	· Dois - ML SYSTEM	DESIGN
	Only Write Down the -	
ii) -		and don't look back There is anything else they want section before moving to next section
(ii)	Be correspond with y + Keep claritying + or make assumptions,	our interviewer throughout let the interviewer know and more forward
	You can't cover both Ack the interview	modeling and deployment in detail.
\ \ \	It is possible that in sections Instead of Jew i) write the fill focus on answer	Herviewer paints to dive deal into few lines about each cection, so outline with into line first past protect of the question into viewer asks instead of looking at time-remaining.
Ĭ.		1

-> Technical Debth (e.g. focalloss, deep hash embedding, Graffage -> Potential Icsues 2-their sofutions thanning-serving stew, -> Trade-offs / rationale for decisions

> Concise: one statement abt each billet pt.

MINGS TO STAND OUT IN ML SYCTEM DESIGN
> 2 Tower: Sampling ANN (Quantization & Indexing)
Single Task ve Multi-task both segression  Single Task ve Multi-task both segression  Single Task: Mide & Deep Model and Factorization machine Deep FM  Single Task: Wide & Deep Model and Factorization machine Deep FM  How watch hutory etc.)
> Features: Context features: sine & cosine fu to encode 23 is chosen
Sparce Features: Embedding Table, Deep hach Embedding  > Loss fn: Focal Loss
-> Trade-offs: - Pointwise, Pairwise, Listwise - rapposit vs NN - Factorization Machine Adv.
> Resys Special issues: D Cold start: user & item 2) Brias - Popularity - Position - Serving - Relevant feature list for the problem + their encoding +intuition
Relevant feature list for the problem + their encoding +intuition

				-	
With.					
•		•		•	
			4		
					•
					•
					•
				•	<i>i</i>
					; -

<b>→</b>	CLARIFICATIONS  Explanation of terms in question: search reb social media  Hately peech
<b>→</b> → →	Think of all user actions for the problem: will inform abt i) the subels contine metaics in the separation online metaics in the substituted of all likes, hide report interest of cyclem e.g modelity of data cultivated, hide report interest of cyclem e.g modelity of data cultivated in the substitute of subscribe followed in the separation of the subscribe followed in the seventh of the se
$\rightarrow$	Sale of System: How many veexs? How many items? Daily Active Usexs? Provide ballpark estimate
<b>→</b>	Performance of system: - Batch Vs. Real-time  - Ros of times undate is needed  (a.g. # of timeline undates)  Poovide ballpark estimate

> Any constraint: Legal/Compliana, Interpretability, Online togining?

-> Objective Fy: wotch Time | Reactions | Click Maximize num of Convections established (sent + accepted (in terms of things ML can offimize not reconcarily phonen 1 Classification only Ranking RecSys (Retoreval Ranking) -> Type of Pooblem:

- Mitti-stage or One-stage - Mitti-tack or lingle tack > Modeling Architecture

TOP THINGS TO HIGHLIGHT TO Hings in this booklem that I want to highlight to gain brownie HED (since time-limit)

KET FEATURES:
Only describe features that seem meet predictive for the problem intuition + encoding (e.g. strength of com. >> people you may know)

ONLINE METRICS: Think of metrics to measure when wer-item interactions
-> Polmary metalcs
-> Secondary metalcs
-> Secondary metalcs -> Counter Metalcs
-> Justify online metaics to high-level business objective
Lan
Justify online metaics to high-level business objective streng online metaics to high-level business objective streng tem impression, click, dwell time, reaction, further action e.g. burchuse comparation match time, event regulation come regularly come regularly to come regularly the commended content vs search in a session time to success (in search)  Time to success (in search)  The glaveries
constanting constant on recommended content is search in a session
want > Time to success (in search)
Level No. of queries
Level) No. of queries  > sessons with chicks (or intended adding. Functions)
( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( )
- Counta Metrica: hide, report, creatch time search time high =
Business : Users who re-violed in the rest x days  (rectantion)  (evel Conversion Rate
Right (referrition)
Business (rectention) level Conversion Rate
(Ang. or) Conversion Rate (Ang. or) Revenue Milt entire Bookmark Rate user base) Bookmark Rate
Satisfaction survey
Language.
FRANDIABUSE USE CONSENT CONTENT
A) - # of harmful content defected + reported by users
FRAND   ABUSE USE COSED:  The Interessions of harmful content detected  The Hof harmful content detected + reported by users  The Y-of narmful content detected harmful but appealed & reversed  - Y-of content detected harmful but appealed & reversed

