



# TRANSFORM PROD DATA INTO DEV DATA

3 PRACTICAL STRATEGIES

2023/06/27

DR. REMI CURA, PRINCIPAL DATA SCIENTIST, CENTAUR LABS  
REMI.CURA@GMAIL.COM



# TODAY'S TOPICS

## BENEFITS OF

- USING A DEV DATABASE
- ... BASED ON PROD DATA
- ... THAT HAS BEEN OBFUSCATED

*A TESTIMONY FROM THE MEDICAL TECH SPACE*



TESTIMONY □ NEED CONTEXT

AKA, HOW CRAZY AM I?



# CONTEXT: WHO AM I

GAME: PUT EVERYTHING IN POSTGRES

PHD : COMPUTER SCIENCE / GIS

VECTOR + TOPOLOGY

POINT CLOUDS (LIDAR, STEREOVISION)

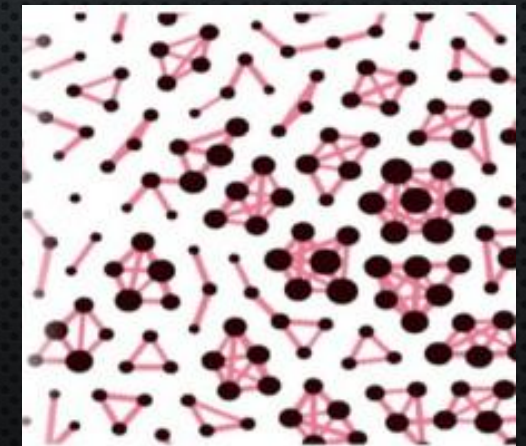
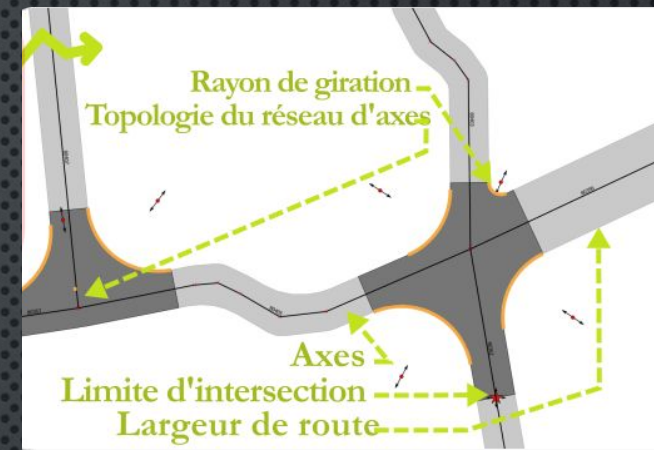
POSTGRES AS A BACKEND? (ML, STREET MODELING)

POSTDOC: PARIS SCHOOL OF ECONOMICS

HISTORICAL MAPS + GEOCODING IN POSTGRES  
(FUZZY PLACE + TIME)

MIT: POLITICAL SCIENCE DEPARTMENT ( LOBBYVIEW )

MONEY IN POLITICS : RECORD MATCHING, GRAPHS  
(LOBBYING+CAMPAIGN+ECON)







# CONTEXT: CENTAUR LABS

[CENTAURLABS.COM](https://centaurlabs.com)

TECH STARTUP (20 EMPLOYEES, 10 DEVS, SERIE A)

