

[Day 5] - "MLops for GenAI" (Unit 5)



Operationalizing
GenAI on your AI
Using MLops
Podcast

- Agent MLops - new frontier
- DevOps: making software dev smoother & faster
- MLops: same principle but applied to ML
 - challenges: data valid, eval. models, drift, reproducibility, etc..
 - Gen AI life cycle phases
 - Discover: no one size fits all model depends on use case & quality, latency
 - Dev & exp: iteration, refining data, eval metrics & going back to rework
 - Eval: can be automated once mature enough (ex AutoSS)
 - Deploy:
 - Govern:
 - Emphasis of tuning and adapting models
 - AI: max model guard - collection of all models.
 - Different from trad. prediction models:
 - multipurpose, not trained for a specific task & display emergent properties
 - input = everything, sensitive to prompt
 - Prompted model comp using prompt template.

Foundational
Model Paradigm

User input

Is model T400 available?

Prompt

You an...

U: can I?

O: Sure....

usr: is the modl?...

out:

Prompt Template

Instructions

You an...

examples

User: can I? Output: Sun...

User input w/
placeholder

usr: {{Question}}

placeholder

Output:

*fixed into foundational m+

Prompt is Data + Code : In Gen AI: also tweaking prompt itself, prompt engineering

Prompt | foundational model

Prompt model

Evaluate

↑
log as experiment & iterate

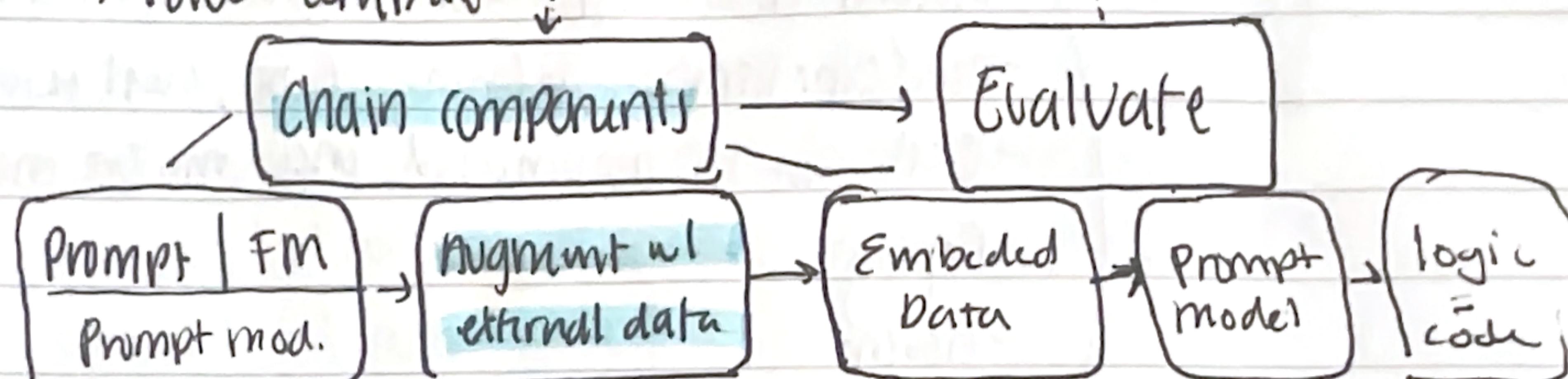
Prompt = Data

Pwmp = Code

- fewshot ex, any knowledge base, prompt = data
- Instruction, template, needs user control etc
* track which versions of prompts work best for which versions of models for reproducibility.

Chain & Augment

- connecting multiple prompted model components together along w/ calls to external API, etc.
* some complex log or exp. + iterate



- Agents use LLM as brain.
- Nature of inputs is much harder to do upfront
 - sol: think about entire chain as one unit
 - eval & versioning = more compn
 - Grounding, major chunk → grit for chaining
- Open source framework
- Adapting to perform better on task / domain
 - supervised fine tuning
 - reinforcement learning from human feedback
 - track all artifacts

Continuous training & tuning

Pwmp: Data Pract.

