

Level 0 maturity : Manual process

- Manual build, train, lione, and deploy model { Jupyter Notebook }
-

Characteristics

- ① Manual | Script Driven | Interactive process.
- ② Disconnection between ML and operations team.
- ③ Infrequent release of models | no CI of models.
- ④ No CI. or CI is ignored.
- ⑤ Deployment \Rightarrow prediction service \times No ML pipeline deployed.
↓
No Cont. Training.
- ⑥ Lack of monitoring of performance.

Challenges :-

- ① Time consuming
- ② Maintenance cost \rightarrow high
- ③ No new ML ideas will pushed easily

- ④ Prone to bugs.

Possible solution:

- ① Actively monitor your model in production.
- ② frequently (depends on the situation) retrain your model.
- ③ frequent experiments with new optimized implementation to produce better models.

for ex: Change or experiment with latest SOTA

Make it more adaptable / Robust

Level 1 maturity : ML pipeline automation.

AIM

→ perform CI of the model in production → CI of model
by automating ML pipeline prediction service

→ To automate :-

→ new data \Rightarrow fresh model based on latest trend

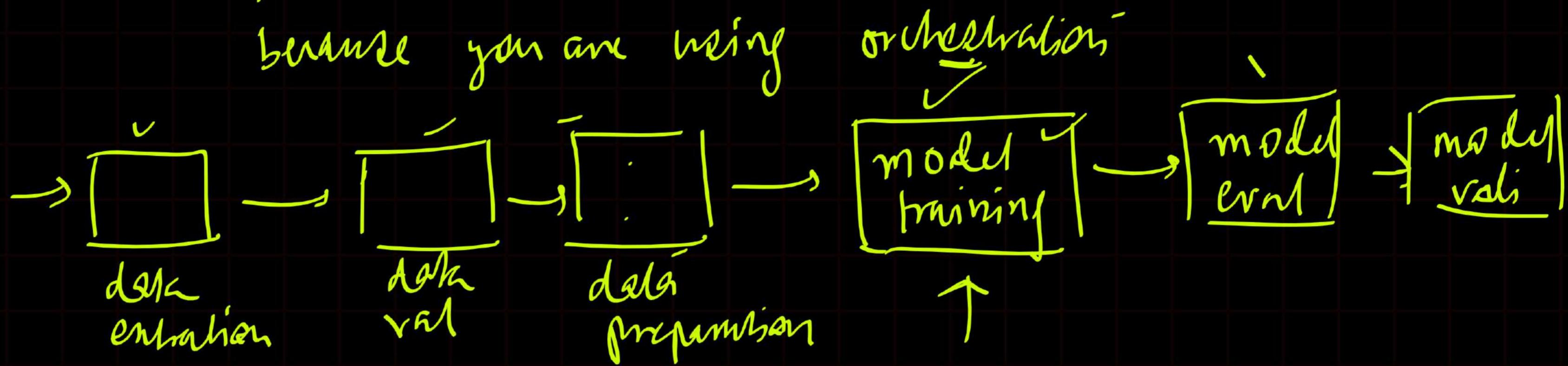
→ data & model validation

→ pipeline trigger & metadata management

→ Pipeline Orchestration

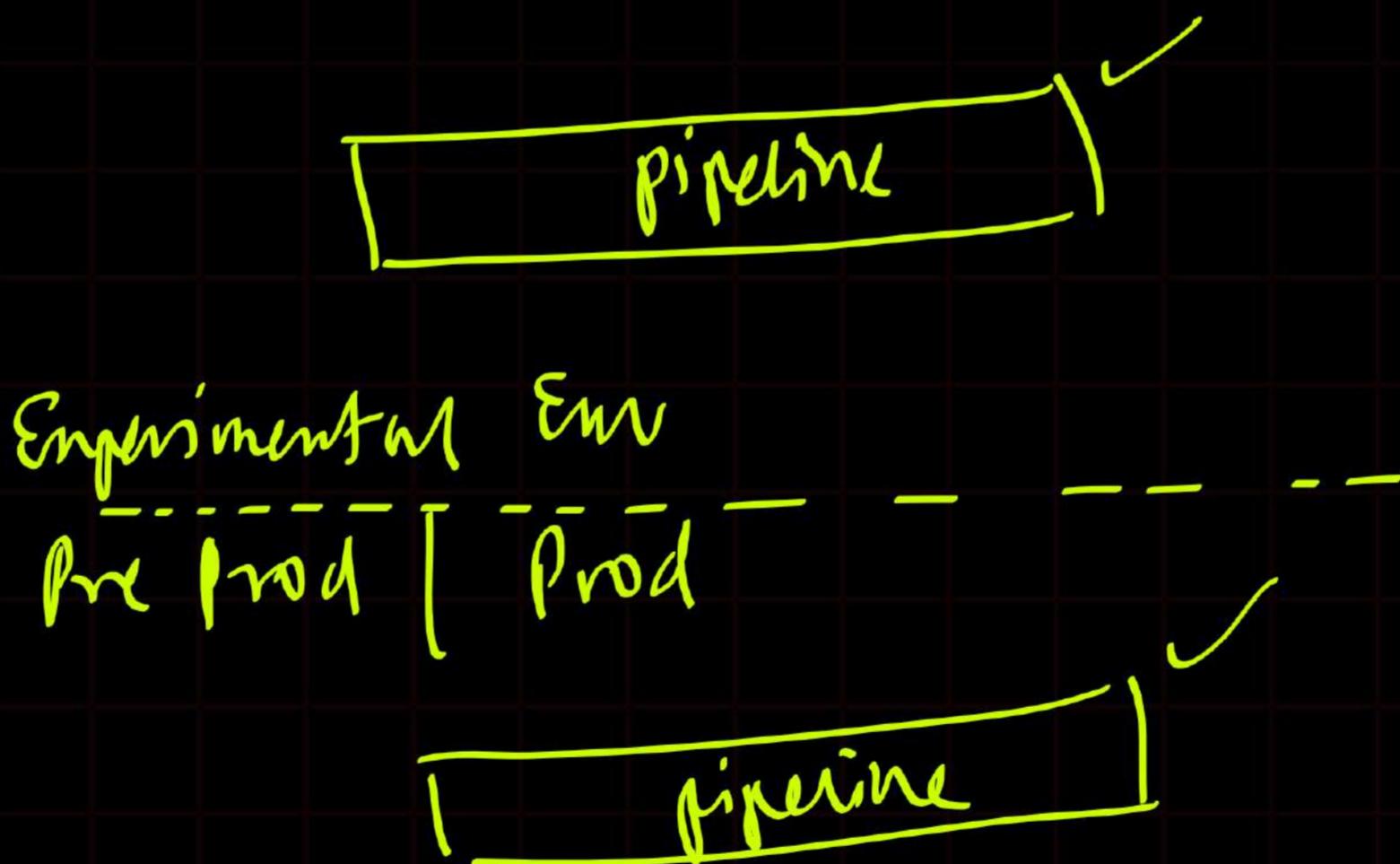
Characteristics :-

1. Rapid Experiment



2. CT of model in production

3. Experimental - Operational Symmetry



4. Modularized code components and pipeline:-

- Reusable
- Shareable

5. CI of model.

6. Pipeline deployment

deploy pipeline → deploy model

frequently used terms :-

① Pipeline Orchestration

ordered execution of components.