WISDOM OF THE CROWD FOR MEDICAL DATA

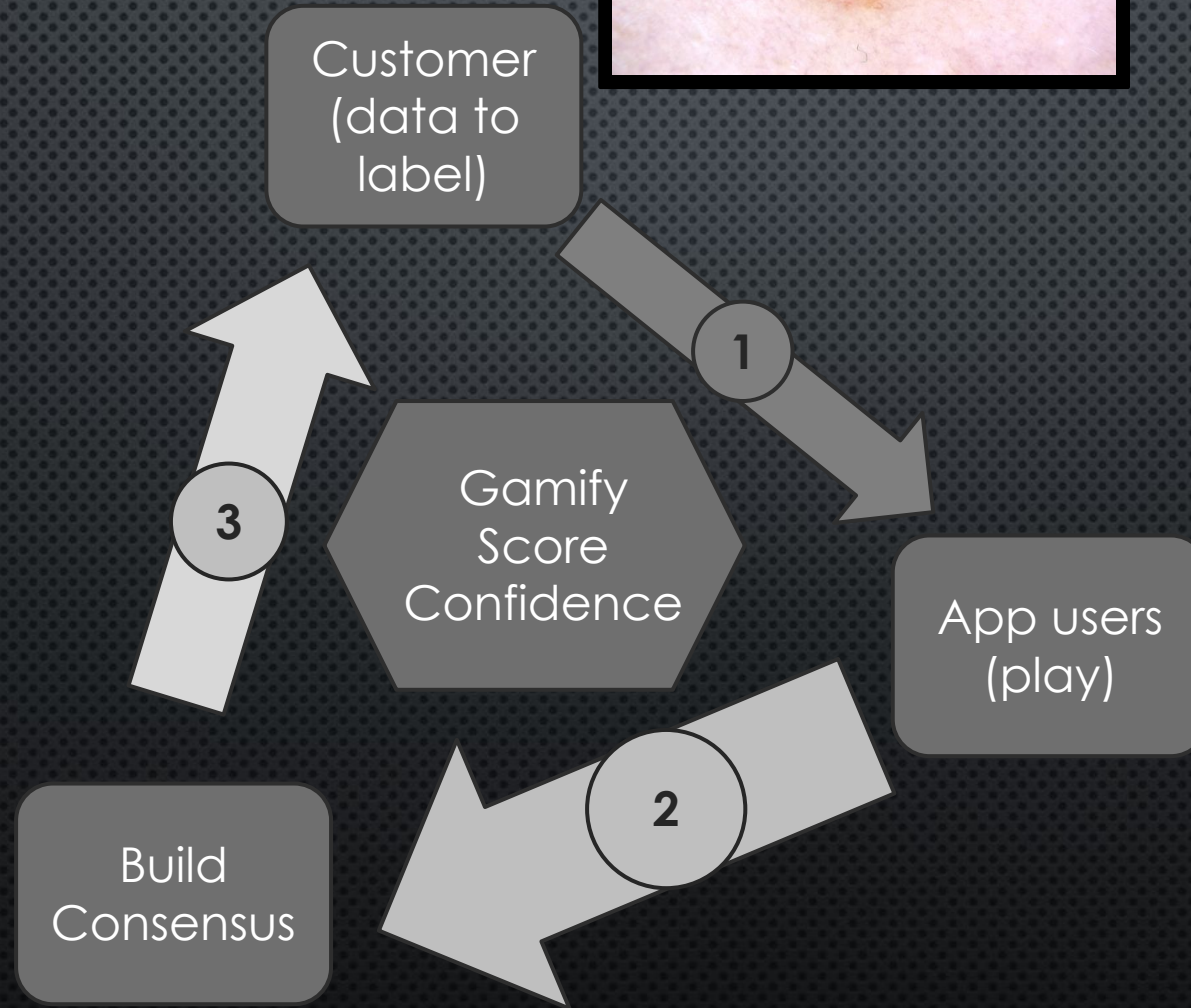
GAMIFIED THROUGH AN APP : DIAGNOSUS.COM

~100 TABLES, 1000 COLUMNS, 300GB, ~400M ROWS



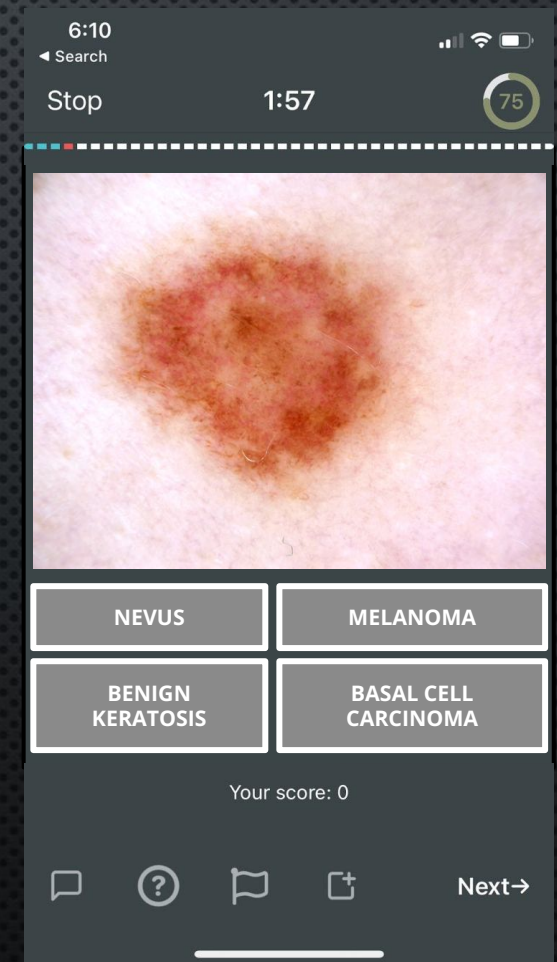
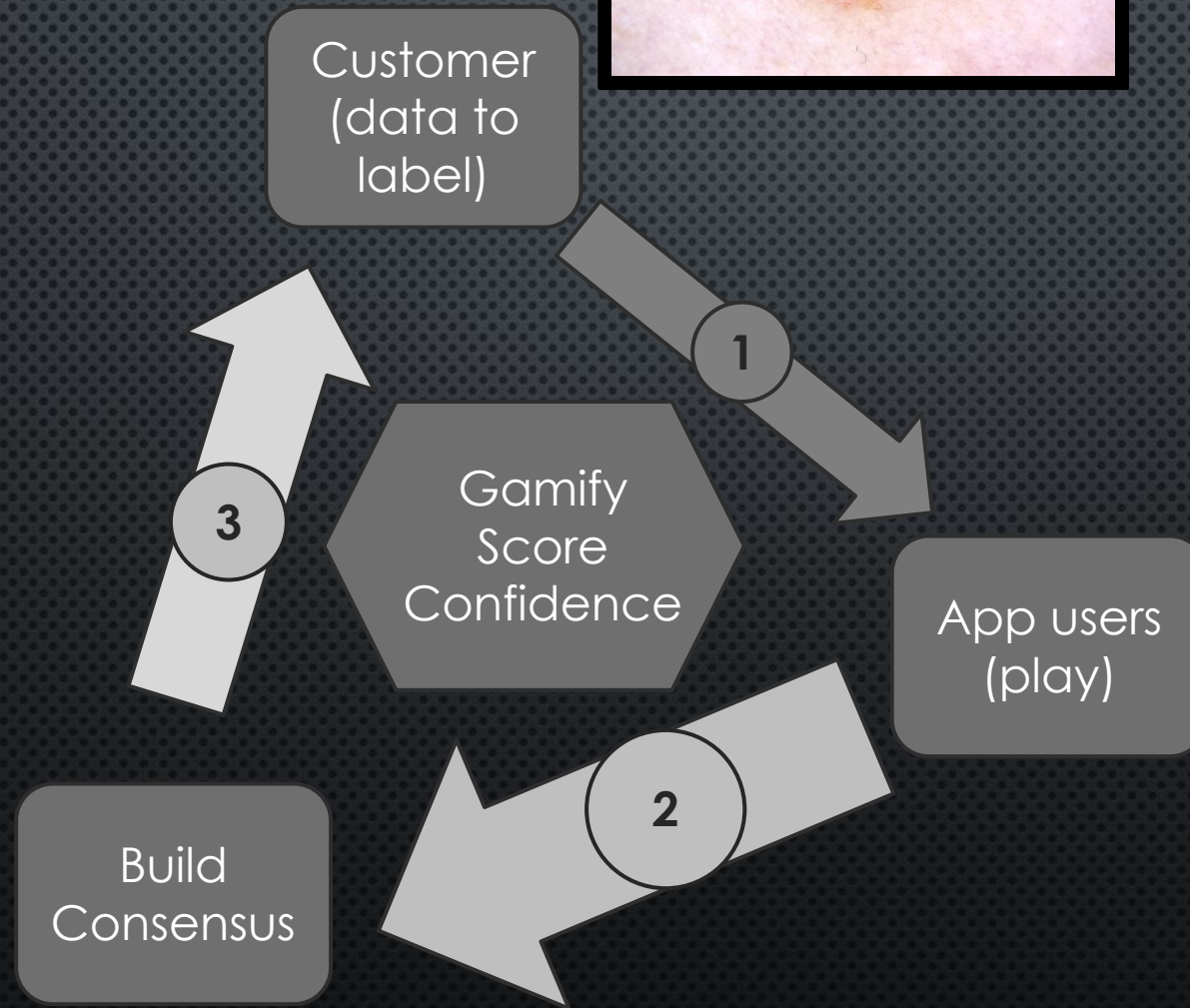


