

- What is the project about? Tax Assistant / Tax 6-pilot for filing
 - Significance
 - Earnings call (Intuit Assist)
 - media coverage (Collaborative Prod. Sys. & Roadmap)
 - My Role
 - Tech Leadership: Arch + hands-on (Driver from AI), code-review, XFN + peers designation for key stakeholders
 - Mentoring & coaching (team size leading: 6 AS + 3 MLEs) (lead of leads)
- Small group of mins*

- Context
- Customer Problem
- Prior Soln (before GenAI era)
- Issues with Prior Soln
- Business Problem: Helpfulness rating 8/10 → 64% conversion ↓ 40M
- Success Criteria: How did you come up with it?

- 1. oneness*
- Tech Strategy in step-wise fashion
 - LLM selection: Tradeoffs: Cost vs Latency vs Accuracy (Hallucination) *→ explain*
 - Tech Levers: Tradeoffs i) Prompt ii) RAG iii) Fine-tuning
 - Architecture Diagram (whiteboard)
 - Eval as a challenge, calc as challenge *→ sign-post*
 - Success Criteria: S2C, Tale North: >95%, >90%, <1y.
 - Sign-Post: Helpfulness, Accuracy, Reliability, Coherence, Hallucination, latency <5sec, cost (Acceptable by LT)

- 2. soft values*
- Challenges: i) Eval → Auto-Eval ii) Calc → Tools
 - Learnings:
 - i) switching LLMs
 - ii) planner (Q&A questions)
 - iii) PEFT, LoRA, learning rate, system prompt, user prompt, globality & quantity, why not fine-tuning, Lora, learning rate, catastrophic forgetting & self
 - SFT Deep-Dive
 - i) Training Data (prompt frequency) *→ same tokenizer*
 - ii) Training Loss, Instability, eval & iteration
 - iii) Hyperparameters: batch size, learning rate multiplier, epochs - innovation (LLM-as-a-service)
 - Deployment & Inf. optimization (Paged Attention, KV Cache, 400G F100 GPU)

- 3. mins*
- Learnings:
 - Business Impact Delivered: 21M users, 0.25/s 10M
 - Next steps: Deployment with ft model

SFT Deep Dive

• Data

- quality over quantity
- quality < AWS Mech Task
- consistency across 3 experts
- quantity: few thousand examples

- Data cleaning: dedupe
 - Data validation:
 - format Prompt: Response
 (fewshot script) | Instructions
 (stamp) prompt | Content
 user prompt | User question
 - json

• Training:

- Full fine-tuning vs PEFT (LoRA, QLoRA)
- Hyperparams: batch size, learning rate multiplier, epochs (0.1 to 2)
Issues: SoM
- Loss instability: learning rate multiplier
- Catastrophic forgetting
- Overfitting: early stopping
• Data Augmentation. (another LLM to create syn. data)
- Out of memory errors:

H100 > A100
faster training
faster inference

• Inference

- Evaluation: LLM-as-a-judge with Human-in-the-loop
- Monitoring in Production
- Inference: VLLM (Dynamic batching + Paged Attention)
- AWS Bedrock, Hosting in Intuit VPC
- Services: K8 + Docker + Autoscale + Load balancer.

Anthropic's API, Bedrock API
Hugging Face (PEFT)

• Software Packages:

Axolo
Llama-factory
Llama Index
Haystack

- ~~BREADTH & DEPTH~~
- BUSINESS PROBLEM & CONTEXT + BUSINESS GOAL
 - YOUR ROLE & CONTRIBUTIONS
 - EXISTING SOLUTION & PROBLEMS IN EXISTING SOLN
 - PROPOSED SOL & ARCHITECTURE
 - CROSS-FUNCTIONAL TEAM COMPOSITION & KEY STAKEHOLDERS
 - STRATEGY & SUCCESS CRITERIA (Sign-Post metrics
or secondary Metrics)
 - TECHNICAL DETAILS
 - High-level Components
 - Challenges & their sol
 - Trade-offs
 - Tools & Technologies used
 - BUSINESS IMPACT
 - LEARNINGS
 - FUTURE PLAN