earlier failure of ML project :-

↳ Give code (code to connect two components of pipeline)

Tools to solve above issue :-

① Apache Beam

② " Airflow → UI

③ Kubernetes (Kubernetes based)

supported
by major
cloud platforms

Advantage :-

① Standard Orchestration
+
abstraction

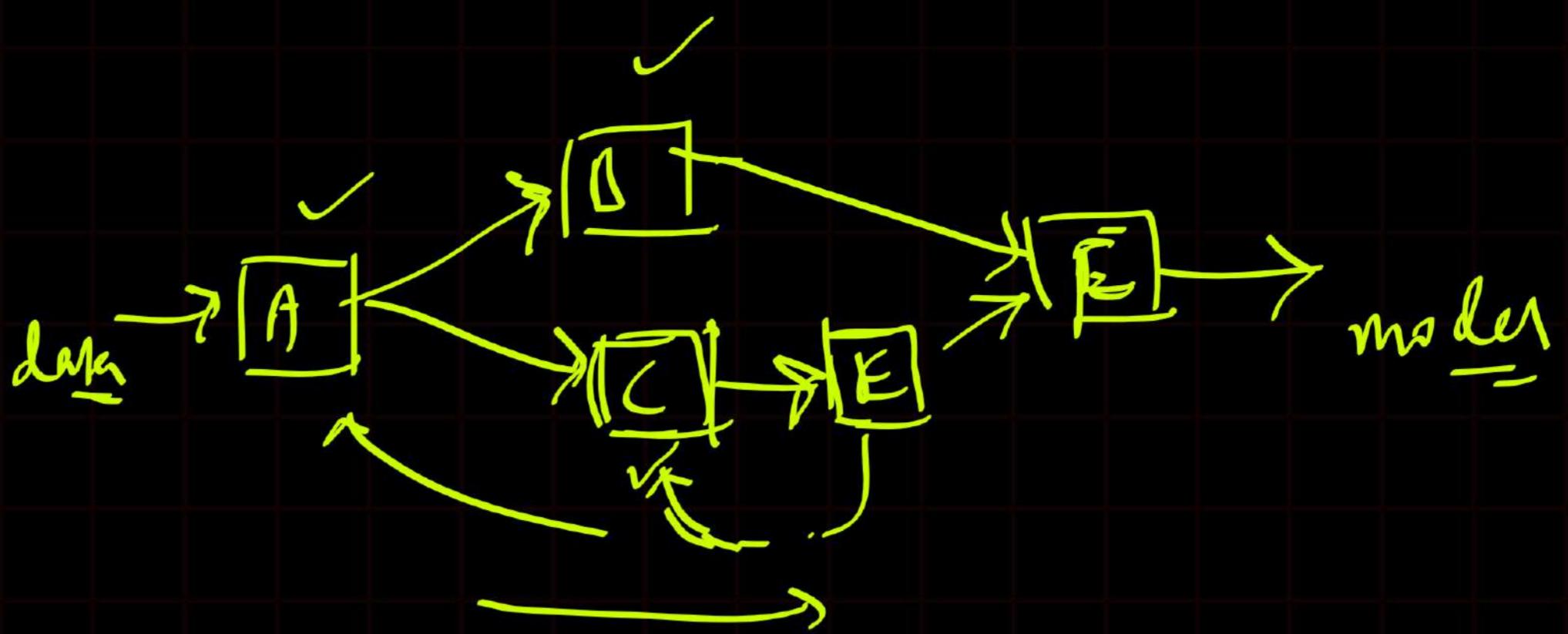
UI

② Supported by many cloud platforms

③ Easy to monitor / debug.

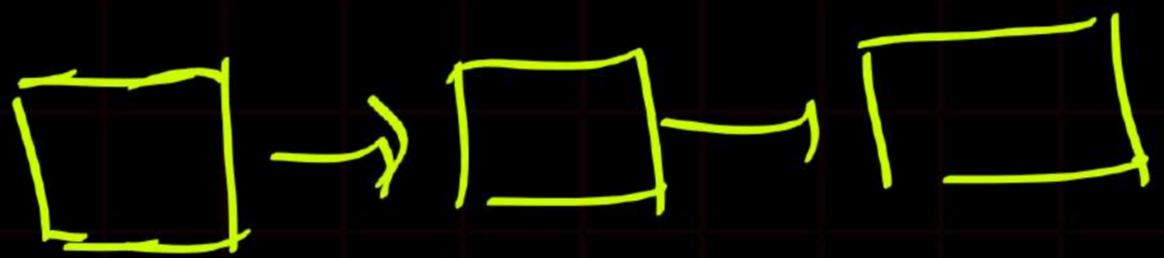
DAGs \rightarrow Directed Acyclic Graphs.

Graphs \rightarrow



cyclic $\rightarrow \infty$ loop

ML pipelines



DAGs \rightarrow one component of all your pipeline
orchestration tools.

Level 1. → important components.

pipeline evolution

new data / live data → new model
version on new data

Requirements to fulfill above evolution :-

① Automatic data validation

↳ Decides whether to start retraining

or

stop execution of pipeline & go for ^{manual} investigation
{ DS }

on following basis :-

(A) Data schema steps :-

↳ data is not in compliance with
expected schema or as per
DSA → Data sharing agreement

Solution :- → skip pipeline & let DS team
investigate.

some last stable model

or

maintenance fix out { Wony }
(as)

⑧ Data Value skew :-

When data pattern or statistical properties changes



Trigger retraining.