# CONTEXT: CENTAUR LABS



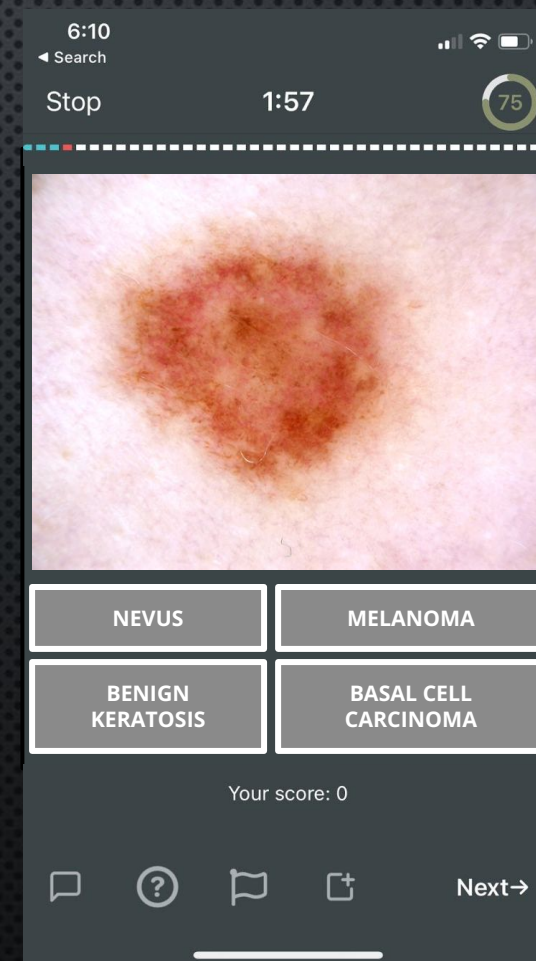
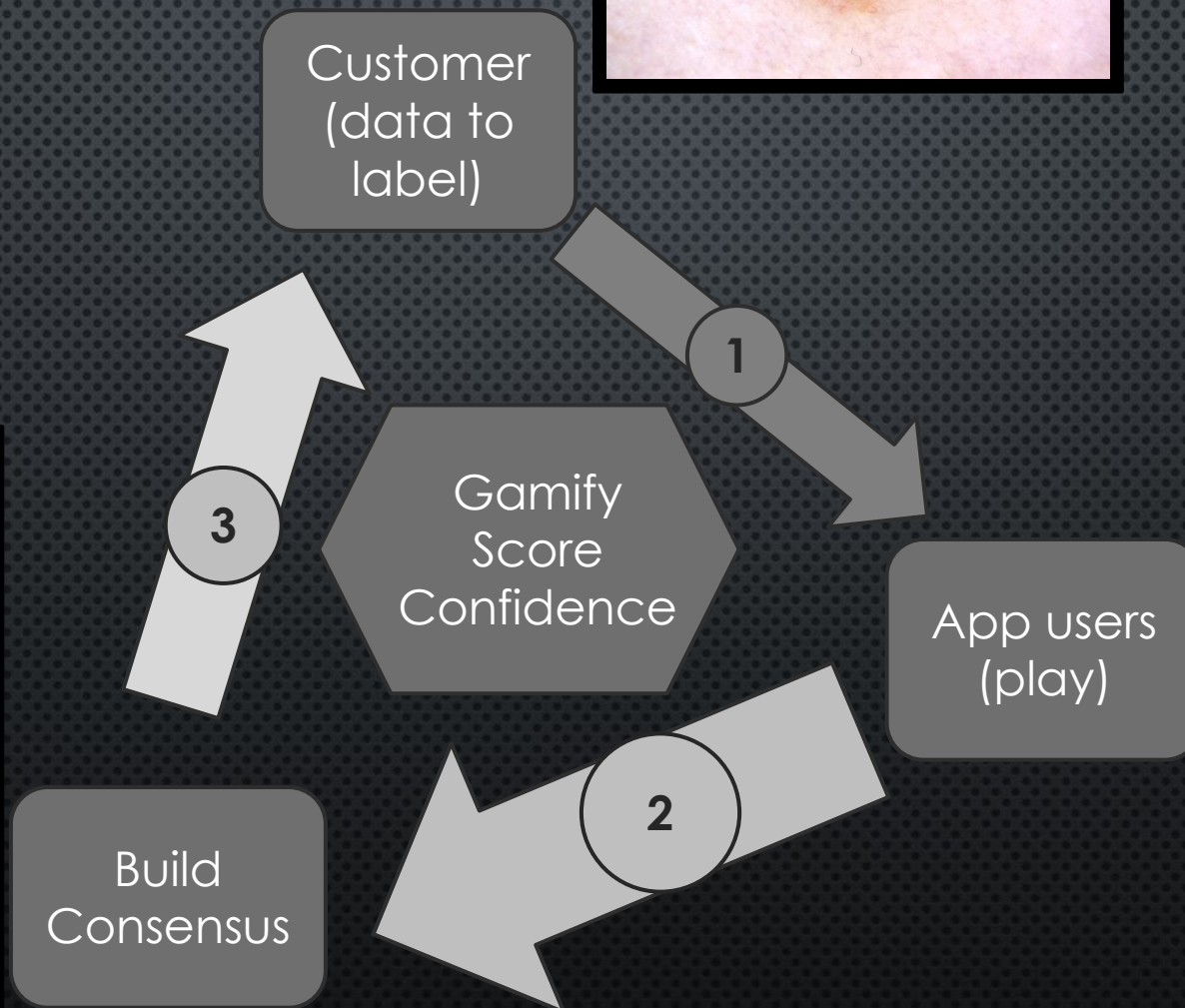


# CONTEXT: CENTAUR LABS





# CONTEXT: CENTAUR LABS



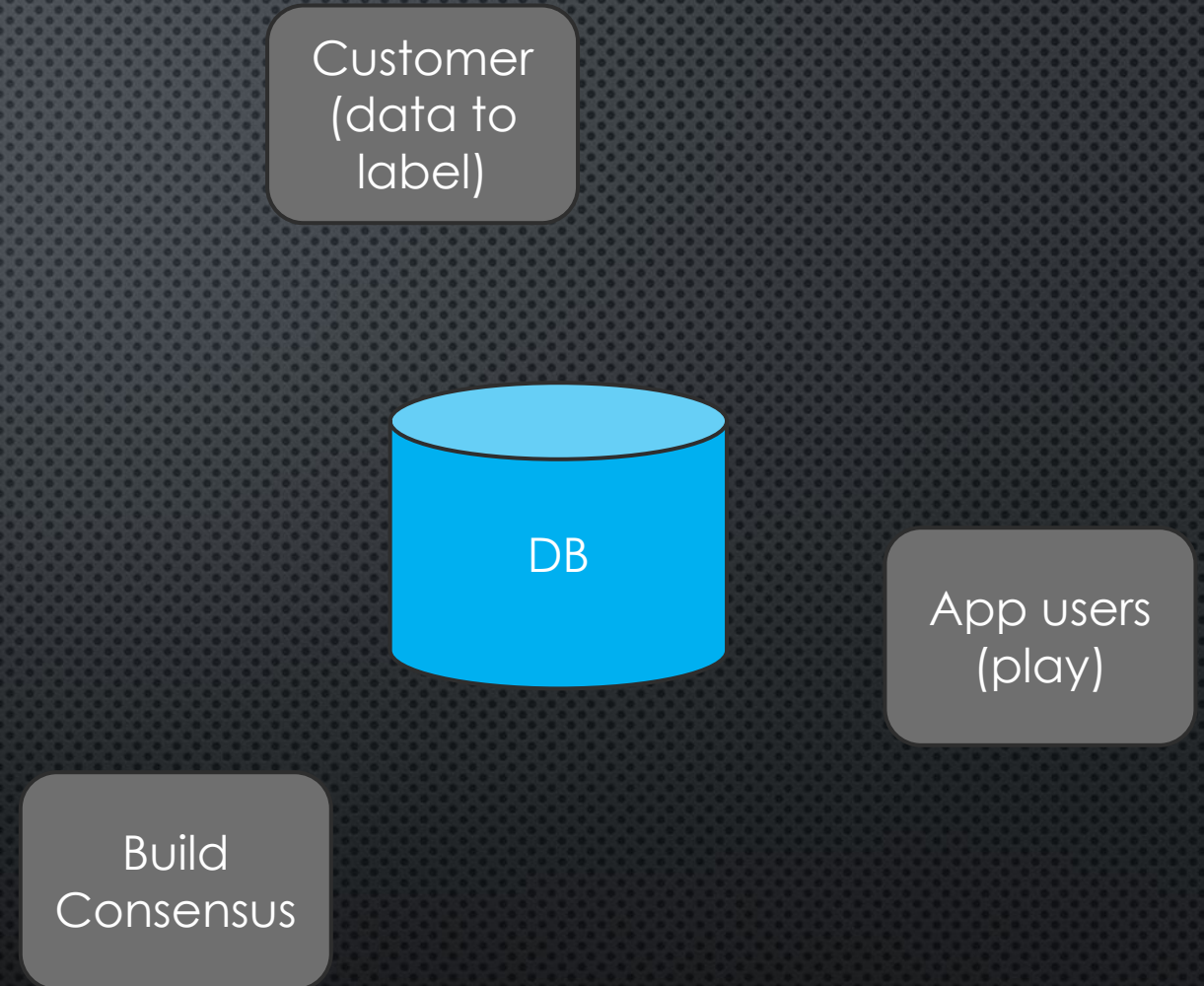


# CONTEXT: SUMMARY

DB-CENTRIC: 8/10 DEVS WRITE SQL

DB IS NOT ONLY DATA STORAGE

- DATA TRANSFORM
- BUSINESS LOGIC
- CONVENIENCY/HELPER FUNCTIONS





WHY A DEV DB?



