O'REILLY®

Bookclub

# Designing Machine Learning Systems

An Iterative Process for Production-Ready Applications

Chip Huyen

* Read a chapter a week
* Livestream discussions every Saturday ▶ YouTube @MLOpsLearners
* Starts Saturday April 13th 2024!

**Chapter 2: Introduction to Machine Learning Systems Design**

**Group 6 (Pablo Discobar)**

Saturday, April 20, 2024

▶ YouTube

https://www.youtube.com/watch?v=q7RhT39LV8E

# Pablo Discobar

Ch2

Introduction to Machine
Learning Systems Design

Studious Pablo discobar with books

# INTRO



**Smit Zaveri**

Senior MLE @ Hitachi Vantara

Currently working on Monitoring data.

 SZ



**Gopichand**

MLR Intern @Samsung R&D

Working on Diffusion models

Let's connect : Linkedin😊

# INTRO

Raghunath Mahakud

SSE@PegaSystems.

Currently working with RAG based architecture using genai.

Linkedin

- Abhiramt_
- Aparna6227
- Divayjindal
- Shyam

# AGENDA

- Recap
- Key-Takeaways with examples
- Impact of GenAI - Personal Opinions
- Breakout Session

## Recap of Chapter 1

1. Success with ML by major companies like google made increased ML adoption across Tech industry.
2. More and more people moved to the ml space.
3. We saw how ML in Academia is very different in ML in industry.
4. How Tech Debt among data-science team can be a huge problem.
5. How ml-algorithm is just tip of the iceberg when coming to ML in production
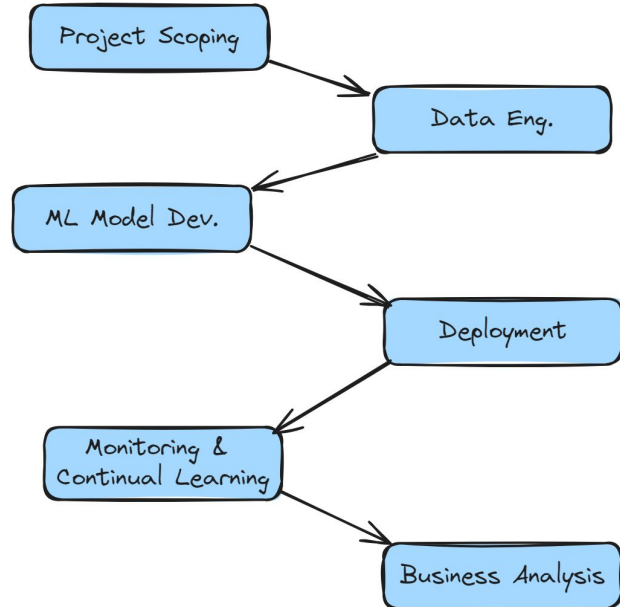
# AI has gone Mainstream !

# Misalignment and Challenges

1. Because of misalignments among Business and ML teams - ML projects failed
2. People lost jobs - Layoffs - ML/Data science teams
3. Next Obvious question!

What's the recipe to get Successful ML project ?

# Planning ML Project - Streamlining with Clarifications



**In a nutshell :**
1. Gather requirements and scope
2. Bring in and process data from multiple sources
3. Develop model guided by ML objectives
4. Make it available to user interaction
5. Monitor data shift and retrain
6. Capture changes in Business Objectives

# [Example] Problem

**Business problem**:      increase conversion (visitors -> new users)

**Business metric**:      new user signup rate

**Success criteria**:      new user signup rate goes up X%

**Assumption 1**:      if visitors like the items they see on the site, they'll sign up

**Solution**:      a recommender to show visitors items they'd like

# [Example] Problem

**Business problem**:   increase conversion (visitors -> new users)

**Business metric**:   new user signup rate

**Success criteria**:   new user signup rate goes up X%

**Assumption 1**:   if visitors like the items they see on the site, they'll sign up

**Solution**:   a recommender to show visitors items they'd like

**Assumption 2**:   like an item == click on it

**ML model**:   predicting how likely a user will click on an item

**ML metrics**:   precision@k, recall@k, logloss

# [Example] Problem - Final Output

- Users keep clicking on items recommended to them
- But the signup start rate doesn't go up



SOMEONE TELL ME, WHAT IS HAPPENING!

What to do when ML metrics are great, but business metrics aren't?