LLM PROJECT

COMPANY / IT VISION
& STRATEGIC PILLARS

: Taxes done for everyone with
✓ 1) Ease
✓ 2) Confidence
3) Putting more money into
the pockets of our customers

CONTEXT: While filing taxes, a customer has multiple questions
e.g. what is child tax credit? → FAQs

Am I eligible for EITC? → Personalized FAQs

Why is my refund \$x? → Outcome Explanation

What is the effect of withholding \$x on my taxes? → Tax Advice

CURRENT SOLⁿ: Static content owned by diff teams

e.g. eligibility
outcome
per tax topic

PROBLEMS IN CURRENT SOLⁿ:

- 1) Static → Lots of maintenance effort
- 2) Don't have ability to ask follow-up questions
- 3) Not personalized (some sections)
(e.g. gen. description of who qualifies for EITC but not whether the user qualifies or not)
- 4) Cannot customize the response for tone, comprehension per user group

BUSINESS PROBLEM &
BUSINESS GOAL

Helpfulness rating
80% → 64%
 $\Delta 16\%$

↓
↓ 50% S2C

CROSS-FUNCTIONAL

STAKEHOLDERS: Product, Core Engg, Analytics, Legal, Finance, Compliance, Tax Analysts,
Security, Content Design, Domain Experts

TECHNICAL STRATEGY:

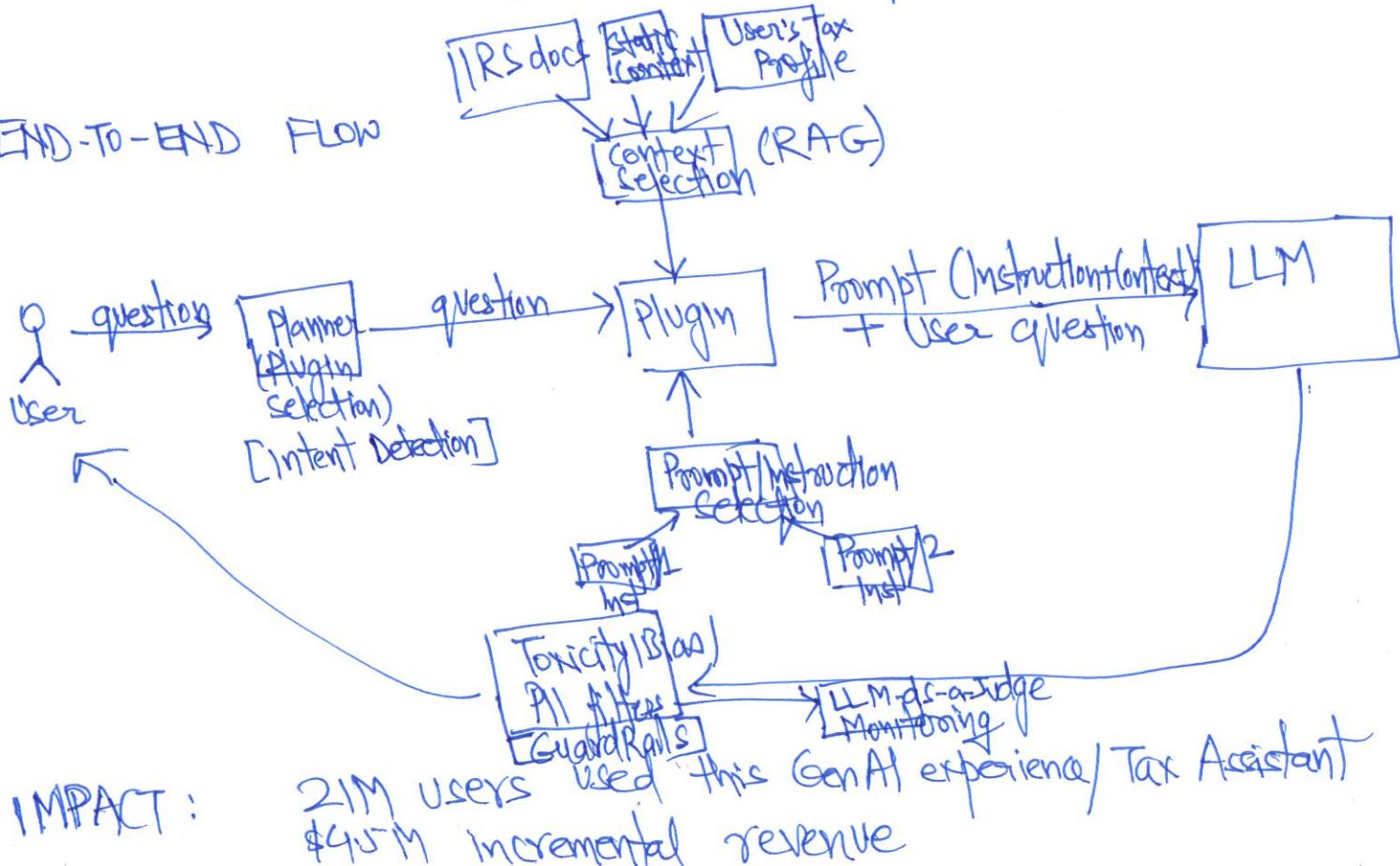
- Step #1: IRS docs : FAQ
- Step #2: Personalized FAQs
- Step #3: Personalized Outcome Calc.