# WHY A DEV DB? SAFER 1/2

- CASCADING SAFETY
  - PARACHUTE: 1/5K CHANCE OF USING RESERVE
    - ~ ODDS OF DEATH BIKING
  - PARACHUTE: 1/220K CHANCE OF DYING
    - ~ ODDS OF DEATH BY LIGHTNING
  - MY OWN DBA STATS AT CENTAUR LABS:
    - DEV DB: ~1 ERROR /WEEK
    - PROD DB: ~1 INCIDENT/YEAR
- EASIER UNDO
  - MEH, WHY BOTHER UNDOING
  - RESET WHOLE DEV DB
    - EVERY 2 WEEKS OR MANUALLY
  - OK I'LL UNDO ... WITHOUT TIME PRESSURE/STRESS



Image



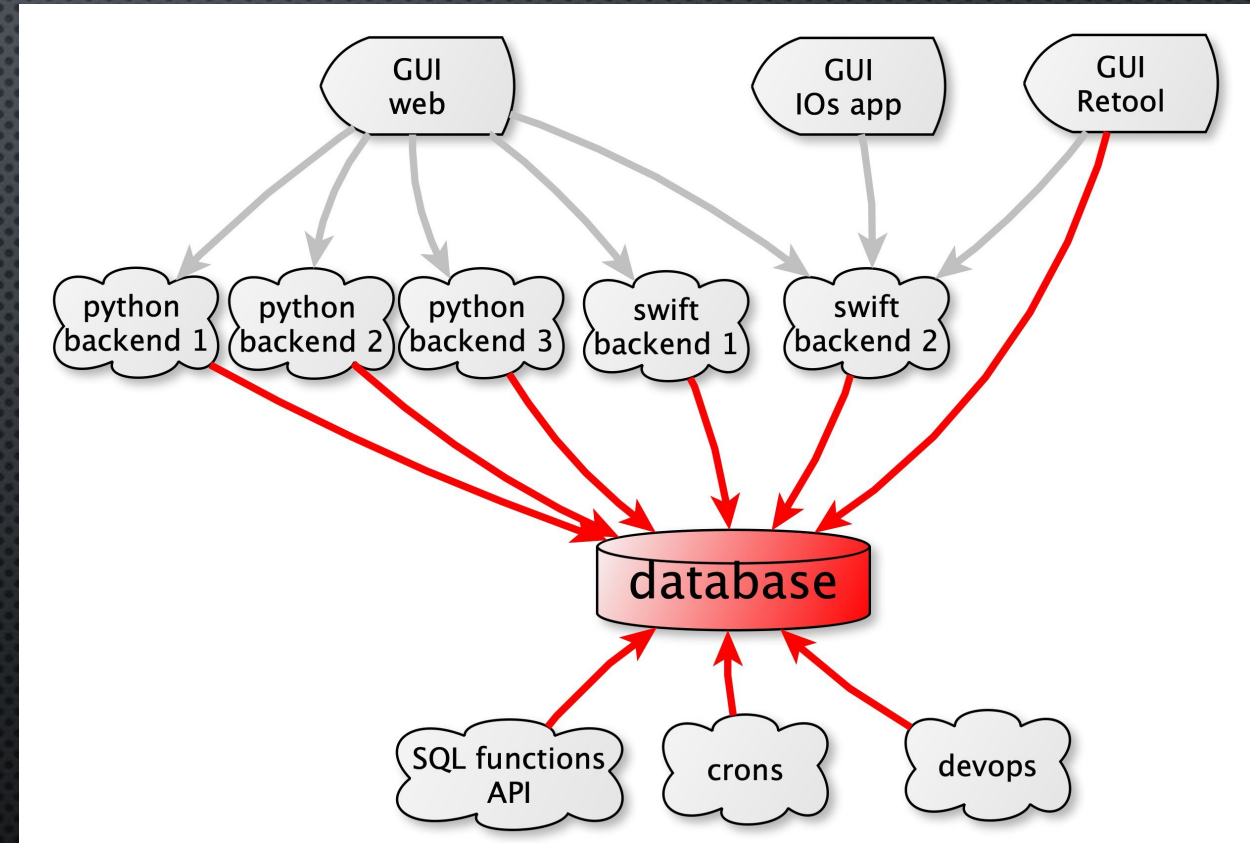
# WHY A DEV DB? SAFER 2/2

- **ANY** CHANGE : HOW DO YOU KNOW IT WORKS?
- TESTING !
  - “THIS FUNCTION SHOULD DETECT BAD DATA”
  - ☐ I NEED BAD DATA TO TEST IT
  - ☐ PUT BAD DATA IN PROD??
  - PROD: ALERTING, LOGGING, ... ☐ DON'T WANT TO TRIGGER THAT
- DEV DB IS THE PERFECT PLACE FOR THAT



# WHY A DEV DB? EMPOWERS DEVS

- 3 GUI's , 5 BACKENDS ...
- ☐ NEED EMPOWERED DEVS
  - WE LEARN BY MAKING MISTAKES
  - QA/ CODE REVIEW
    - MORE POWERFUL TO REVIEW SCHEMA+DATA THAN CODE IMO
  - PERMISSIONS
    - 3 DEVS HAVE DDL PERMISSIONS ON PROD
    - 9 HAVE DDL SKILLS AND PERMISSIONS ON DEV
- ☐ FIRST IDEA != BEST IDEA
  - FREEDOM TO EXPERIMENT





# WHY A DEV DB? CONFIDENCE = SPEED



Image



Image

ON WHICH PATH CAN YOU GO FAST?

WHY?

☐ CONSEQUENCES



# WHY A DEV DB? CONFIDENCE = SPEED



Image



Image

**PROD** ☐ **IN USE!**

**ANY** QUERY CAN SLOW IT /BRING IT DOWN

- DDL, DELETE, UPDATES
- SELECT (RESSOURCES, LOCKS)

REDUCE CONSEQUENCES ☐ FASTER DBA WORK



# WHY A DEV DB? TAKE AWAY



## DEV DB:

- SAFER
  - CASCADED SAFETY
  - TESTS
- EMPOWERS DEVS
- MORE CONFIDENCE, MORE SPEED

*WHAT DEV DATA?*



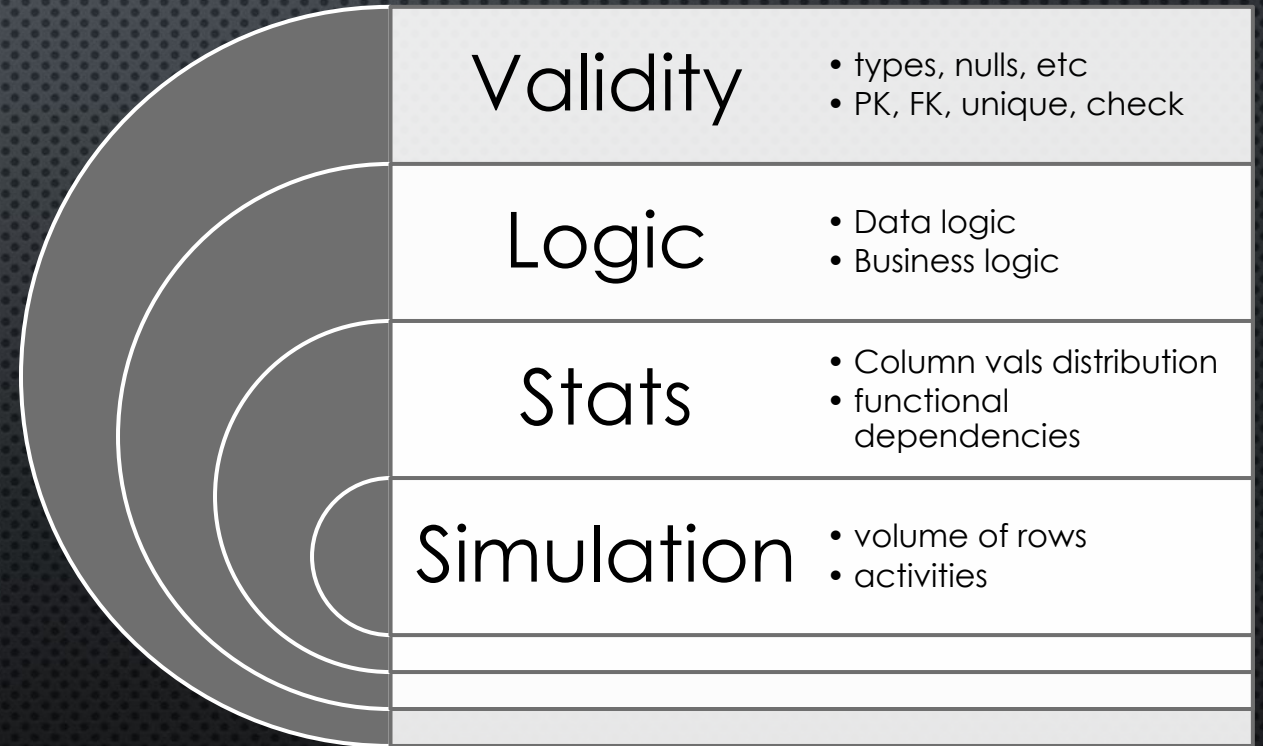
# WHAT DATA FOR THE DEV DB?