- tune model periodically
- synthetic data
- Model quant. to manage costs
- Gen AI builds prototypes
 - challenge: wider ranges of input
- Hard to materialize

• essential to safeguard AI apps against Adversarial attack

- Gen AI: focus on operationalizing a compute sys. for specific use
- FM: centers around user accessibility
 - * continuously monitor for drift & skew

• Make easy trackable & traceable & accountability

• Grounding entire system. (data, cost, prompts)

• Now dealing w/ autonomous systems.

— require SDIA frameworks & trust: tool orchestration + registry

Agents (OP)

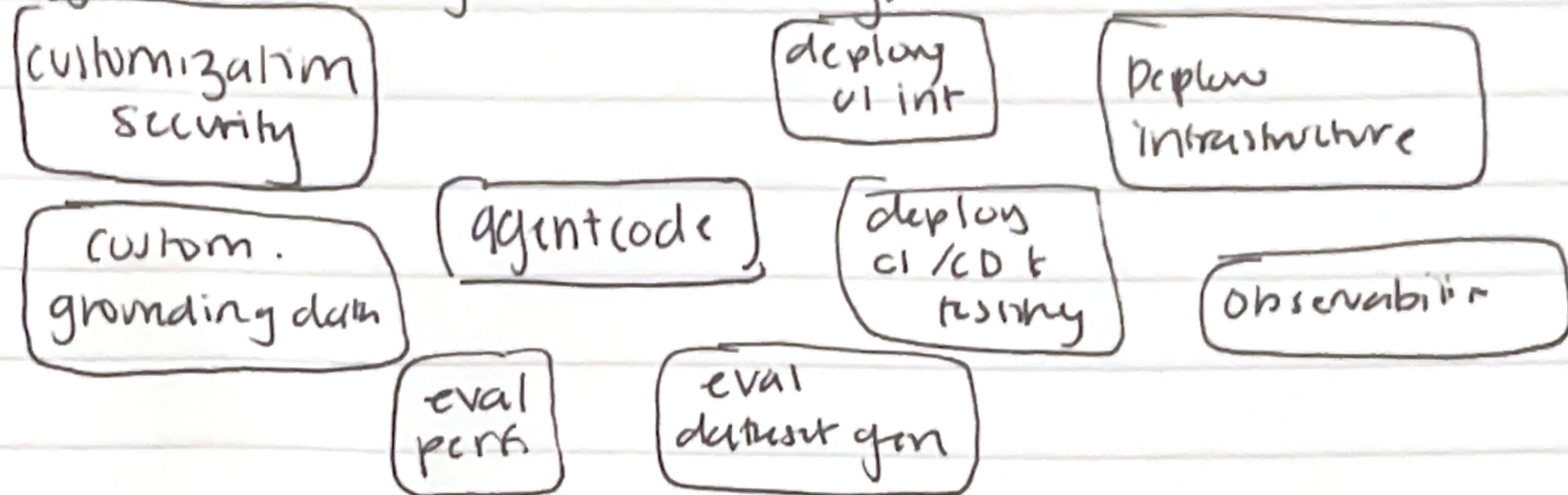
CHIIKAWA



[Live Stream]

Presentation

- Why it takes long to dev. an agent?



Repo=Quick

~kaggle agent-starter-pack

Star

Pop Quiz

- not a core practice of MLOPs? → C. training model from scratch
- Prompt temp. in context of GenAI? → Set of instructions (B)
- Purge & chaining? → B. avoid hallucination & maintain history
- Eval = crucial step in dev of GenAI → B. measure quality better.
- which AI offers numnrt effect at eval. jobs in prod., skew, and drift? → B. Pipelines.

[Caption]

Goal: Create notebook demo. use case using GenAI capabilities: Agents, vector databases, embeddings, or all above.

+ Bonus: create blog./YT vid.