(Issue: LLMs not great at math/calc
Soln: plugin dev)

Step #4: Coordinating above plugins together
in chat via Planner
(to switch seamlessly b/w topics)

SUCCESS CRITERIA:		
REQUIREMENTS	i) Accuracy	>98% (Compliance business)
	2) Relevance	>90%.
	3) Coherence	>90%.
	4) Hallucination	<1%.
	5) Latency	<1 sec
	6) Cost	(Acceptable by SLT)

END-TO-END FLOW



COMPONENTS DETAILS

1. LLM Selection : Accuracy Vs Cost Vs Latency

Evaluated GPT3.5, GPT4, GCP Palm2, Claude Instant, Claude 2

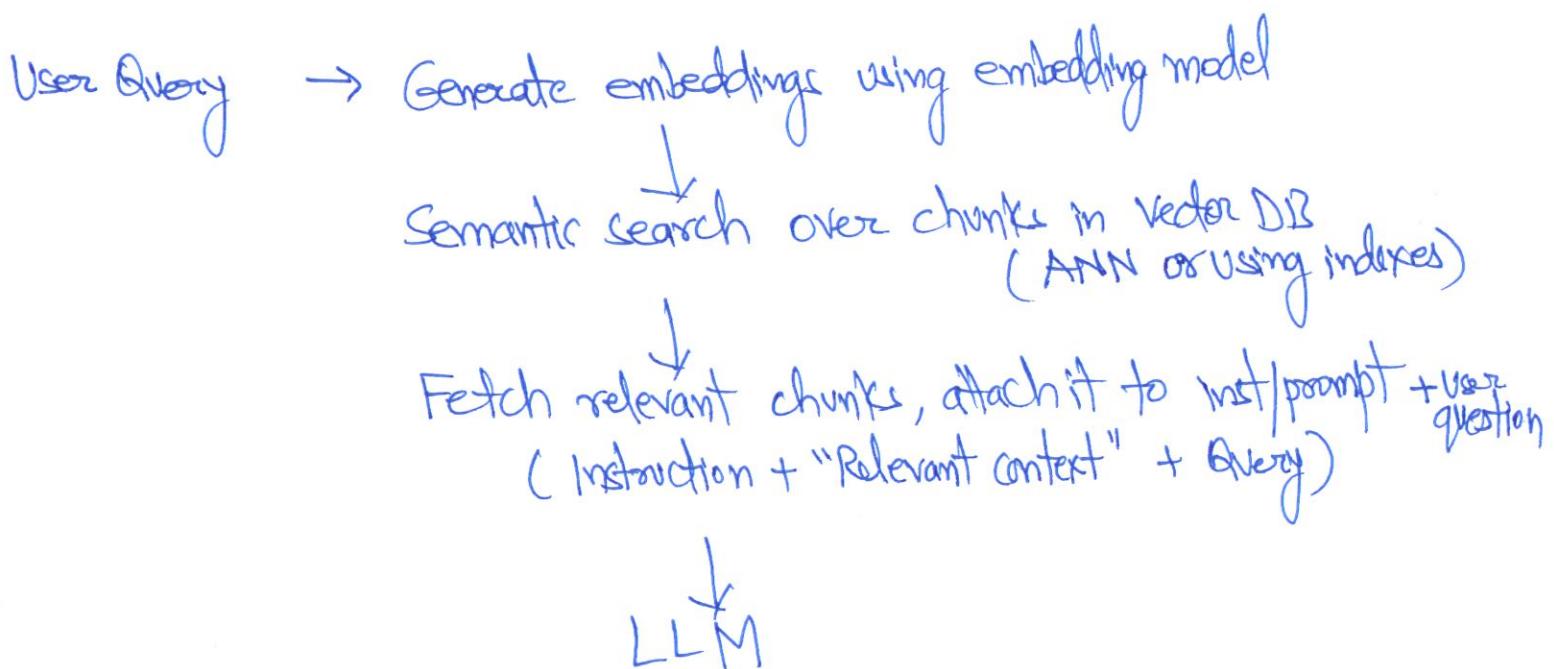
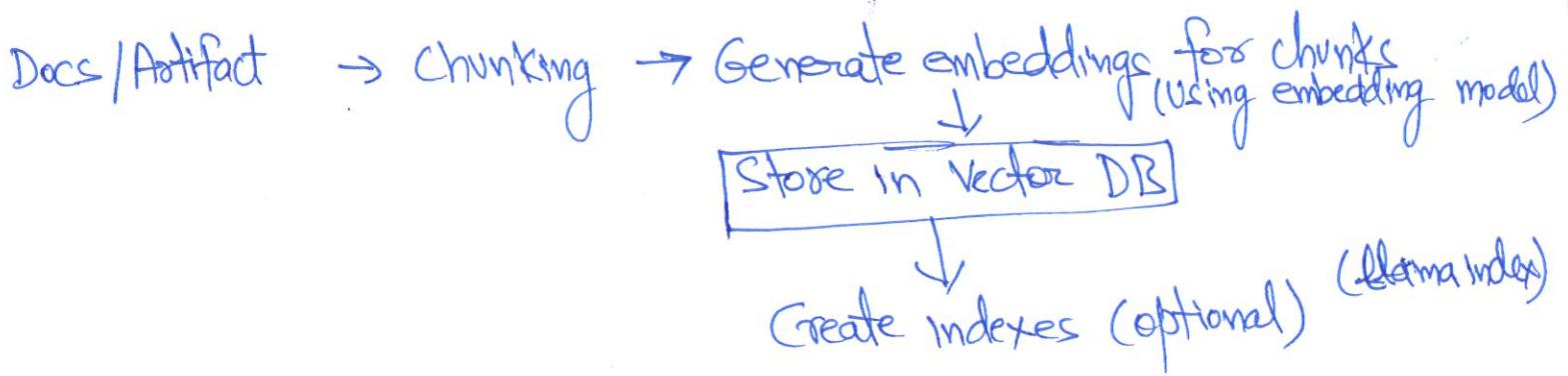
Winner: Claude Instant
(
Latency < 1 sec
Accuracy ✓
Cost ✓
)