Requirement

② Model Validation

↳ After successful model training on new data -



Evaluate & validate your model before putting it into production.

Offline steps
before putting
your model

Online steps
model is putting
after deployment

offline steps :-

- Ⓐ Get evaluation metrics using your trained model on Test data set.
↳ know its prediction quality.
- Ⓑ Compare above metrics with current model in production
- Ⓒ Performance of your model must be consistent on regional / cluster samples

④ Test the model deployment for infrastructure compatibility

pickle, h5, pt, pt

REST API
gRPC

inferring
n-infering

Online model validation { when your model is serving in production }

⑤ Canary deployment

↳ roll out for subset of users or servers

e.g. (what)

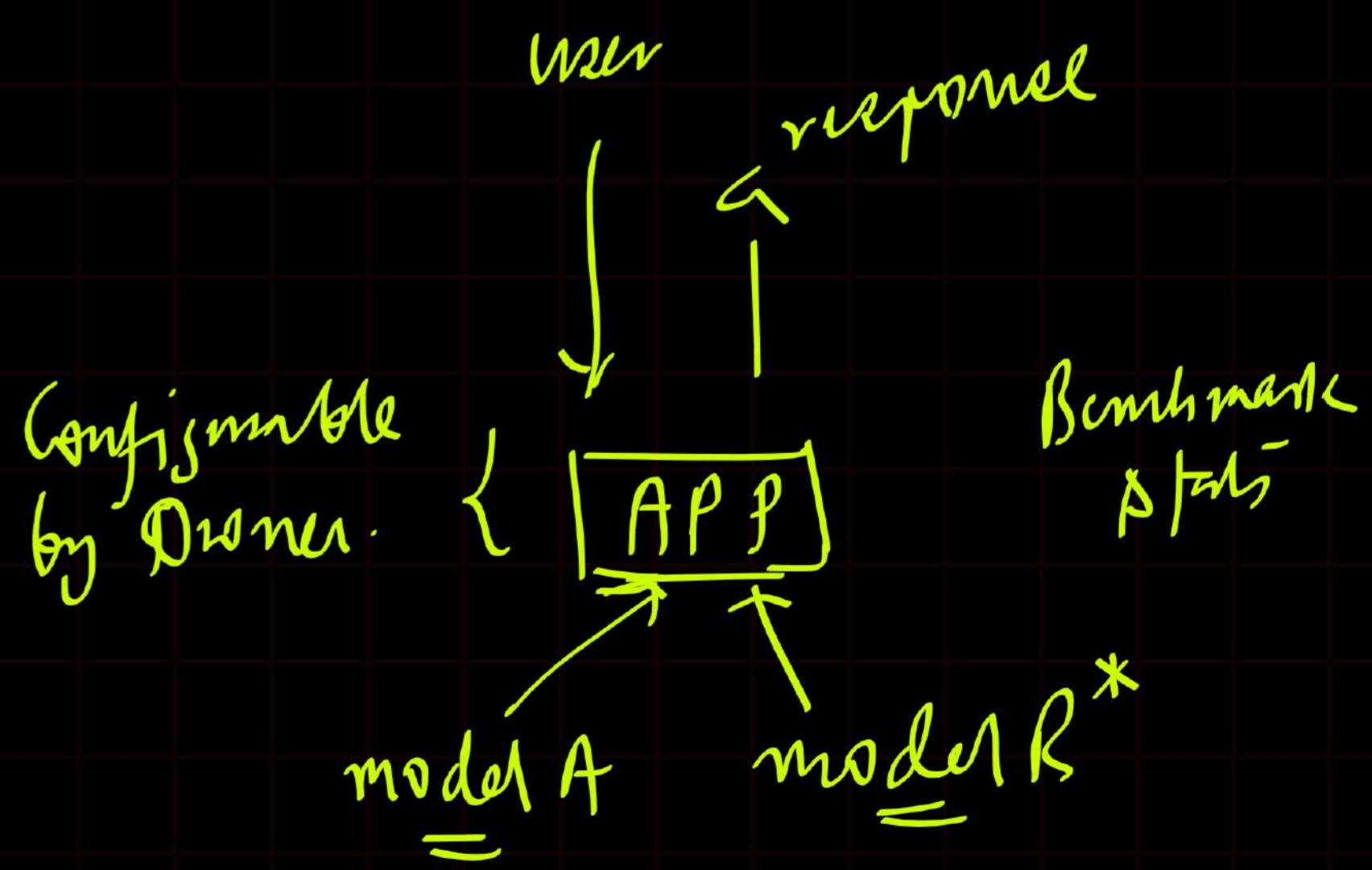
Tesla autopilot

HAL → Tejas, LCH, Dhruv

⑥ A/B testing setup

Advantage: you get info of data on which model A & B performed well.