SYNTHETIC VS PROD DATA



# SYNTHETIC DATA

- SYNTHETIC: INSERT MADE-UP STUFF
  - WHAT IS “MADE-UP”?
  - DEPENDS ON USE CASE
    - **ALL** : VALIDITY
    - **DEV**: LOGIC
    - **DATA SCIENCE**: STATS
    - **DEVOPS**: SIMULATION





# SYNTHETIC DATA

- PRO:
  - PERFECT FOR TESTS !
  - FOR EACH COMPONENT: MAKING UP RARE DATA
    - TESTING SQL INJECTION ROBUSTNESS
    - BAD GEOMETRIES COLLECTION
  - FOR ALL COMPONENT: INTEGRATION TEST
    - COVER "BASIC"/ "MOST COMMON" SCENARIO

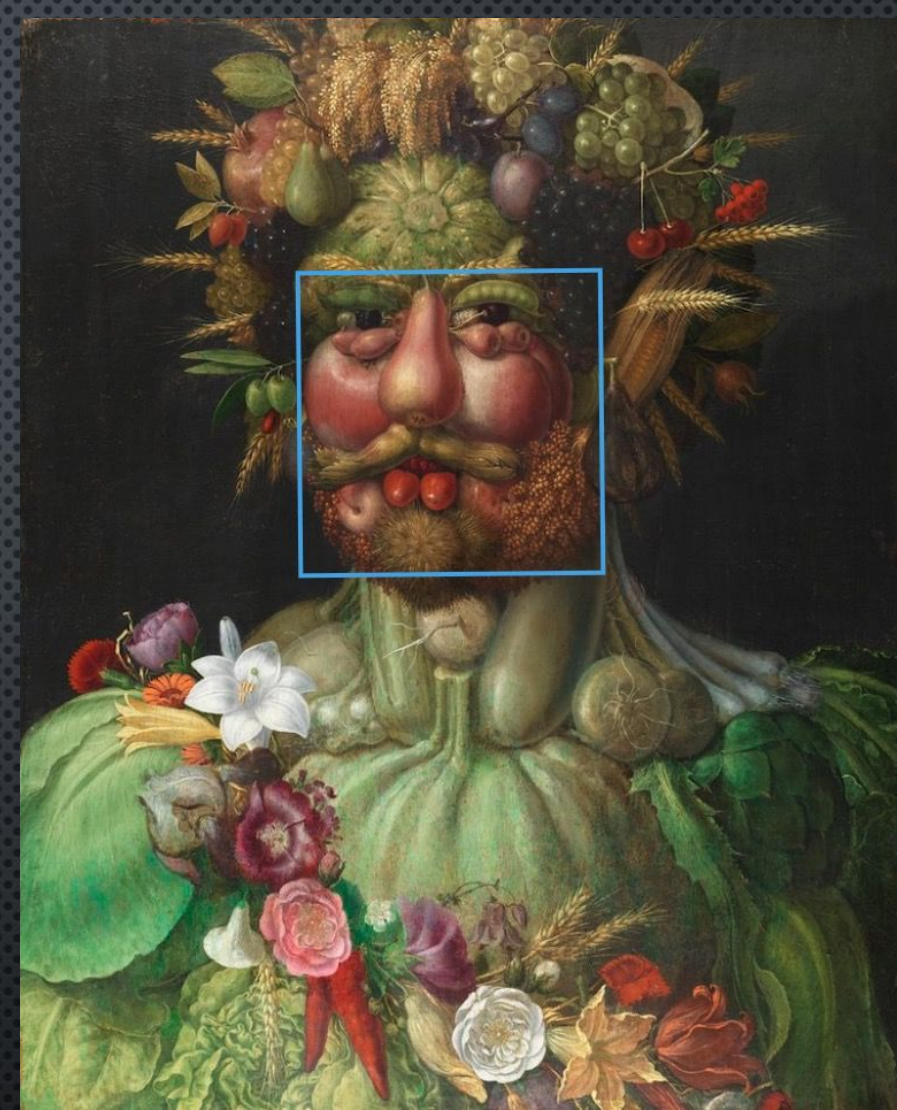


Image with face detection service



# SYNTHETIC DATA

- CONS:
  - TRADEOFF “REALISTIC” VS “PAIN TO GENERATE”
    - “REALISTIC” IS VERY VERY HARD:  
CHATGPT: REALISTIC TEXT GENERATION .. FOR 1 COLUMN ...
    - COMPANIES SELL THIS SAAS
  - ANOTHER PIECE OF CODE TO MAINTAIN
  - YOU WILL NEED SOME PROD DATA (NOMENCLATURE TABLES)