# Good ML metrics != good business metrics

What to do when ML metrics are great, but business metrics aren't? 😱

**Re-evaluate assumptions**

- **Assumption 2**: <u>click</u> on item == <u>like</u> an item
    - Use secondary metrics to get more insights into what's happening
        - Time spent on an item, items saved / liked / shared, bounce rate
    - Choose other signals
        - **Long CTR**: only count the clicks where users spend at least X secs     [YouTube]
        - **Purchase-through-rate (PTR)**: <u>Purchase</u> an item == <u>like</u> an item                [GrubHub]
        - **Add-to-bag (ATB)**: <u>Add to cart</u> == <u>like</u> an item

# Good ML metrics != good business metrics

What to do when ML metrics are great, but business metrics aren't? 😱

**Re-evaluate assumptions**

- **Assumption 2**: <u>click</u> on item == <u>like</u> an item
- **Assumption 1**: if visitors like the items they see on the site, they'll sign up
  - **Assumption 1a**: visitors don't sign up unless they have to
    - Solution: limited activities for visitors - sign up to do more     [Medium, Glassdoor]
  - **Assumption 1b**: the simpler the signup process, the more people sign up
    - Solution:  one-click sign up with Google, LinkedIn

13

# Good ML metrics != good business metrics

What to do when ML metrics are great, but business metrics aren't? 😱

**Re-evaluate assumptions**

ML only takes you 50% of the way there

# What does all this mean?

**ML solutions** must solve **business problems**.

**Success criteria** must be based on **business metrics**.

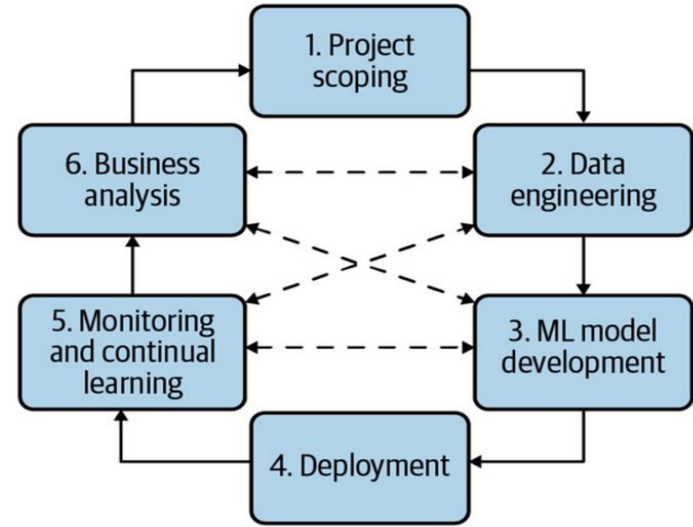It's essential to understand what **assumptions** we're making

**Choosing both Business metrics and ML metrics are HARD.**

# What does all this mean?

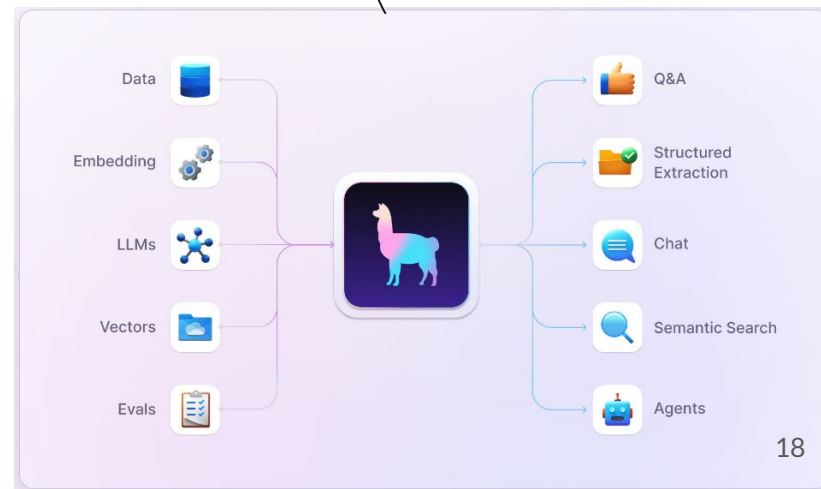ITERATIVE PROCESS

&

NON-LINEAR IN NATURE

# [EXAMPLE]

**Chat GPT fastest growing app with 100 million active users**

- OpenAI's ChatGPT has taken off in record-breaking fashion. According to OpenAI, even Chat GPT, the most talked about AI tool, surpassed 10 lakh users in mere five days, after its launch in November,2022.[8]

- If we compare it to its peers, then it took Instagram nearly 2.5 months to reach 1 million downloads and Netflix around 3.5 years to reach 1 million users.[9]

- By January 2023, ChatGPT hit 100 million active users, making it the fastest-growing application in history.

