

Common Pitfalls when building

GenAI Applications - Chip Hogen Date: 1/1

- ① Use GenAI when it is not needed
(If all you have is a hammer, everything looks like a nail)
- ② Confuse "Bad Product" w/ "Bad AI" (UX is hard)
- ③ Start too Complex (Shiny object Syndrome,
Be Vigilant when adopting
any abstractions)
- ④ Over-index on Early Success
 - Easy to build a demo, hard to build a product
- ⑤ Forgo human evaluation: Supplement LLM-as-a-judge w/ human eval. Improve System over time

Common challenges

- Accuracy / latency tradeoff
- Hard for agents to differentiate b/w ^{similar} freely
- Hard to completely understand customers intent
- Hard to create Unit tests
- Reliability from PPI providers (finedicty)
- Safety & Compliance

(6) CrowdSource Use Cases for GenAI

- Individuals might be biased toward the problems that immediately affect their day to day work instead of problems that might bring the highest ROI
- Without an overall strategy that considers the big picture, it is easy to get sidetracked into a series of small, low impact apps and then come to the wrong conclusion that genAI has no ROI

12 step checklist

- ① Cost
- ② Modality: text, image, Audio, Video
- ③ Customisation options
- ④ Inference options: realtime, batch etc
- ⑤ Latency
- ⑥ Architecture
- ⑦ Performance benchmarks : leaderboards
- ⑧ Language
- ⑨ Size and Complexity
- ⑩ Ability to Scale
- ⑪ Compliance
- ⑫ Environment

Building a Gen AI Platform

- Chip Huyen

Date: 1/1

Jan 18 2025

D. Model API : Query → LLM → Response

Step 1. Entrance Context

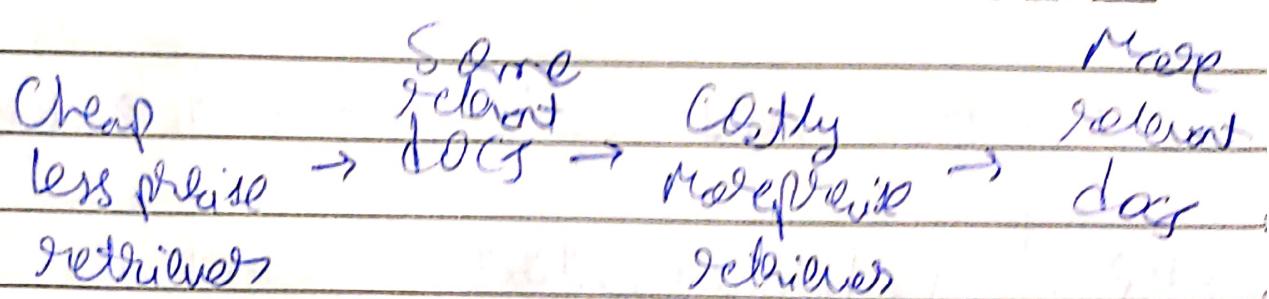
- The more relevant → less the model relies on internal knowledge → Better Ans info in Context ↗ Hallucination
- Retrieval
- RAG - Generation
- LLM augmented
- chunking : Splitting docs to manageable chunks.
To decrease overall context sent to LLM

Retrieval approaches :

- ① Term based : ex: Keyword Search
(Sparse vectors) BM25 (improves on TF-IDF)
- ② Embedding based : data → emb model → emb
(Vector Search) Query embedded, compared w/ emb data
- ③ Hybrid Search : ① + ②
 - Useful for exact words, names, acronyms etc.

Date: ___/___/___

Semantic
hybrid
Search



- Exact position of relevant chunk do not matter as much as they do in Traditional Search, as long as they are present

Ensemble: Combine different rankings together
HS to get a final ranking

RASH but better
Rankers
James Briggs)

- We lose some information if it is compressed to create a vector
- Say top-12 is 3, but one relevant chunk is in 17th position, that chunk will not be sent to the LM

and 3-16 are filling in the context window (limited)
irrelevant chunks

- So we want to retrieve all relevant chunks but not send all of them to LM → So Reranking

- Transformer is Computationally intensive.
Cosine Similarity is not
- Generating emb for docs & Storing them is doing the hard part first so that Cosine Similarity retrieval is Easy (Fast) at query time
- In Reranking, test is sent to LLM (docs, query) so info is uncompressed. So ranking is better but Costly.