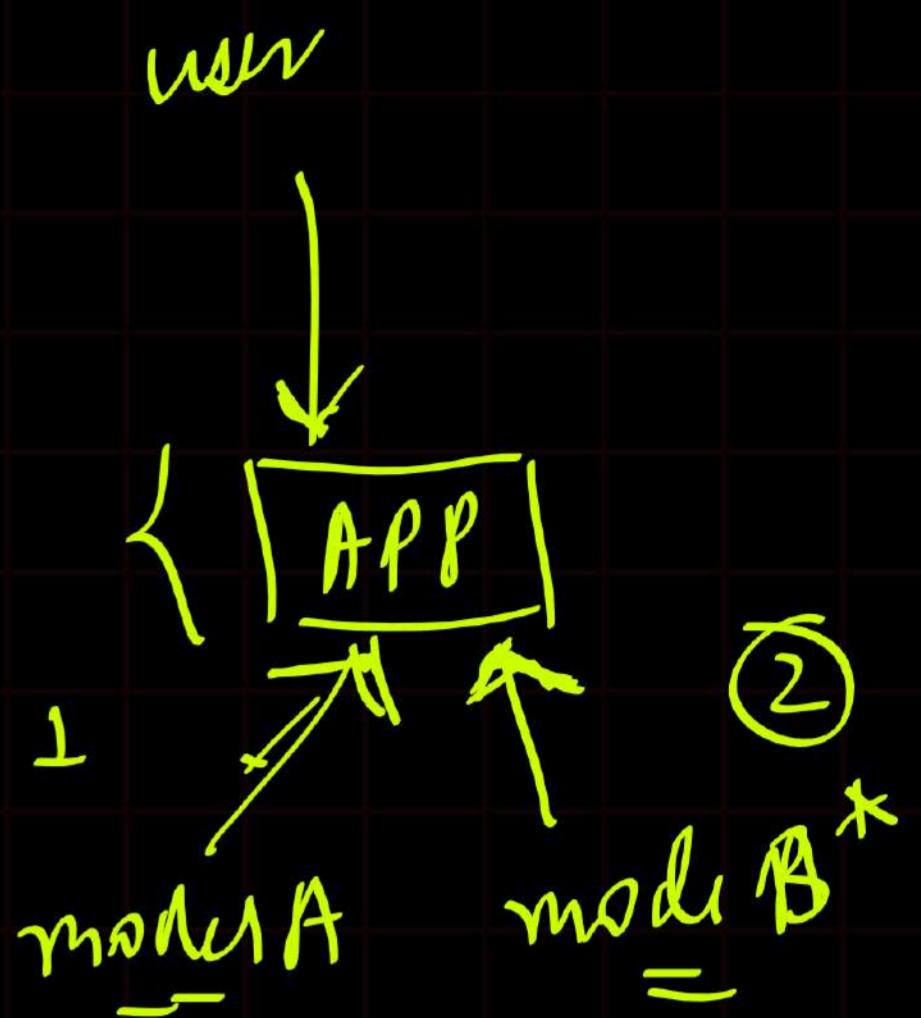
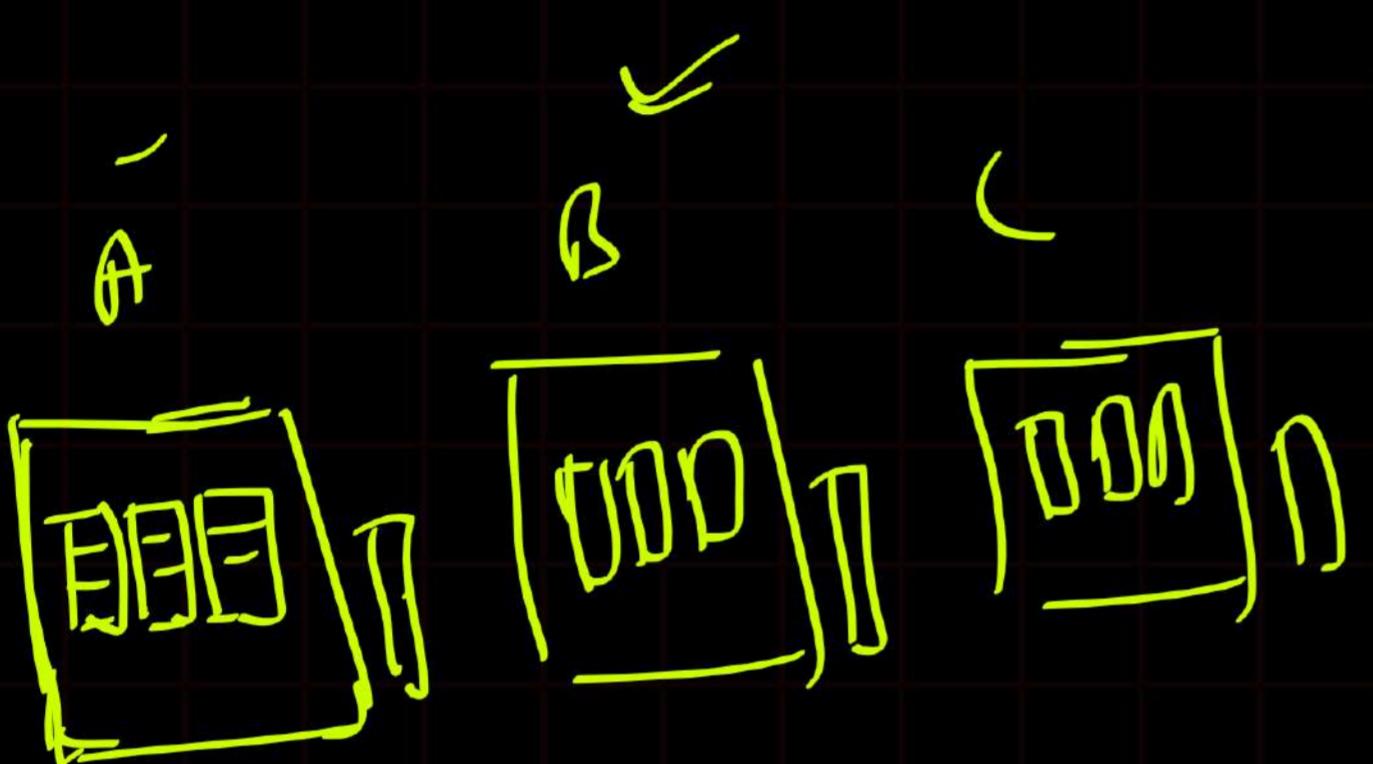
Disadvantage: Not a dynamic routing



③ Multi Arm Bandit deployment

Algorithm

Dynamic algorithm is
favour of better performing
variation
{ Dynamic routing }



Advantage: Dynamic routing

Disadvantage: You'll never get to
know kind of input
data stati which was
responsible for better
model performance

② feature store

XX Central Approach
↳ provider better control

↳ central depo
↳ for? → standard definition
storage & access of feature
for training and serving

feature
(x y z) → n y z
n y z

how? → provider API for
→ high throughput batch serving
→ low latency real time serving
for feature values.
→ supports training & serving workloads