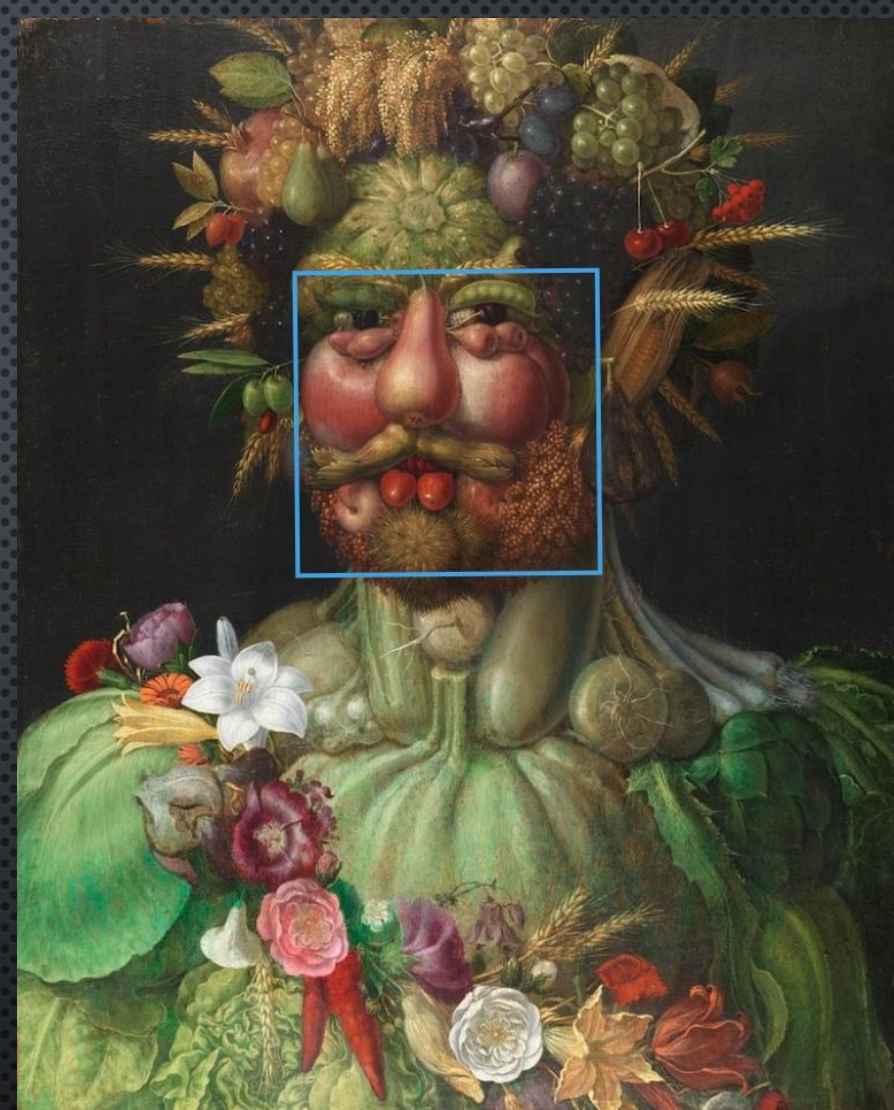
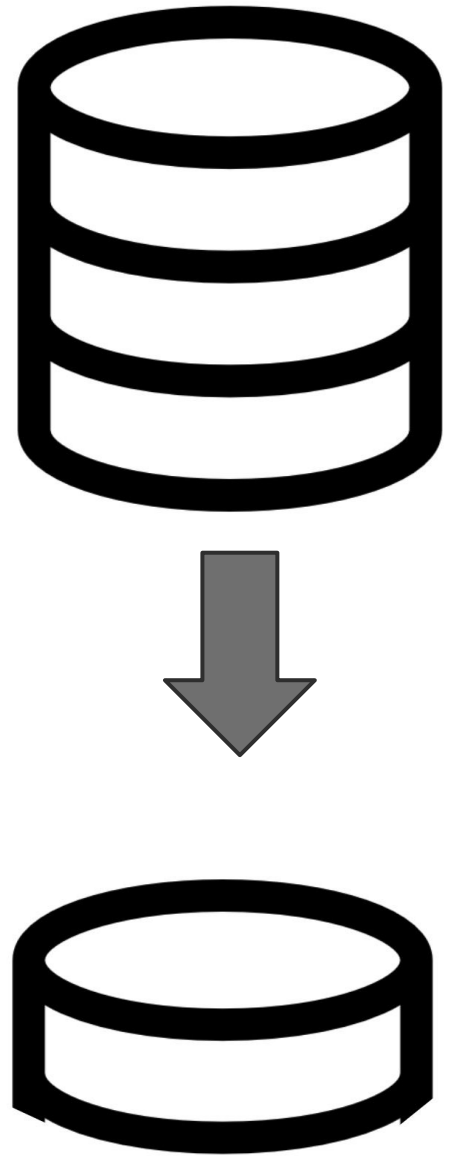


Image with face detection service



# PROD DATA

- PROD DATA: COPY DATA FROM PROD
  - STILL NEED TO FULFILL PK/ FK
    - ALL DATA DEPENDENCIES FULFILLED
    - RIGHT ORDER
- PRO
  - ULTIMATE “REALISTIC” DATA
  - DEBUG USING ACTUAL DATA TRAIL
- CONS
  - WHAT TO COPY?
- BIG SECURITY LIABILITY
  - MUST BE OBFUSCATED / MASKED / SCRAMBLED