2. INSTRUCTION / PROMPT : - Diff. prompt engg. techniques for diff. LLMs
e.g. GPT vs Anthropic

- Zero-shot prompting, few shot prompting
Chain of Thought etc.
- Iterative & labor intensive
- security Testing for prompt injection/
prompt leakage
- Prompt/InstSelection based on later
question intent / topic

3. CONTEXT (RAG) AUGMENTATION:

- Augmenting tax specific info
- Latest tax info (LLMs are stuck in time
due to their training data cutoff)
- Overcomes context length limits
- Cost effective (due to lower cost of
less no. of input tokens)
- Biggest lever before fine-tuning
& after instruction/prompt



4. EVALUATION :
- Automation
- 0) Human Evaluators
 - 1) Syntax related metrics e.g. BLEU, ROUGE → Human evaluation
 - 2) Semantic metrics b/w user question & response (cosine similarity b/w embeddings) → Human eval.
 - 3) Using syntax & semantic metrics as features & (evaluation verdict (by humans), build model → NLP is TASK-SPECIFIC AS labels ANOTHER LLM AS A JUDGE
 - 4) i) Regression : "Expected Response" Known
 - Using (i) checklist (ii) steps followed by tax experts (iii) corrected response , code a regression prompt (iv) few-shot examples for another powerful LLM as a judge
 - ii) Online : "Expected Response" not Known
 - CheckL-generated response bounded by "context" using another LLM as a judge

5. FINE TUNING:

SFT

RLHF

① - Cannot fine-tune on ~~user~~ data → compliance reasons

SFT: ② - Fine-tuned open source & models available on AWS

Bedrock e.g. Falcon using synthetic data
but it did not meet accuracy bar ($>98\%$).
fair accuracy

RLHF :- Partnering with Anthropic (specialist)

(providing data & assisting in crowdsourcing
human evaluation)

- Can describe about Reward Modeling Issue
(ordering responses → Many human labelled
needed)

- CHALLENGES:
- 1) Context length limit (initially) : RAG
 - 2) Switch from GPT to Claude Instant
- prompt engg. challenges
 - 3) Automated Eval.
 - 4) Planner (oos questions)
 - 5) Calc. → plugin

- TRADEOFFS:
- 1) Latency + cost vs. Accuracy : (LLM Selection)
Claude Instant accuracy lower than GPT4

- 2) Prompt ✓ vs fine-tuning : Initial phase RAG ✓
(tax law can change)

Later: (combo of RAG + fine-tun)

WHAT COULD HAVE BEEN DONE DIFFERENT? / LEARNINGS

- 1) LLM to be fixed first (prompt, context, eval. tightly coupled with LLM)

2) Accuracy definition adapted to LLM world

3) Very hard to get tax accuracy > 90% (since we are in compliance business)

4) Automated Eval's a necessity for faster iterations

YOUR CONTRIBUTIONS:

- 1) Strategy & Roadmaps
- 2) E2E system Design / Architecture (Genes as well)
- 3) Implementation (initial soln) : c.g. RAG using llama-index
Prompt engg
Regression for fine eval
- 4) Leading E2E development + deployment + integration with other subsystems (cross-collab)

PROMPT

Vc

FINE-TUNING

PROS:

- No data to get started
- Smaller upfront cost

- Learn new info
- Reduce hallucination
- Lower cost per request (after fine-tuning)

CONS:

- Hallucinations
- Less consistent

- More high quality data requirement
- Upfront compute cost

WHAT FINE-TUNING CAN DO?