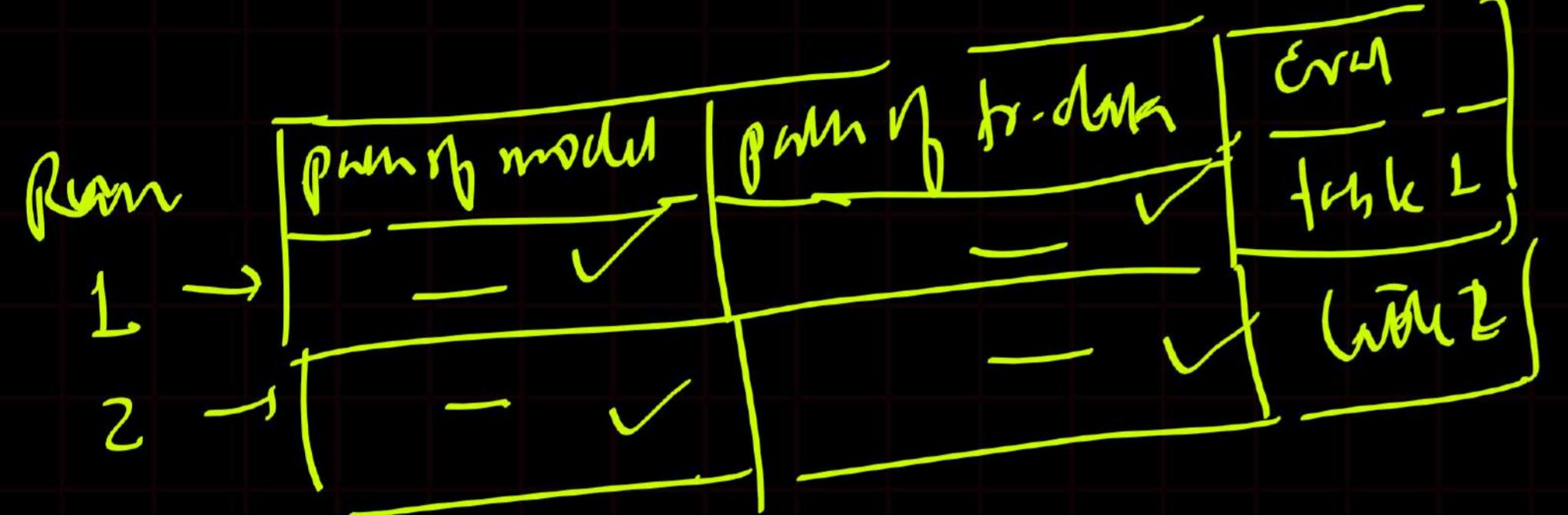
→ benefits? flaws?

- ↳ avoids similar features, different definition.
how it maintains feature & relevant metadata
- ↳ same up-to-date feature value
- ↳ avoids training-serving skew.
by using feature store as a central source
for experimentation / CT / online serving.

③ Metadata storage -

↳ ^{records} information about each execution of ML pipeline

why?



↳ To help in :-

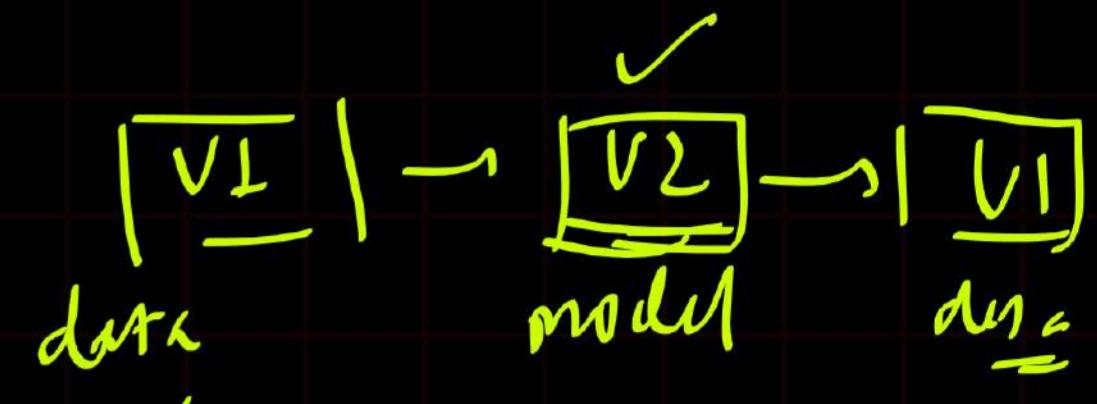
- ① Reproducibility
- ② Comparisons
- ③ Artifacts lineage
- ④ Debugging errors & anomalies

What if records in every step?

↳ (i) pipeline components & component versions executed.

e.g. dwc

md5 :



(ii) start | end date & time | Time duration for execution of each step.

(iii) Enricher (new) ✓

(iv) parameters passed to the pipeline

(v) pointers to the artifacts produced in each step

↳ path to prepared data
validation anomalies
computed stats

↳ keep in debugging
the failed step in
ML pipelines

(vi) pointer to prev. stable model in case rollback

(vii) Model eval metrics for every training & testing set

④ ML pipeline trigger.

↳ train on new data in production.

When you can trigger the training?

↳ On demand: Manual execution

↳ scheduled: if new data comes at regular intervals then scheduled trigger is fine.

↳ Availability of New training data:

↳ Model performance degradation: on noticeable change in performance drop.

↳ significant change in data distribution.

in experimentation → EDA → if anticipated → fix the computation steps

→ monitor

When level 1 is suitable?

- No freq deployment of new implementation of pipeline
- Manual testing of few pipeline is sufficient
- Manual deployment of new pipeline
- Submit fixed code base for pipeline to Dev team to deploy your ML pipeline into target environment
- New model based on new data { not on new ML idea }

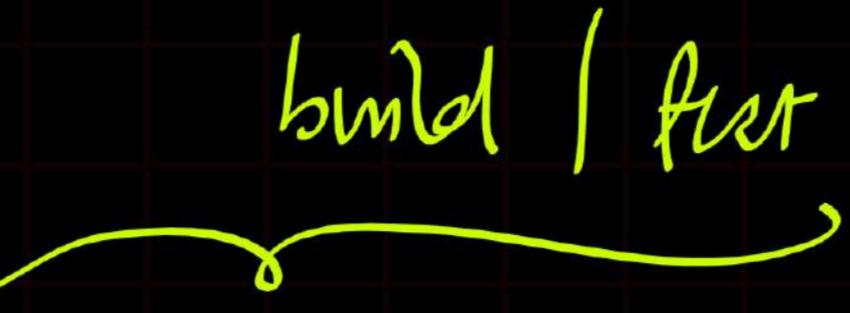
Challenge:

- Time consuming
- No freq updates

Solution:

↳ CI CD setup to automate
build | test | deployment
ML pipeline

Level 2 : CI/CD ML pipeline automation.