Ex:

1000	$\xrightarrow{\text{retrieve}}$	25	$\xrightarrow{\text{rank}}$	3
total docs		retrieved docs		most relevant docs

- These retrieved chunks are sent to the LLM, to answer queries.

RAG w/ tabular data

- back to
chip's
article
- ① Text2SQL: SQL created from User's query + table schema
 - ② Run this SQL query
 - ③ Pass results + query to LLM, get response

Intermediate Step to determine which tables for query may be required if all table schemas don't fit in LLM's Buf

Date: ___ / ___ / ___

Agentic RAG

- A workflow that can incorporate external actions is called agentic
- term based retrieval
emb based retrieval
SQL extraction
Web Search
- Action that a model can take to augment its context

Query Rewriting

- User might ask ambiguous questions, especially in follow up, as they might have provided context in previous Q.
- So last Q should be rewritten to ensure it makes sense on its own
- Typically done using AI models:
 - "Given the following conversation, rewrite the last user intent to reflect what the user is actually asking"
- If identity resolution is required (How about his wife), then it gets complicated. #learnthesmarterway Rewrite model not dialogue

Step 2: Put in Guardrails

- Protect users & developers, place it where there is potential for failures.

Input
G/R → leaking private info to external APIs
Model Jailbreaking

If query contains sensitive info → Block Query
Mask PII, Send, Unmask

- Guardrails such that insert, update, delete queries cannot be generated automatically

- Define Out of Scope topics & predefined phrases to reject user query

Output
G/R → Output Quality Measurement →

- Empty responses
- Malformatted responses
- Toxic responses
- Hallucinations
- Sensitive info (trained or retrieved)
- Broad Risk & other bad signals

Petty while
balancing cost
& latency

human fallback
Splitter
Outrills
Angry
Sentiment

After n
turns #learnthesmarterway

Guardrail
tradeoff

① Increased costlier to latency
risk than

so better to have guardrails

② OP GR might not work in
stream completion mode

③ Self-hosting mode → data not sent to
3rd party

But then, cannot rely on guardrails
provided by 3rd party

Step 3

Add Model Router & Gateway

Specialized models: Better response, ↓ Cost

Intent classifications & Avoid out of scope queries

Router → next action predictors; ask for classification

Model gateway → Interface w/ models (open & closed) in several manners
Access Control & Cost management

Fallback Policies for rate limits / PPI failures
Load balancing, logging, Analytics

ex: Portkey, litellm

Step 4: Reduce Latency of Cache

- Prompt Cache (at inference API)
- Many prompts have overlapping text segments (ex System prompt).
 - This can be cached & reused so that it is processed just once (during first query).
 - Supported by OpenAI, Claude and Gemini (Called Context Caching)

- Exact Cache
- System stores processed items for reuse later when exact items are requested
 - Ex: Do not gen Product Summary if already exists
 - Avoid redundant vector search by checking if incoming query is in Vector Search Cache
 - Useful for queries having multiple steps or time consuming actions

- Semantic Cache
- Allows reuse of Similar query results
 - You should have reliable way to check if it is Similar
 - might not be useful as many of its components are prone to failure

Step 5: Add Complex logic & Write Actions

Complex logic

- Ability to Plan & decide what to do next
- Loops & Conditionals

work actions

- Can help automate Complex Plans
Ex: Customer outreach Workflow
- Prevent injection: Social Engineering on AI
- Ex: Leaking Sensitive data from db,
Corrupting data

Orchestration: Workflows & Pipelines

Observability

Metrics: model metrics: Accuracy, toxicity, hallucination
retrieval quality

System metrics: Uptime, mem, throughput etc

Develop metrics to monitor different ways
in which models can fail

Logs

Log everything

Process

Request's execution Path through Components