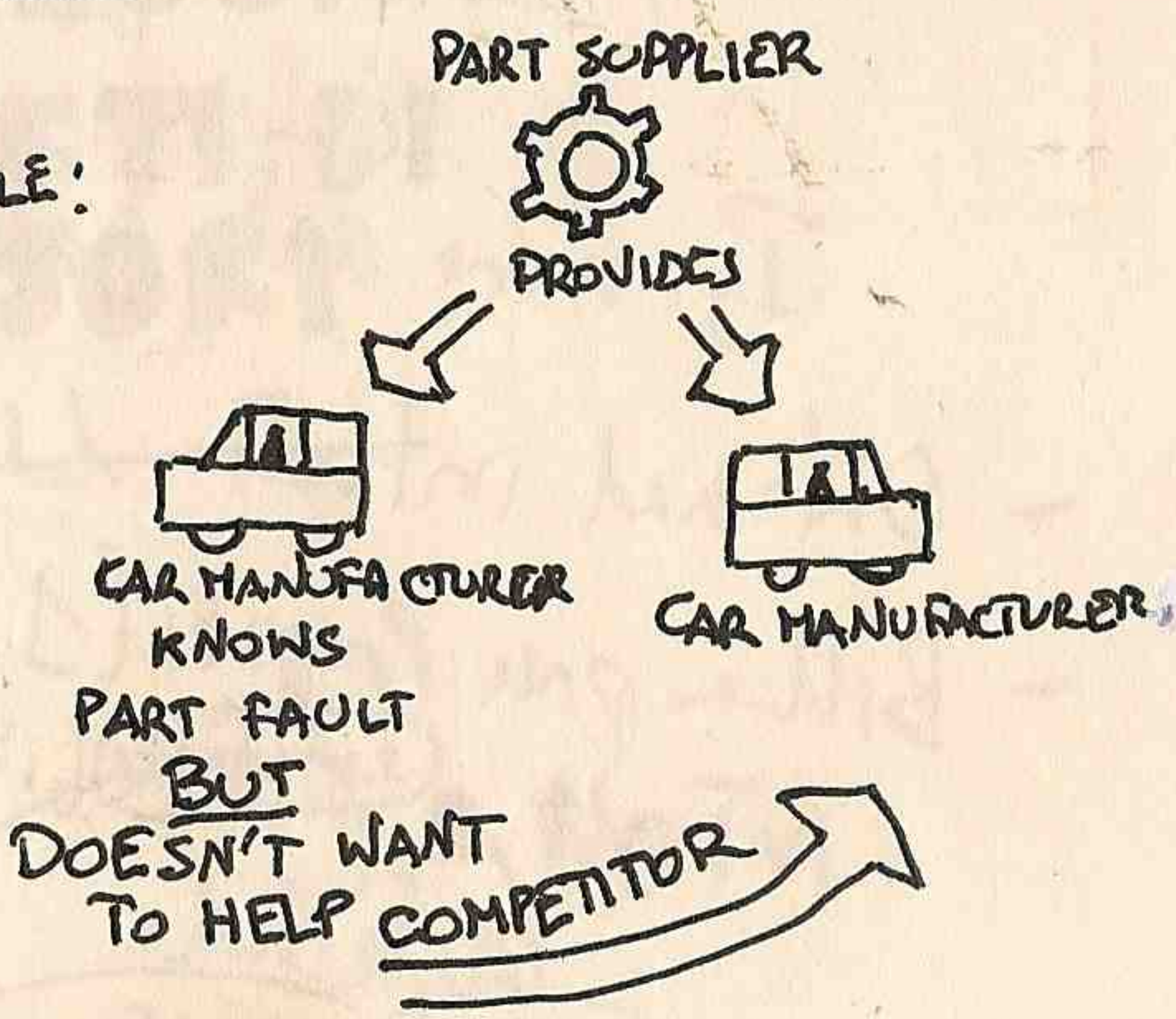




PROPERTIZATION || NOT GRANULAR



EXAMPLE:



EXCLUSION RIGHT || HAVE PROBLEMS  
SUI GENERIS

+ TECHNOLOGY? ⇒ LABEL (CC)  
⇒ ENFORCEMENT (DRM)  
CREATE A GLOBAL SURVEILLANCE INFRASTRUCTURE



WHAT CAN WE DO?

1. SHARING OBLIGATIONS (FRAND) ⇒ SEE ITA-GOOGLE CASE
2. "OPEN DATA" MANDATES ⇒ DATA PRODUCED AFTER GOVERNMENT FUNDING MUST BE PUBLIC