- ↳ rapid & reliable updates of ML pipeline
- ↳ automating  build / test / deployment of ML pipeline
 - ↳ beneficial for DS team \Rightarrow They can now focus new R & D. / new ideas
- ↳ New idea can be brought into production easily

CI/CD steps

- (i) Dev | Envt:
 - i/p: Raw data
 - process: experimenting with new ML ideas.
 - o/p: Assume code for ML pipeline steps.
- (ii) Pipeline CI:
 - i/p: prev o/p
 - process: testing,
 - o/p: Packages (containers) executable
- (iii) Pipeline CD:
 - i/p: prev o/p
 - process: Deployment packages into the target environment.
 - o/p: Deployed ML pipeline with new ML idea
- (iv) Automated triggering | CT:
 - i/p: prev o/p
 - process: trigger training in prod. based on schedule or
(other examples discussed)
 - o/p: Trained model → model registry.

(V) Model CD :

↳ i/p : prod o/p
perform : pick the suitable model from model reg.
↳ & integrate it with prediction service (prod).
o/p : pred. service.

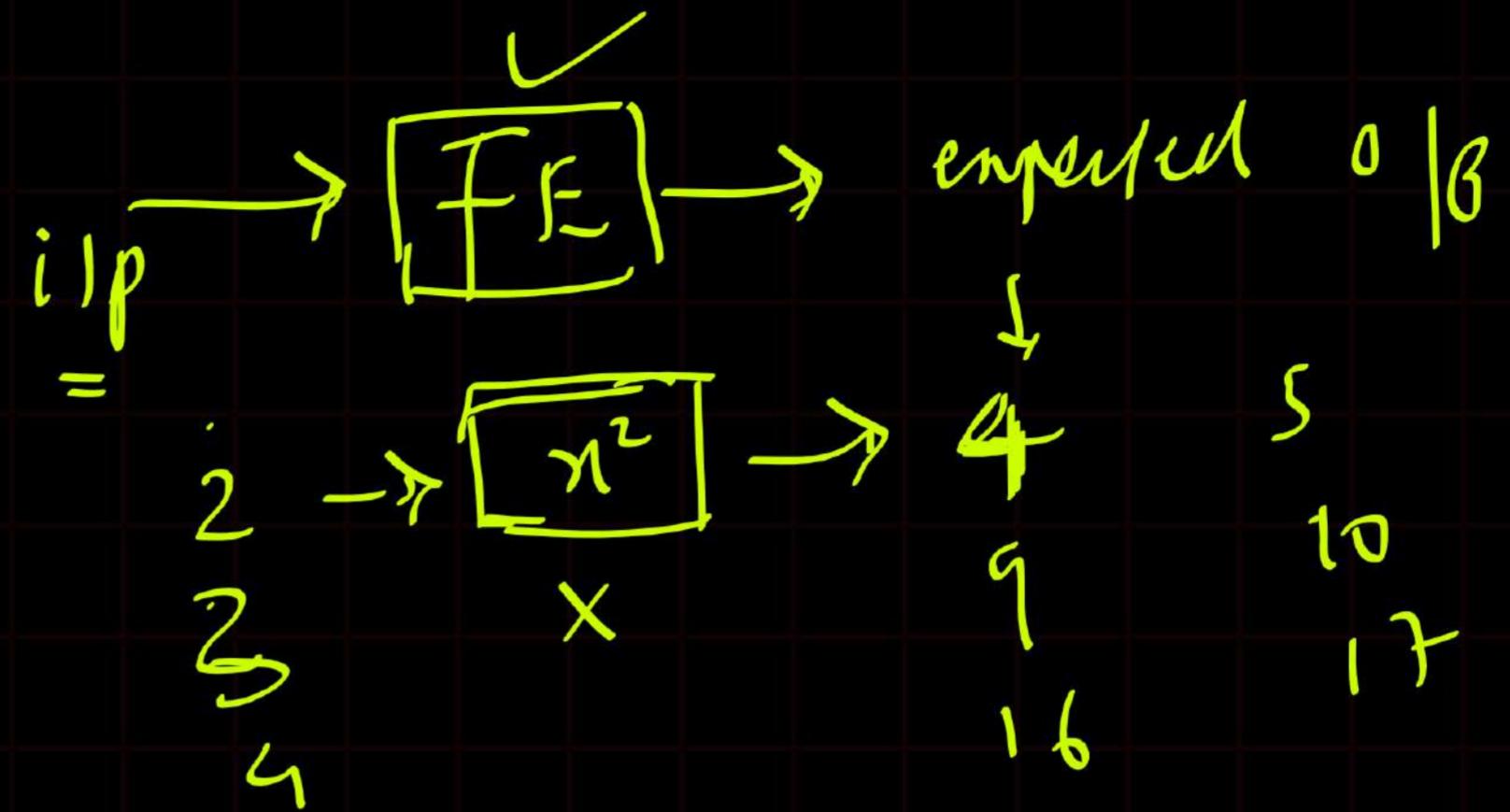
(vi) Monitoring :

perform : collect stats of the model performance on
new data

o/p ↗ trigger execution of pipeline in prod. { minor change }
 ↗ trigger the new experiment cycle { major change }

move on (I) ?

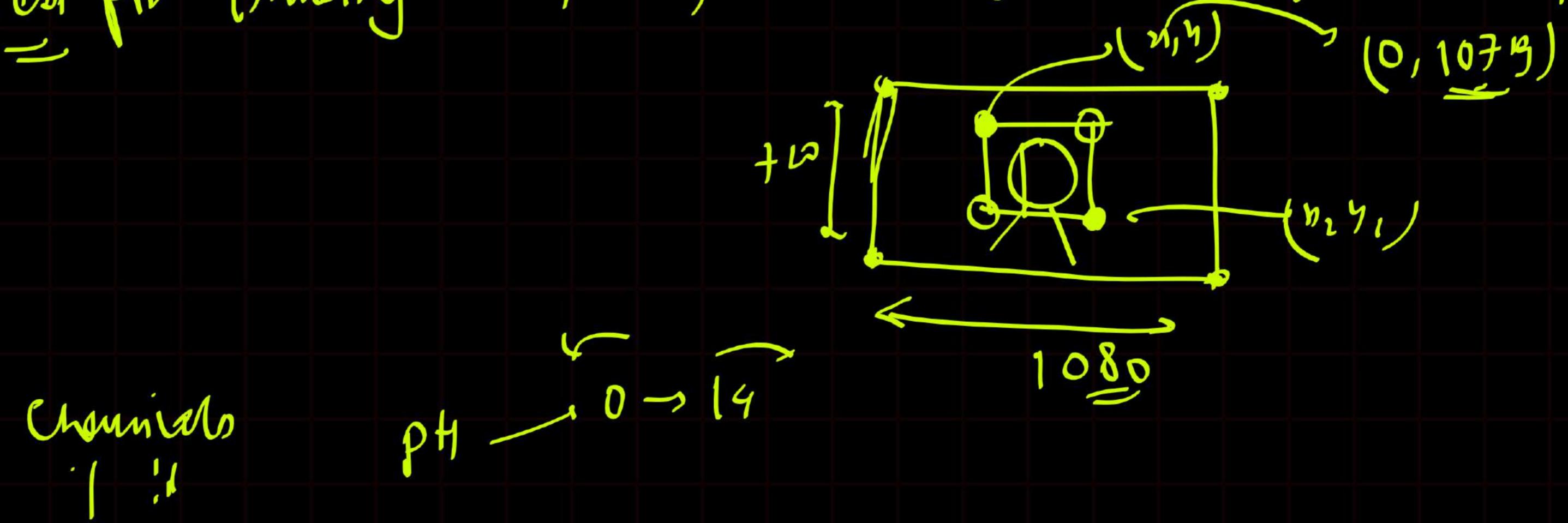
- ① unit test for feature engineering logic?