-	,	•	•	-	,
					•
			4		
			; i		
			•		

LABEL CRITERIA AND TRAINING SET -> Annotated/Muman labeled Vc Existing onboarding/engagement Proc/ Cons: High Quality - Ve - hide - report - the pression to than t sec -> +ve - Chick + duel time >t sec - Reactions 3 may be weighted by all. relactions - burchase -registered for event -conn. sont + accepted - Session-based: items in context window > Graded Relevance for pairwise or latorise LTR - only impression, no dick = 0 - Click = 1 - Click + dwell +time > t sec = 2 - Click + dwell +time > t sec + further engagement = 3 -> Delay or Unavailability of labels: - Homan labellery booklem e.g. whether yet any seturn ever x days of days of any seturn ever x - one classification > Popularity brac: Restord pop. Hems by a cortain / of training set > Poivacy? need anonymization in data
> Favorece Bias for gendur, age race, underrepresented got represented in data (e.g. 9) by vol.)

· ~ . ^	
Ranking Classification: - Cluster -ve items - Pob. based campling	Instead of balancing use i) balanced cross-entry or befor ii) FOCAL Local
- Hard - Ve samphing (select - Ve items that are samplar to the items) - undersample - Ve set	, - <i>)</i>
- over sample the set - class with (e.g. scale par	tw
Cand. Generation i) in-batch random-ve gloods) In 2 tower: ii) in-batch-ve g combinity Hard -ve	notion
- CANDIDATE GENERATION  - Labels & Sampling	and dates
- How to kind condidates	(based on some herateric or Imp. feature)
- Why multiple and generators Metaics: (Recal) focuse	) Dhoreth
-> TRAIN-TEST CPLIT:  - VISUAlly by time frame/days, some mo le	eakage
- only downcample take care dustry val & test set feet the	ne come distribution of the &-ve labels
-> COURCING TRAINING DATA: Unstream pipeline the book of using Deeg	

FEATURE- ENGINEERING > Identify Actors and their hictorical interaction with "all of other ITEM: - main properties of item e.g. poice, location, duration, cuidne - modality of dataery. text, video, image > embeddings
- engagement: likes, shares, etc. 3 rate + reputation pagerank, subscribed Scounter features - freshwers/age Location Position - Demographics: age, gender, lang., country, city, accountage USER: - User-item historical interaction - Identify some imb. properties of item and wer's interaction history with these proposties of "all" Hem Counter and rate features - seq. features of items - User-other actor hichookal interaction
(2) TEM FEATILIDE SPARSE USER & ITEM FEATURES! (can quoid them since they come at huge, Computational Istorage Cost > embedding take size but Deep mach embedding can solve for HI Imp. properties of Hem e.g. director id,

NOTE: ONLY the above features go into 2 tower so that isolation to maintained for creating user and Hem embeddings. User Tower: User dense toward Them Tower: Item (dense toward)

CONTEXT! Specific to user-item pair: time of day of week, aggregate ed 2 towers architecture · . . . of media consumed on mobile feature) at user level CROSS-FEATURES: "Specific" User- item - other Ador interaction - "specific" User-item interaction e.g. how many times user how watched the genre of more that A to be make - specific" was -other ador interaction e.g. parts liked by user created by auth Note: The above features go into Ranking Classification where specific user-item pairs are ranked classified (along with user & item features: dense + sparce) NOTE: In case of Ranking/Classification after 2 tower N/W, the features used are: 1). Candidate generation 1). Candidate generation score (dot product blu user & item embeddings) 2) Aug. embedding of wer-item historical interaction e.g. watch, history, like history (seq. of items) 3) (nos- Featives.

## ENCODING FEATURES

GPARCE FEATURES: EMBEDDING LAYER Isave -> Using embedding layer for sparse features with high condinality results in evalueding loopup table getting huge and shows the porcess of training inference significantly e.g. user is mapped to embedding victory (128 dim.) results in embedding lookup table of size (2000,128) 1) Hashing to reduce size of lookuptable  $\mathbb{Z}_{p/p}$  : 2) Using rate & counter features

Summary statistics count summary statistics count

By there: country time many

maryin 2 months min min\_in\_works 3) Deep Hash Embedding Morandization or Standardization or Birming NUMERICAL FEATURES: -one-hot encoding: it coordinality is low CATEGORICAL FEATURES: [ Feature-hashing: when lot of values/high coordinality Adv: can handle unseen categories -Embedding: usually jointly learnt with DNN smood (extract layer to convert shores one-het encoded vector to dence embeddings) - 4f 2 tower used in cond. generation phase then fetching, embeddings for items and have SEDWANTIAL FEATURES: avg. them ( (e.g. watch history)

- Attention layer and then concatenating the off

- Embedding layer and then aggregation the

TIME BASED FEATURES; Sine & cosine for so that 23 is charator

LOCATION BASED (Lat Long): Haversine Distance

IMAGE, VIDEO, ANDIO: Pre-processing & using bre-tained

TEXT-BAJED:

Word-level: Word2/ec (not context-chare)
Clause or centencer/evel: pre-trained models like
BERT embeddings

Model

Rawline Model

Multi-stage

Or Ore-stage multi-tack

(retrieval, filtering, ranking, re-ranking) in actory to injust the injust of generaling cardinates + limiting their number

Ist stage filtering based

(first degree connection park wheen), pages people followed, growns joined

events in the skimile todius,

ads for specific control age, gender)

Single Tack VS Multi-tack

Mine & does not req. User Jitem Dow- features

The model

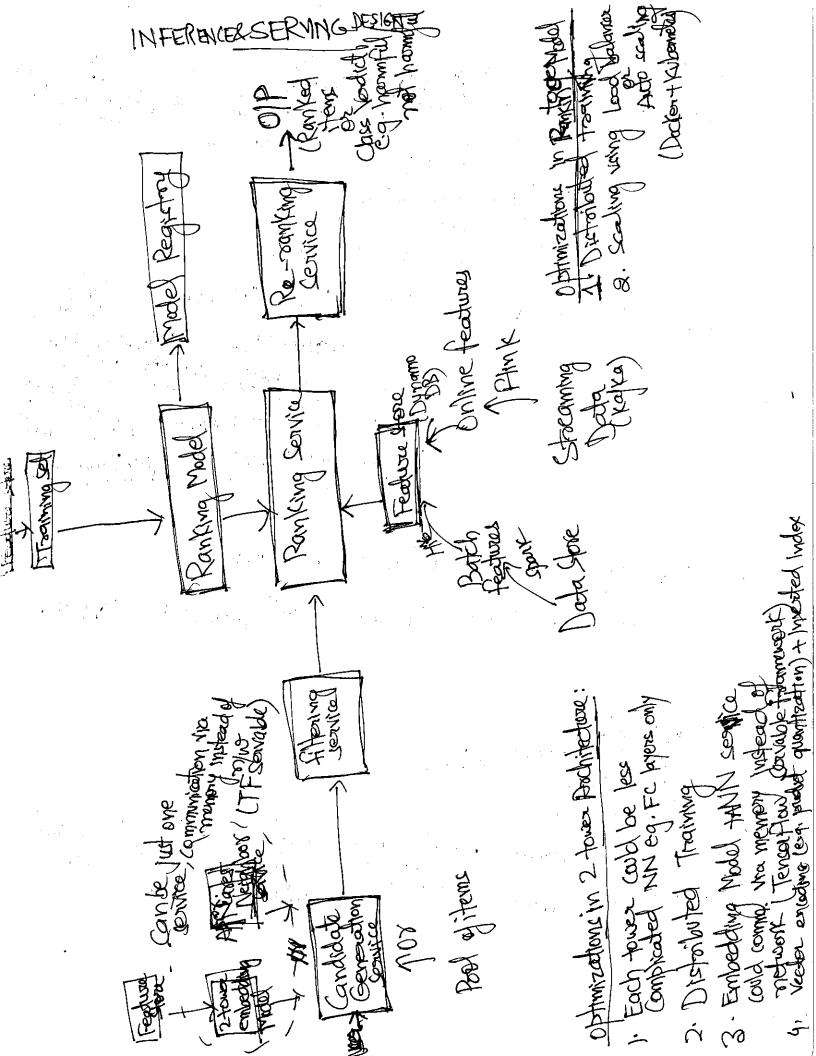
The stage of the > Inference is usually Real-Time
> Search System: D-tower: 1) User-away tower
> Search System: D-tower: 1) User-away tower

> Tor Ranking Classification Model to Bet the store in ranking phase)
> Tor Ranking Classification Model to Start with boseline
of Baseline model Rule-based exterm: Start with boseline
(more accurate but more difficult
> Pointwise, Pairwise, Licture: Pour 2 Gus (but more difficult
> Pointwise, Pairwise, Licture: graded relevance to implement strain) Fraining data: graded relevance -> 2) Afriter chaosing pointage, primarily the decidon Choosing NN: I date is unstructured & embedding by the of sparae leatures is needed I continual learning I training & inference (slower than r: Focal loss Vimbalanced + multiple modalities in data) Balanced Cross-Entropy (Vinbalanced data) -> Training pipeline, managing model version, Distributed Training 2. Less # of labels for fine-tuning -> For NLP or mage based models: fine-toning on tack specific