[Source](#)

17

# [EXAMPLE]

RAGs to riches



GPT-4

D When was your last knowledge update?

My last knowledge update was in April 2023. If there have been any significant developments or changes in the world since then, I might not be aware of them.

Is this conversation helpful so far?



Data

Embedding

LLMs

Vectors

Evals

Q&A

Structured Extraction

Chat

Semantic Search

Agents

[llamaindex](#)

# [EXAMPLE]



(a) Example jailbreak via competing objectives.

(b) Example jailbreak via mismatched generalization.

Figure 1: (a) GPT-4 refusing a prompt for harmful behavior, followed by a jailbreak attack leveraging competing objectives that elicits this behavior. (b) Claude v1.3 refusing the same prompt, followed by a jailbreak attack leveraging mismatched generalization (on Base64-encoded inputs).

Jailbroken: How Does LLM Safety Training Fail?

# What does all this mean?

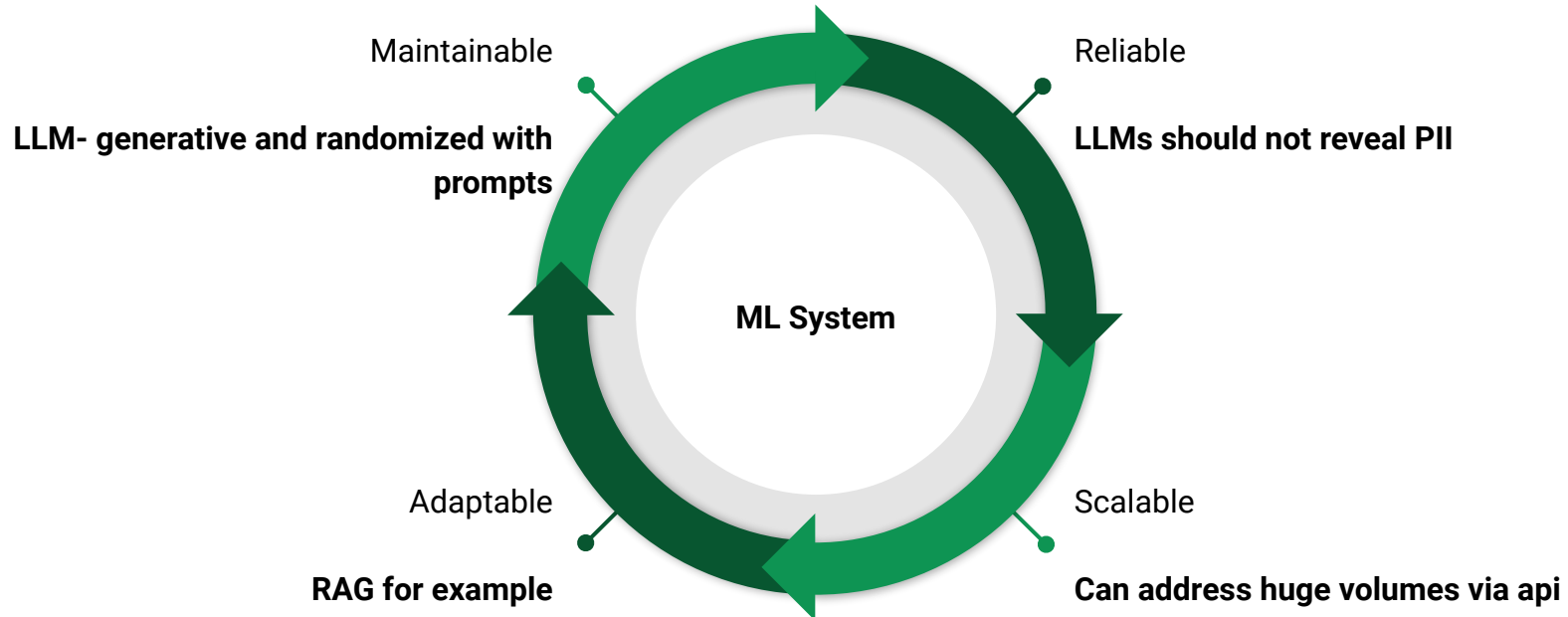**ML Systems Fails Silently !** even with LLMs

**ML models stale overtime !** How often they should be re-trained ?

- Continual Learning

How to test the new model ?