- Behavior Change / Improve consistency
style, tone, format
- New Knowledge



General Purpose → Specialized LLM

INSTRUCTION FINE TUNING:

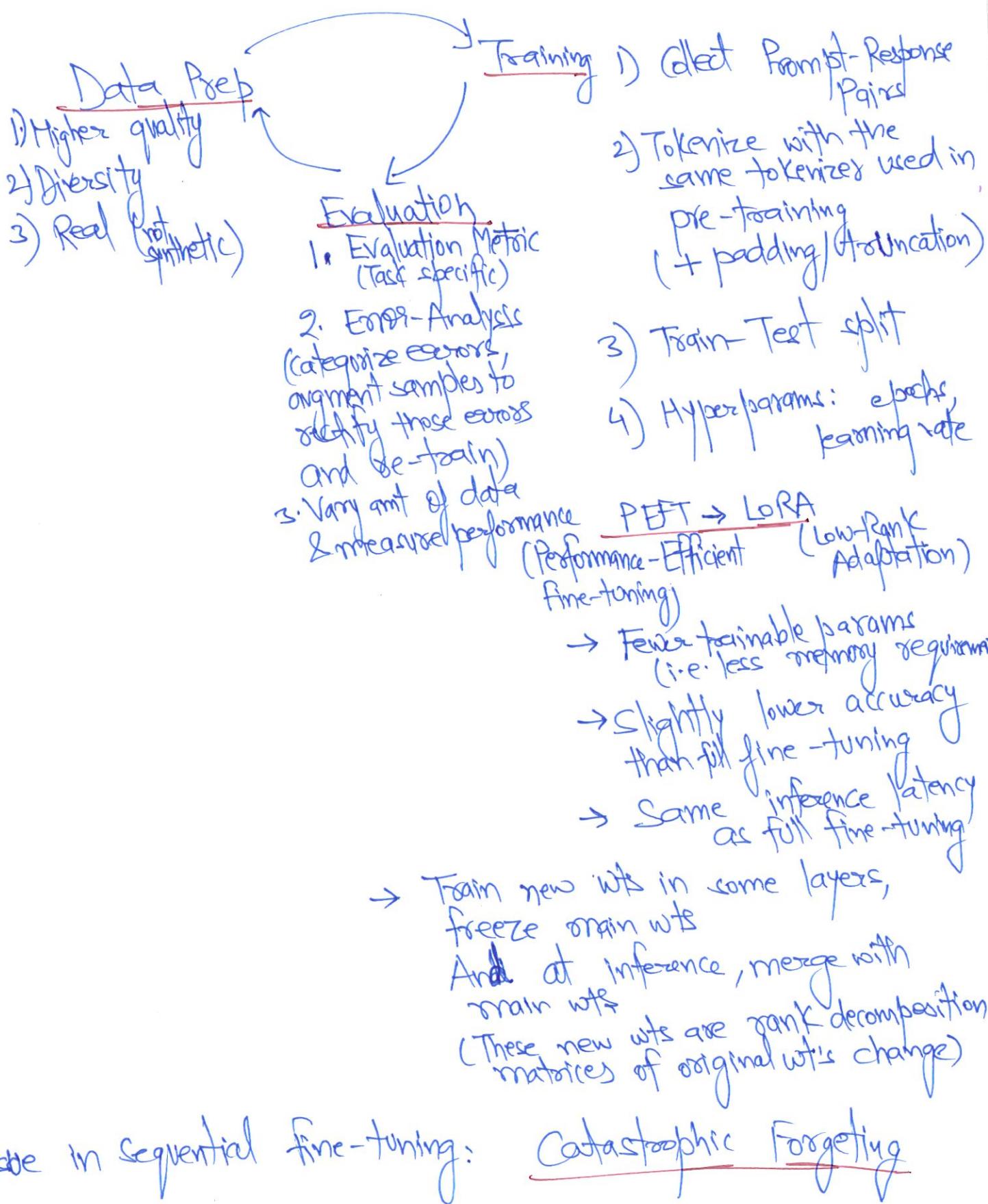
Data Prep: Instruction - Response Pair
(optional I/P)