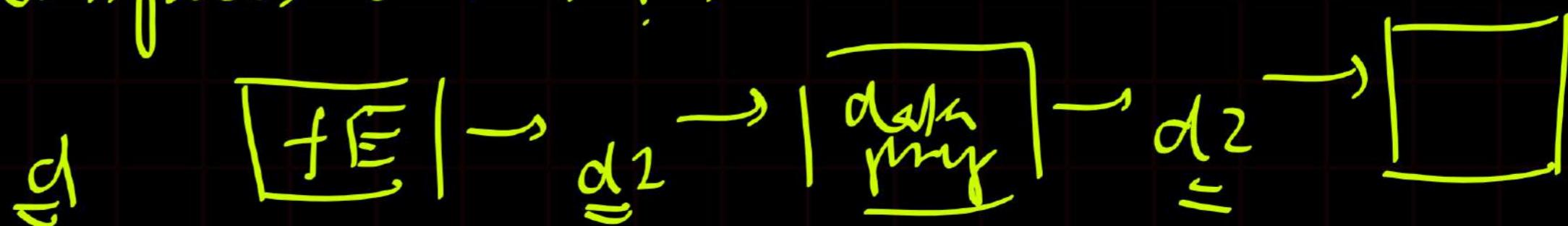
- ② unit test for methods implemented in source code.

- ③ Test for model training convergence.

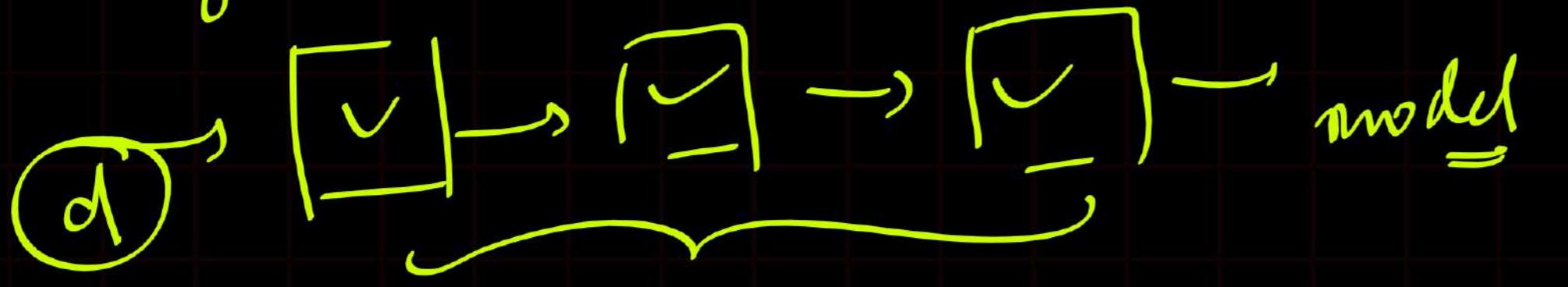
- * ④ test for training NaN, divide by zero, large no. $0/f$ from model.



- ⑤ Testing each component when it's generating the expected artifacts or not!



⑥ Interpretation tool -

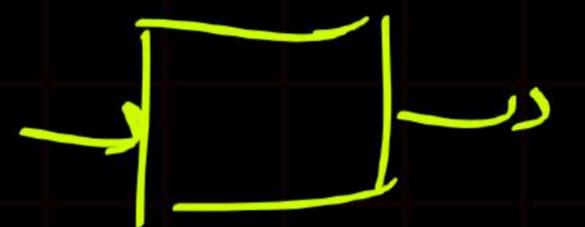


color images

hyper spectral Images

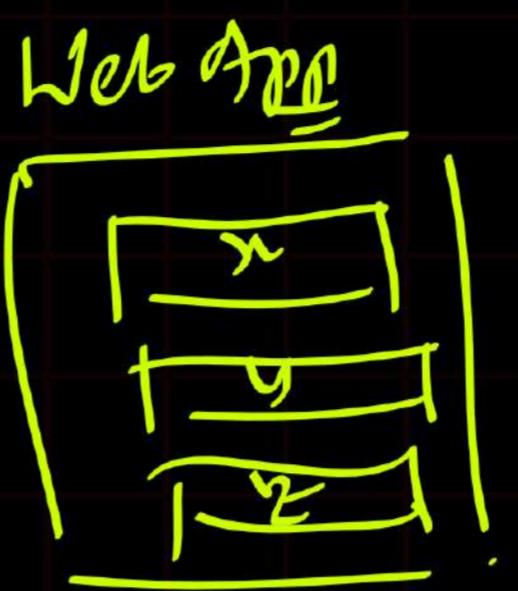
③

③+, mRE



How CI can be more smooth process?