Safety Issues with LLMs

# What does all this mean?

Maintainable

**LLM- generative and randomized with prompts**

Reliable

**LLMs should not reveal PII**

ML System

Adaptable

**RAG for example**

Scalable

**Can address huge volumes via api**

21

# [EXAMPLE]

activity. The IC3 also strengthens the FBI's partnerships with our law enforcement and industry partners.

The 2016 Internet Crime Report highlights the IC3's efforts in monitoring trending scams such as Business Email Compromise (BEC), ransomware, tech support fraud, and extortion. In 2016, IC3 received a total of 298,728 complaints with reported losses in excess of $1.3 billion.

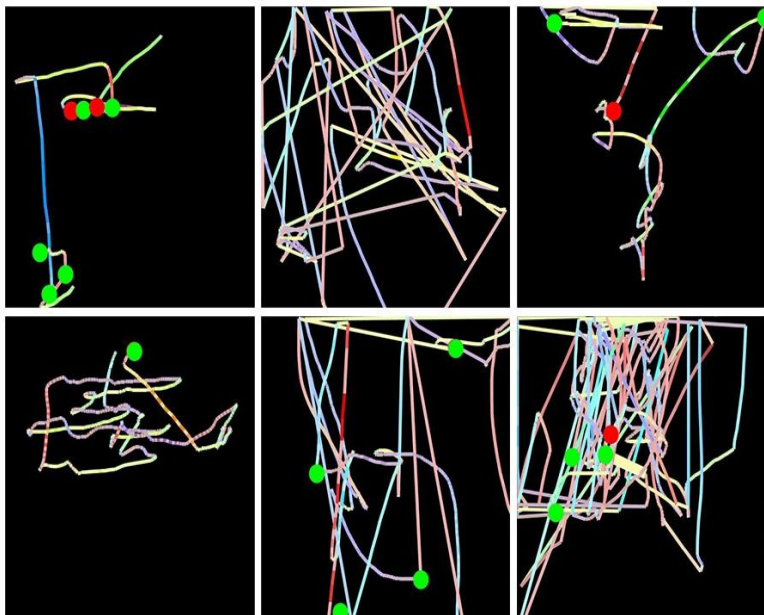This past year, the top three crime types reported by victims were non-payment and non-

Source

habits, education level, and familiarity with a service or system.

Habits and behaviors are very difficult to change and if we can identify legitimate users by their typical behavior patterns - we can detect anomalies on a totally new level. Same goes with fraudsters - ability to identify and quantify behavior patterns of cyber criminal will allow us to uncover and neutralize threats that may be undetectable by other means.

The question is: can we recognize the user - or a class of users by some unique ways they use their input devices, such as mouse or keyboard? Enter Behavioral Biometrics: the field of study related to measure of uniquely identifying and measurable patterns in human activities.

# [EXAMPLE]



| _time | _raw | X + Y coordinates of mouse movements data in Splunk |
|-------|------|-----------------------------------------------------|
| 2017-03-31 15:43:00.894 | {"site":"***","referer":"***","type":"mousemove","info":{"x":407,"y":903,"button":0,"time":1491000180894},"session":{" |
| 2017-03-31 15:43:00.892 | {"site":"***","referer":"***","type":"mousemove","info":{"x":427,"y":901,"button":0,"time":1491000180876},"session":{" |
| 2017-03-31 15:43:00.869 | {"site":"***","referer":"***","type":"mousemove","info":{"x":451,"y":895,"button":0,"time":1491000180869},"session":{" |
| 2017-03-31 15:43:00.861 | {"site":"***","referer":"***","type":"mousemove","info":{"x":479,"y":887,"button":0,"time":1491000180860},"session":{" |
| 2017-03-31 15:43:00.852 | {"site":"***","referer":"***","type":"mousemove","info":{"x":500,"y":883,"button":0,"time":1491000180852},"session":{" |
| 2017-03-31 15:43:00.844 | {"site":"***","referer":"***","type":"mousemove","info":{"x":524,"y":877,"button":0,"time":1491000180844},"session":{" |
| 2017-03-31 15:43:00.837 | {"site":"***","referer":"***","type":"mousemove","info":{"x":544,"y":868,"button":0,"time":1491000180836},"session":{" |
| 2017-03-31 15:43:00.828 | {"site":"***","referer":"***","type":"mousemove","info":{"x":562,"y":858,"button":0,"time":1491000180828},"session":{" |
| 2017-03-31 15:43:00.821 | {"site":"***","referer":"***","type":"mousemove","info":{"x":578,"y":855,"button":0,"time":1491000180820},"session":{" |
| 2017-03-31 15:43:00.812 | {"site":"***","referer":"***","type":"mousemove","info":{"x":588,"y":848,"button":0,"time":1491000180812},"session":{" |
| 2017-03-31 15:43:00.805 | {"site":"***","referer":"***","type":"mousemove","info":{"x":595,"y":843,"button":0,"time":1491000180804},"session":{" |
| 2017-03-31 15:43:00.796 | {"site":"***","referer":"***","type":"mousemove","info":{"x":602,"y":837,"button":0,"time":1491000180796},"session":{" |
| 2017-03-31 15:43:00.788 | {"site":"***","referer":"***","type":"mousemove","info":{"x":605,"y":832,"button":0,"time":1491000180788},"session":{" |
| 2017-03-31 15:43:00.783 | {"site":"***","referer":"***","type":"mousemove","info":{"x":608,"y":826,"button":0,"time":1491000180783},"session":{" |
| 2017-03-31 15:43:00.772 | {"site":"***","referer":"***","type":"mousemove","info":{"x":610,"y":822,"button":0,"time":1491000180772},"session":{" |
| 2017-03-31 15:43:00.766 | {"site":"***","referer":"***","type":"mousemove","info":{"x":613,"y":816,"button":0,"time":1491000180766},"session":{" |