# SYNTHETIC DATA / PROD DATA TAKEAWAY

- DATA FOR DEV TAKEAWAY



- SYNTHETIC

- GREAT TO MANUFACTURE TESTS
- “REALISTIC” IS HARD



- PROD

- AS “REAL” AS IT GETS
- MUST BE OBFUSCATED

- *HOW TO GET PROD DATA INTO DEV DB? □ 3 STRATEGIES*



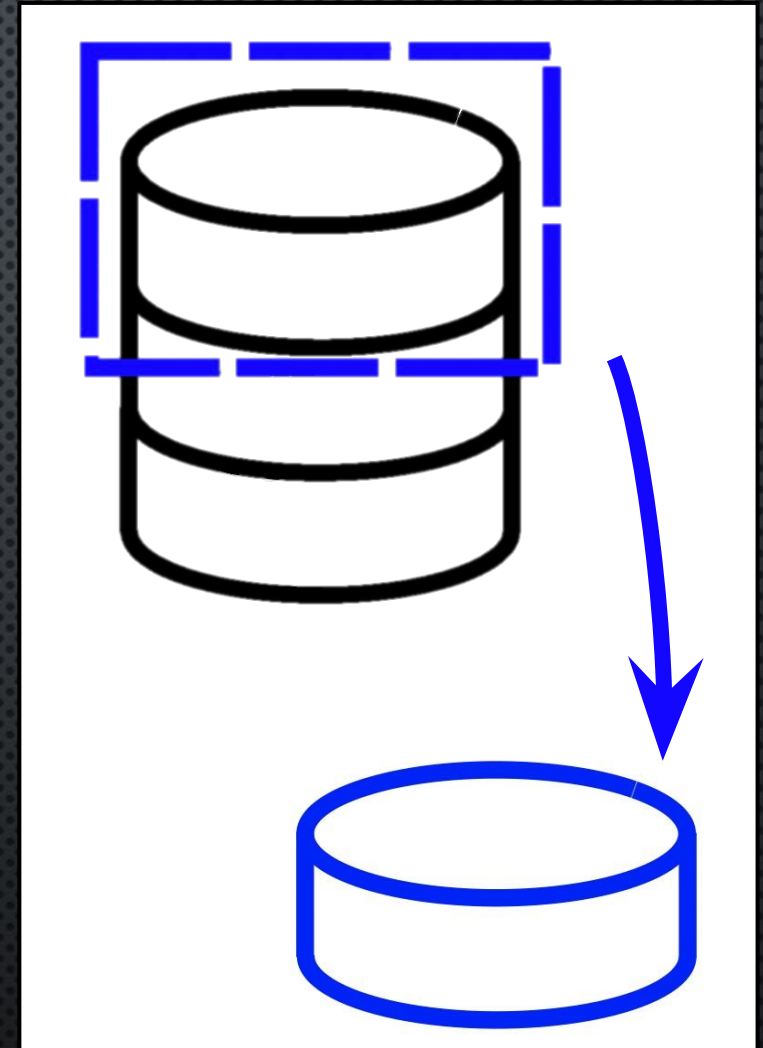
# GETTING PROD DATA INTO DEV DB

3 STRATEGIES



# PROD → DEV : PARTIAL INSERT

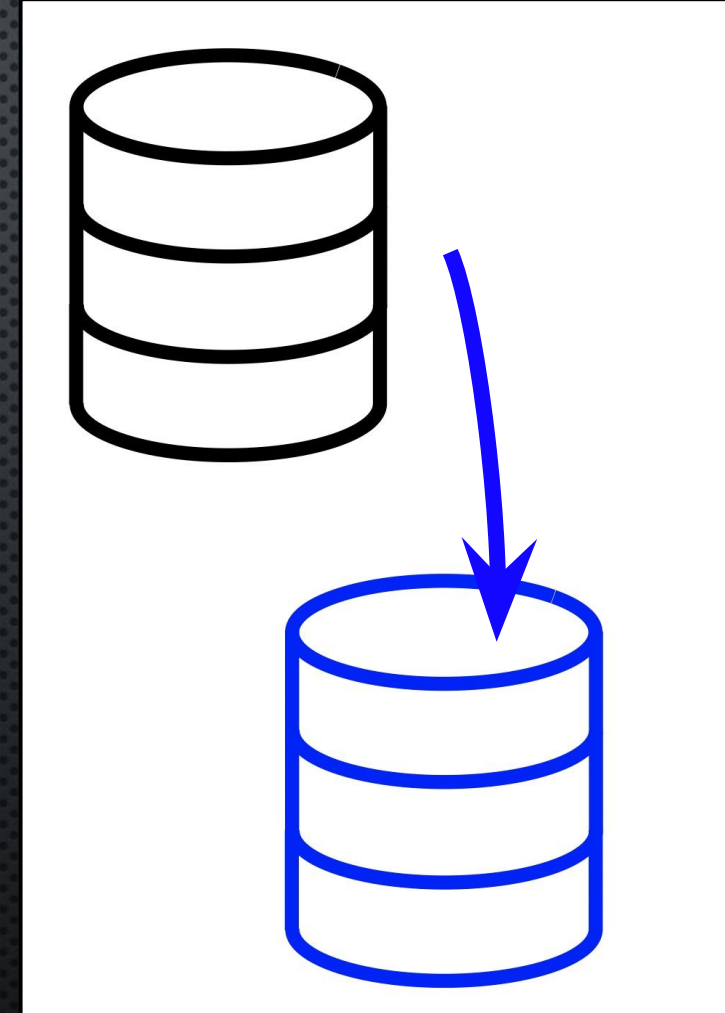
- PARTIAL INSERT
  - PROD: DUMP SCHEMA
  - DEV : CREATE SCHEMA
  - DEV: CREATE FDW
  - DEV: SQL: SELECT FROM FDW + INSERT
- PROs:
  - GET ONLY WHAT YOU NEED
    - FAST, CLEAN
    - UNLIMITED UNDO (DELETE ALL THEN RE-INSERT)
  - CAN COPY VERY FRESH DATA FOR DEBUG
- CONS:
  - HARD TO WRITE:
    - SATISFY FK (INSERT RIGHT THING IN RIGHT ORDER)
    - ALSO NEED STATIC DATA
  - HARD TO MAINTAIN:
    - CHANGE TO THE PROD SCHEMA → CHANGE PARTIAL COPY SCRIPT





# PROD → DEV : CLONE

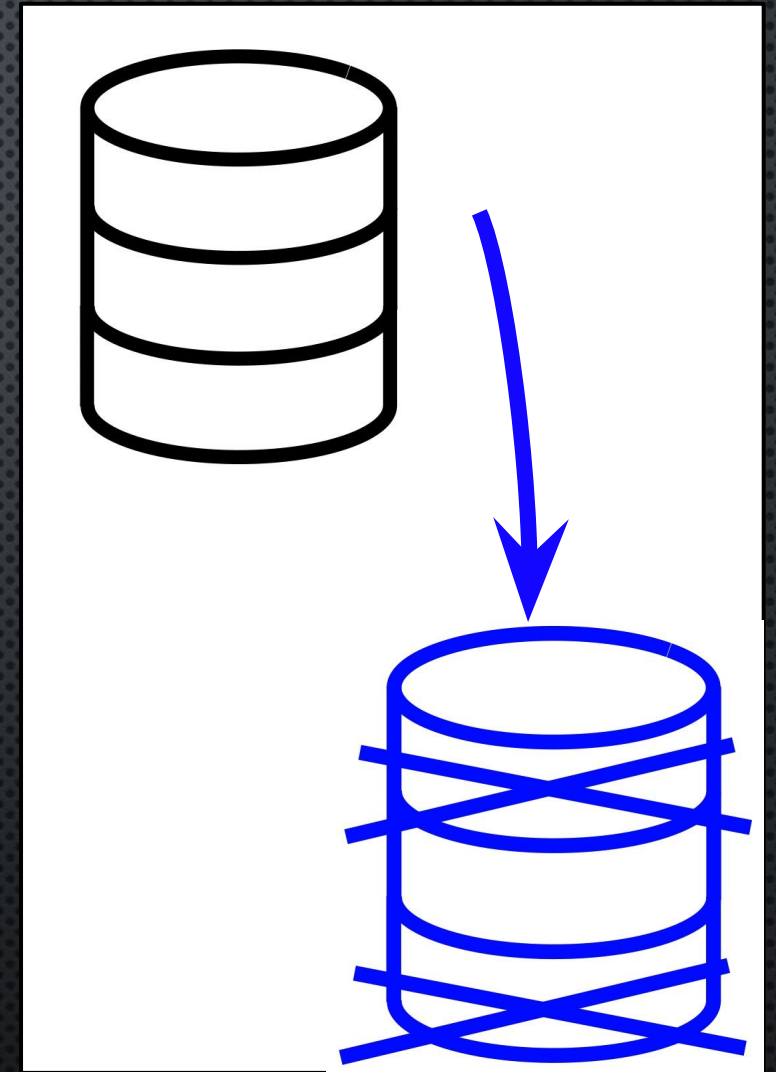
- CLONE
  - PROD: CLONE
  - DEV : INSTANTIATE
  - DEV : PARALLEL VACUUM
- PROs:
  - SIMPLEST TO MAINTAIN
  - ALL THE DATA
    - HIGHER CHANCES TO CATCH BUGS/MORE REALISTIC
- CONS:
  - ALL THE DATA
    - NEED BEEFIER INSTANCE, SLOW OPS
  - MORE DATA TO OBFUSCATE
  - LESS CONTROL (SECURITY)





# PROD → DEV : CLONE + DELETE

- CLONE THEN PARTIAL DELETE
  - PROD: CLONE
  - DEV : INSTANTIATE
  - DEV VACUUM
  - DEV: CHANGE ALL FK TO "ON DELETE CASCADE"
  - DEV: ADD INDEXES FOR FAST DELETION
  - DEV: PARALLEL DELETE
- PROs:
  - KEEP ONLY WHAT YOU NEED
  - SCRIPT TO DELETE VERY EASY TO MAINTAIN
- CONs:
  - INDEXING FOR FAST DELETION: LONG + A PAIN
  - DELETE CAN BE VERY SLOW IF LOTS OF DATA