e.g.: Add 2 numbers

Input: 2, 7

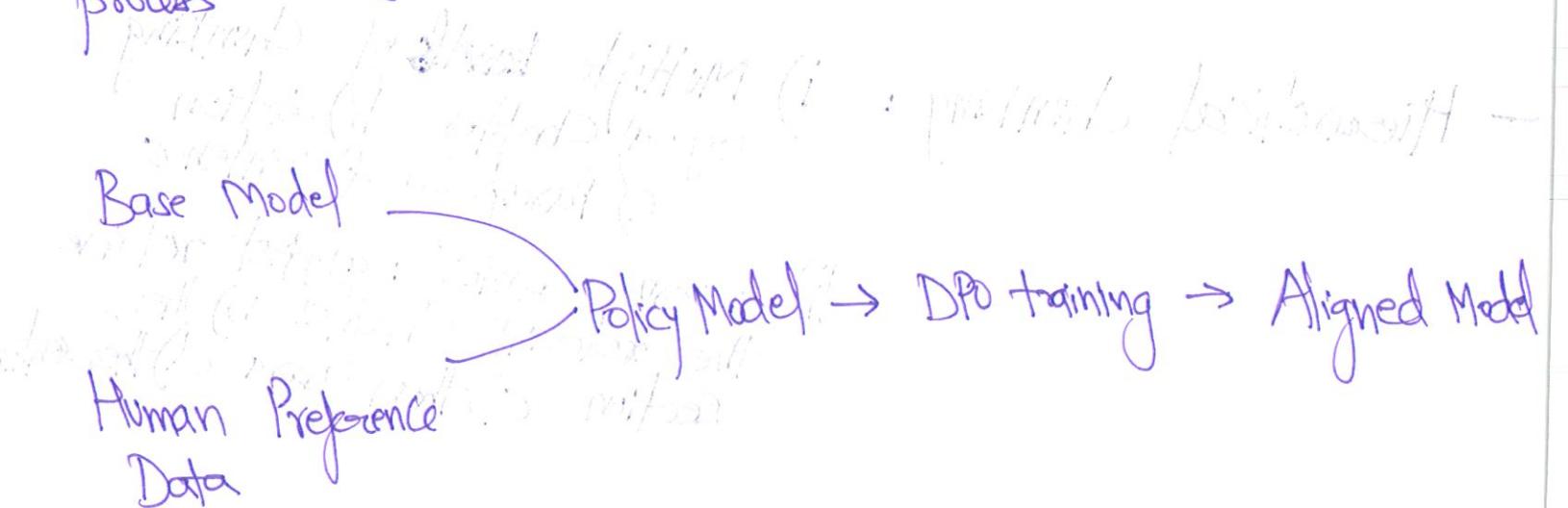
Response: 9

FINE - TUNING



DIRECT PREFERENCE OPTIMIZATION

- DPO is an alternative to RLHF
- Training data has good & bad responses (preference pairs) for the same prompt
- DPO reformulates the RLHF objective into a classification problem
- DPO makes the model learn to maximize the probability gap b/w preferred & non-preferred responses (e.g. cross-entropy loss) while the regularization term prevents it from deviating too much from the base model
- DPO implicitly learns a reward function within its training process



CHUNKING STRATEGIES IN RAG

- Fixed length : Cons → Same topic / content can be separated in diff. chunks (ideally should be in one chunk)
- Sliding window / overlapping window : Tries to solve above issue
- Paragraph / section based
- Semantic chunking :
 - i) Divide into sentences / paragraph / section
 - ii) Find similarity based on embedding of these sent. / para / section to all other sent. / para / section and iii) cluster them together and label the cluster with the topic of the cluster
- Hierarchical chunking :
 - i) Multiple levels of chunking
e.g. a) Chapter b) Section
 c) Paragraph d) Sentence
 - ii) Query comes : a) first retrieve the relevant chapter b) then section c) then para d) then sentence