# [EXAMPLE]

## Group classification

The first task was to prove that deep learning network can be trained to recognize mouse movements of two distinct groups of users: regular customers of financial information services business from non-customers while accessing similar pages.
Our guess is that patterns of behavior of people who are first time visitors are somewhat different from members who are generally more familiar with the portal.

It takes certain degree of learning for a "stranger" to understand the structure of a portal. And such learning curve comes with a mouse activity patterns that might be different from patterns exhibited by regular customers who are in general more efficient in finding and getting access to the information they need.

The architecture of neural network we've chosen roughly resembles successful VGG16 architecture for image recognition. Standard VGG16 architecture was further optimized for specifics of the dataset of non-natural images as well as for a limited size of our dataset.
With input of total 2,000 images for training and 800 images for validation (1000 + 400 for each class) - it took about 2 minutes for such neural network to be trained to achieve about 81% of validation accuracy:

**Unique way of determining User Behaviour using Mouse movements**

# What does this mean?

- How you frame the problem makes ML development easier or tough!
- Inter-portability of Regression ←→ Classification tasks
- Decoupling Objectives can save you from re-training of the Models
- Out of the Box thinking can give you unique and very effective solution
    - [https://www.splunk.com/en_us/blog/security/deep-learning-with-splunk-and-tensorflow-for-security-catching-the-fraudster-in-neural-networks-with-behavioral-biometrics.html](https://www.splunk.com/en_us/blog/security/deep-learning-with-splunk-and-tensorflow-for-security-catching-the-fraudster-in-neural-networks-with-behavioral-biometrics.html)

# POLL

# Impact of GenAI

- Model Development
  - API - more accessible , faster experimentation , more applications
  - Fine Tuning → Model Compression becomes very important
  - Huge compute and cost challenges and hardware and other issues -
    [Open Pretrained Transformers - Susan Zhang | Stanford MLSys #77](#)
- ML objectives -> evaluation became harder
- Business taking inverse route , to use GenAI coming up with use-cases to
  align with it rather than other way around

**Traditional Machine Learning (ML) Tasks:**

- ➔ Primarily used with structured data.
- ➔ Designed for specific tasks: Regression, Classification, Clustering.
- ➔ Requires feature engineering and explicit labels.
- ➔ Output: Numerical or categorical.
- ➔ Examples: House price prediction, email spam classification, customer segmentation, True casing

**Current Language Model (LLM) Tasks:**

➔ Primarily used with unstructured data (text).
➔ Designed for natural language processing tasks: Text generation, Translation, Summarization, Question-answering.
➔ Learns patterns from data without needing explicit labels.
➔ Output: Sequence of words or sentences.
➔ Examples: Generating human-like text, translating languages, summarizing documents, answering questions conversationally.
➔ Plan and execute with RAG for autonomous task.

# Thank YOU & Shout OUT!

# Breakout Session

**What I learnt with this book-club experience?**

**Discussion Topics**

- Any thoughts on Continual Learning for LLMs?
- How does Data Eng. changed with LLMs - especially with RAG, tokenization etc? - less feature engineering and more quality control.
- Guardrails for LLMs
- Challenges with deploying LLMs in-house and scaling them for internal Org?