- Other water of generating card: popular items, trending items - Cold start in 2 tower architecture: (no interaction date) or in general forcerd, generation medite into to get embedding new users: Only use were medite into to get embedding lind top K rearest neighbor Herne (2 Hower architecture accepts new users) - new items: if item profile available, use 2 tower to get item embedding and stone in embedding table and use it while Anding top 1 nearest neighbor items (2 tower can work with new items it features available) give pool & cons of Matorix Factorization & 2 Tower - In case of Peccys Matory Factorization 2 lower - Training efficient Courtly to toal - Inference faster (as embedding the leatures need are static & combe por completed) to be found formed into combeddings at quay fine Any not needed if and placed in the second section on war 2 item for factor Ann : The dor exceeding engineering invested index index (here such a horizonte)

Tor factor Ann: The dor exceeding engineering invested index (here such a Norvigable HNSW (here such a Norvigable LSH (Locality sensitive Maching) - De rolle Content - Re-sanking Service: - Dressity (ang. parmoise similarity blue ranked list of items) - Freshness t De-dwp on the many content or reported their discount the score por contain ant = modulity and the

	OFFLINE METRICS	
-> MoHi-stage:	Retrieval (and Generation	phase: Recall metrics
·	Parking: NINGOL MAP	e.g. Kecall & K Palasion & K
	Ranking: ND(60, MAP (predsion) Freekhness, Di- metrics)	versity.
> Classification	on? Precision of Recall PRI (sometimes have to dec	AUC, F1, AUCROC ide b/w precisions recall



- How to mittgate serving bias steedback look issulfablea-took
- How to reduce later of the case teal-time of pop Rasys).

- Pre-computing feature | embeddings - Model Compactsion - Use complex feature-engg. - Use more nardware infra Kafka + Flink -> Dynamo DB (an fine. - Infrastala Design: Data storements Hive (offlind features)

Model Registry bjed storage: to store model

Relational DB: to store model E-9. Sagemaker - APE: to retnew metadate and metadate and - COST: Spot ve dedicated instances (for forming)

Cost Monitoring via dashboord (weekly, monthly, projections etc.) ATB - Docker & Kutherneter

Testing - Shadow Made

Testing - Canary Deployment > gradually solling out an undate to a small

- Canary Deployment > gradually solling out an undate to a small

- Refreshing model to the entire were base

MONITORING & u 
At the very least: talk about 1) Model Metrics

(if label is available)

If label is rot-available

Vipog rate etc. 2) Operation metales 3) Training-serving skew

Log Model Request + Response > sanity thack & future

1. Ranking - Interleaving

2. Classification - Usual AIB test setup with MDE, sample size and baric steps

-> Novelty Effect

## SPECIAL ISSUES IF NOT MENTIONED

- 1. Serving bias 2. Popularity bias 3. Position bias

  - 4. Cold Stoot problem

-> Calibration

Problem exploration Modeling approaches Infrastructure design Training/testing Deployment

Strengths Gave an overview of the template used Asked about the framework Mentioned what type of actions Asked about where this is surfaced eg homepage Candidate mentioned maximizing both the send and the acceptance Identified the top objectives as multiobjective on click + acceptance Mentioned users see new results on each refresh Mentioned a bunch of online metrics to measure the system Mentioned user engagement after x days as additional metrics Talked about the main actors in the system Mentioned about interaction features Mentioned network features eg number of common friends Also mentioned second hop friends as candidates Proposed a heuristic way to threshold the multihop candidates On mentioning cold start came up with an alternative recall source to mitigate this Mentioned categorical feature encodings Mentioned using sparse id and deep hash embeddings -> technical depth Mentioned deep and wide models and factorization machines as they are good with feature interaction Mentioned using graph sage - graph NN for ranking Mentioned using focal loss for ranking Mentioned using recall specific metrics Mentioned using NDCG ranking specific metrics Mentioned training data and serving have no distributions

Areas to improve Mentioned 10 times refreshed the system without justifying Better answer is ask about latency of search What is the goal? User experience. How do users interact with the system? How much personalization? Scale, users, DAU, items, QPS, explainability, retrain frequency, latency How much data do we have? Whether data is available or existing onboarding flow Need to talk about primary metrics, secondary metrics and guardrail metrics Justify business metrics to target, connect business needs to ML decisions Different characteristics in the industry + risks of a complex design Ask about the end to end product flow of the system How fast should the model be? Is precision or recall more important Can include more graph based features eg page rank Can mention collaborative filtering on candidates Need to mention cold start eg using friends of friends Can mention filtering friends in candidate selection based on user affinity When mentioning embeddings need to mention what we embed and how Can use graded labels instead of 1 or 0 No need to explain focal loss in detail - need better explanations eg importance weighting Its very important to talk about training data -- how to source it/pitfalls etc Need to talk about ab testing/canary deployment Important to touch upon biases -> eg position bias, popularity bias etc

•