# PROD → DEV : TAKEAWAY

- 3 STRATEGIES

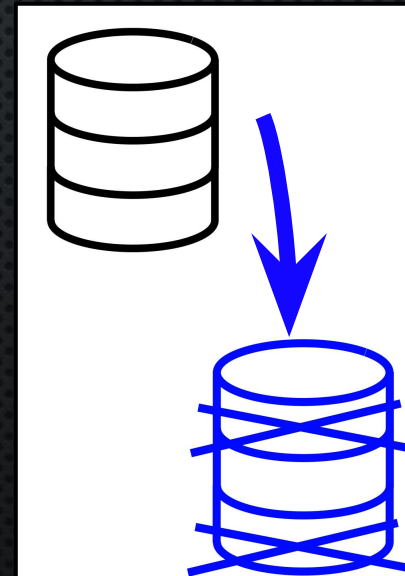
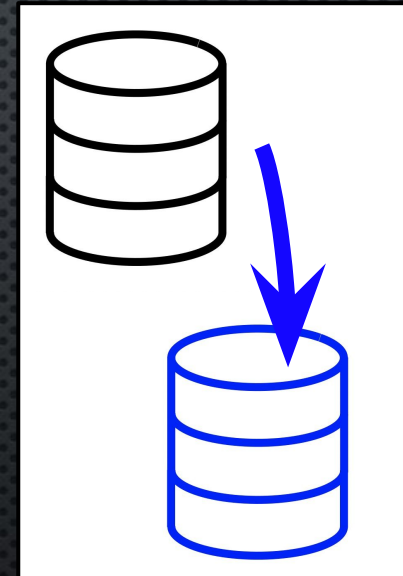
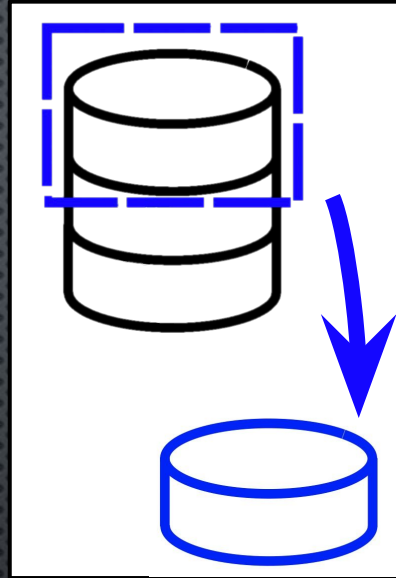


- PARTIAL INSERT
- FULL CLONE
- CLONE + DELETE

IN ALL CASES: PROD DATA IN DEV ?

- REGULATION (DATA BREACH IN WAITING)
- CONFIDENTIALITY (CUSTOMER NAMES, ...)
- PRIVACY (PASSWORD, EMAILS)

□ WE NEED TO OBFUSCATE IT!





OBFUSCATING PROD DATA



# WHAT IS OBFUSCATION

## OBFUSCATION:

- KEEP DATA USEFUL BUT DISGUISE IT
- MAINTAIN (CONSTRAINTS, STRUCTURE, FORMAT)
- 3 MAIN METHODS
  - KEEP ORIGINAL CONTENT, BUT HIDE IT
    - ENCRYPTION: (CAN BE UNDONE)
  - REPLACE ORIGINAL CONTENT, BUT CAN TRACE BACK
    - TOKENIZATION, (DETERMINISTIC) UUIDS
  - REPLACE ORIGINAL CONTENT, RANDOM
    - SCRAMBLING



risky

Random is annoying  
for test



# OBFUSCATION WITH UUIDS 1/3

## USE DETERMINISTIC UUIDS

- `UUID_GENERATE_v5()` : TEXT TO UUID
- `'I LOVE PGDAY.DE'`  $\square$  `'44F37065-35FF-5C70-8650-38A5BC931556'`
- `'44F37065-35FF-5C70-8650-38A5BC931556'`  $\square$  `'I LOVE PGDAY.DE'`
  - CAN'T BE UN-HASHED
  - USING PROD: CAN BE MAPPED (ANNOYING)
- IF INPUT IS UNIQUE, OUTPUT IS  $\sim$  GUARANTEED TO BE UNIQUE
  - OK FOR PK AND FK AND UNIQUENESS
  - INCLUDING COMPOSITE PK!
- BUT ... CHANGES FORMAT



# OBFUSCATION WITH UUIDS 2/3

HOW TO PRESERVE (SOME) OF THE ORIGINAL FORMAT / PASS CONSTRAINTS?

□ USE SEVERAL UUIDS !

- EMAIL:

- ILOVEPGDAY@CONF.DE □ UUID1@UUID2.DE

- FILE PATH

- /ILOVE/PGDAY/CONF.DE □ /UUID1/UUID2.DE

- AWS S3 CUSTOM TYPE

- (BUCKET, FILE\_PATH, REGION) □ (UUID1, UUID2, UUID3)

YOU GET THE IDEA



# OBFUSCATION WITH UUIDS 3/3

## SOME DETAILS:

- MAY HIDE MORE OR LESS
  - CUSTOMER EMAILS: HIDE EVERYTHING
  - PLAYERS EMAIL: KEEP DOMAIN IN CLEAR
- OBFUSCATE A PK ☐ NEED TO UPDATE ALL TABLES IN THE RELATION AT ONCE
- HEAVY DB WORK (NEED VACUUM ETC.)



# OBFUSCATION FOR REAL

- USING UUID IS ENOUGH FOR US
- YOU CAN DO MUCH BETTER AND MUCH MORE COMPLEX
  - SEE THE POSTGRES ANONYMIZATION EXTENSION (DALIBO)
  - (VENDORS OPTIONS AS WELL)



# WHAT TO OBFUSCATE? 1/2

WAIT ... WHAT SHOULD BE OBFUSCATED IN THE FIRST PLACE??

- REGULATION
  - DATA BELONGING TO CUSTOMER
  - PATIENT HEALTH INFORMATION
- PRIVACY
  - PASSWORD, API KEYS, ...
  - EMAILS
  - DISPLAY NAME, ETC.
- CONFIDENTIALITY
  - CUSTOMER NAMES !

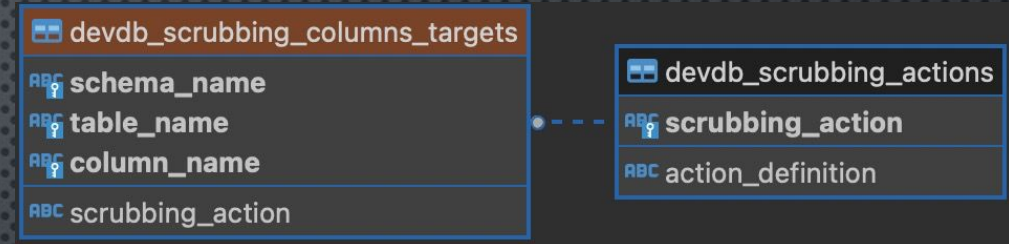
THAT'S A LOT! HOW TO FIND/KEEP TRACK?



# WHAT TO OBFUSCATE? 2/2

## FINDING COLUMNS WITH CUSTOMER NAME

1. LIST OF ALL TEXT COLUMNS
2. LIST OF ALL CUSTOMER NAMES
3. REGEXP MATCHES
4. MANUAL VALIDATION
5. STORE WHAT SHOULD BE DONE WHERE
6. (GENERATE CODE BLUEPRINT)
7. CORRECT CODE
8. PARALLELIZE, WRITE TESTS



<code>schema_name</code>	<code>table_name</code>	<code>column_name</code>	<code>scrubbing_action</code>
providers	customers	name	<a href="#">scramble_string</a>
providers	customers	email	<a href="#">scramble_email</a>



# SOME DEVOPS

*ALL CREDITS TO JON CORTEZ OUR SENIOR DEV OPS*

OUR DEVOPS STACK:

- JENKINS PIPELINES TO CREATE/DESTROY THE DEV DB
  - TERRAFORM TO BOOK INSTANCE + MANAGE PARAMETERS
    - BASH SCRIPT TO ORCHESTRATE
      - BASH SCRIPTS TO OBFUSCATE
        - SQL SCRIPTS / COMMANDS



TAKEAWAY MESSAGE



# TAKEAWAY

## *TESTIMONY:*

- A DEV DB HELPS
- USING PROD DATA COVERS MANY USAGES
- SAFER WITH OBFUSCATING
- ESPECIALLY RELEVANT WHEN POSTGRES IS MORE THAN STASHED DATA
  - POSTGIS ON MEDICAL IMAGES
  - SQL FUNCTION API (>100 FN)
  - ANALYTICS, CACHING, ...



# THANK YOU

*THANK YOU VERY MUCH TO THE COMMUNITY*

QUESTIONS