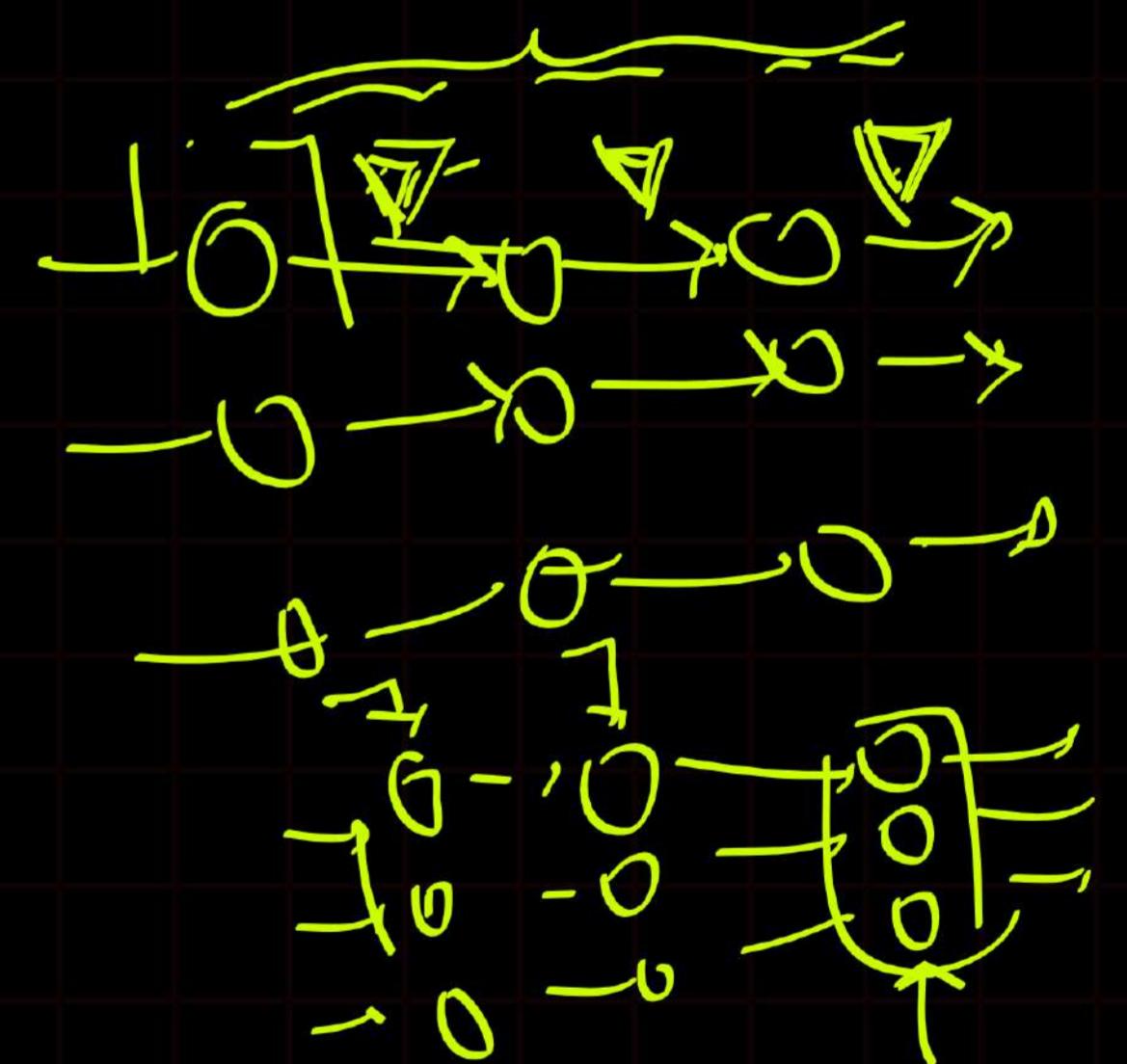
(1) Verify compatibility of model with target infra before deployment



\boxed{x} \equiv $x \rightarrow$ model
 \boxed{y} \equiv $y \rightarrow$ model

Placeholder, \rightarrow
memory!
compute
CPU execution TPU execution

500 MB \rightarrow 400 MB



(2) Testing of pred. service API

API \hookrightarrow JSON

(3) Test for pred. service performance.

↳ Load testing \rightarrow queries per second, model latency

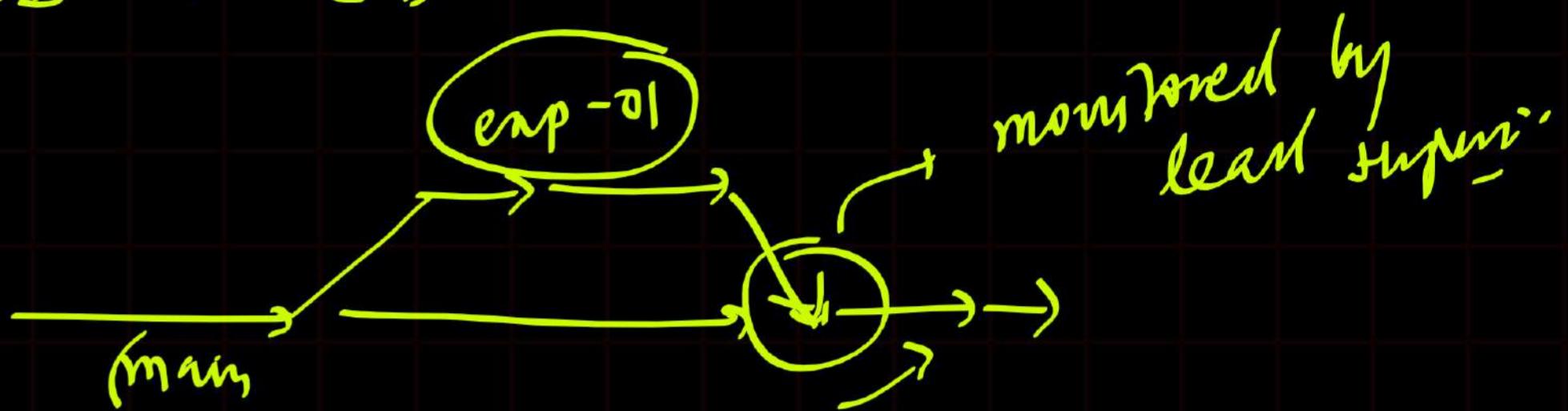
Deployment type:

Automated: You push the code in dev. branch of your GitHub repo. → CI → CD

semi-automated:

Manual:

after several successful iterations on pre prod.
deployment is done manually



Summary Level 2:

- ① Adaptive nature.
- ② less involvement of DS team in maintenance because of automation
- ③ less time consuming
- ④ More freq updates
- ⑤ Scalability.
- ⑥ Maintenance , Easy.

When you should use or go for AI(MLOps) ?

- ① Experimental phase POC ? X
- ② When you have deployed a webapp ~~test~~ + model
your user are growing ? ✓