

Topic 9

Traditional collaborative filtering vs LLM enhanced recommendation

Xian Haotian Lan Tiancheng Jia Xiaoran Wang Liming Li Rongchao

PROBLEM FARMING

TASK DEFINITION

Learn a scoring function to rank items for a user under context:

$$r^{ui} = f(u, i, C)$$

POINTWISE

PAIRWISE

LISTWISE

$$\min_{\Theta} \sum_{(u,i) \in \Omega} \ell(r_{ui}, \hat{r}_{ui}) + \lambda \|\Theta\|_2^2$$

$$\max_{\Theta} \sum_{(u,i,j)} \log \sigma(\hat{y}_{ui} - \hat{y}_{uj}) - \lambda \|\Theta\|_2^2$$

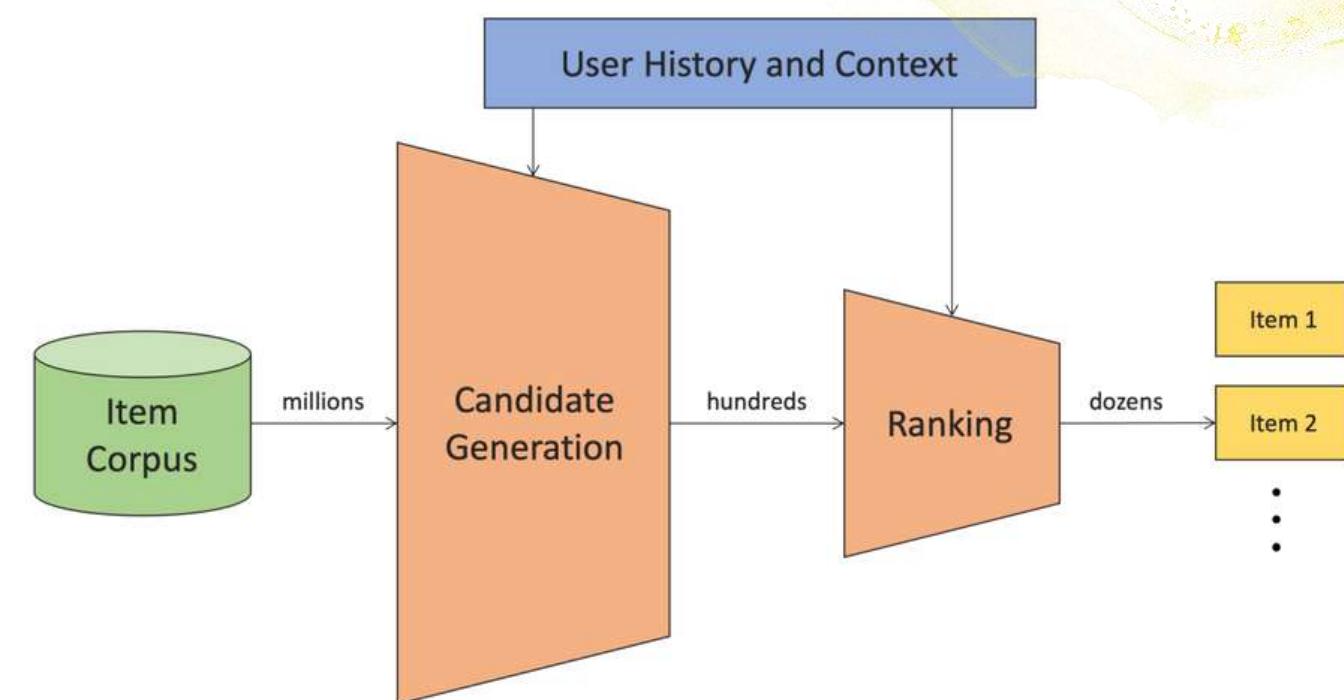
$$\min_{\Theta} \sum_u \mathcal{L}(\pi_{\Theta}(u), \pi^*(u))$$

OBJECTIVE
FAMILIES

Three-stage pipeline: recall candidates, rank with CF/graph/sequence models, then re-rank under user intent.
LLMs augment—not replace—scorers: semantic embeddings unify vector spaces, reasoning extracts constraints and explanations, and orchestration coordinates retrieval, ranking, and RAG grounding for controllable, traceable recommendations.

Image suggestion

Simple funnel diagram (“User History and Context → Candidate Generation → Ranking”). Good visual for Recall→Rank→Re-rank framing; add a “Re-rank” box in your slide.



ASSUMPTION:

USER-ITEM PREFERENCES ARISE FROM A FEW LATENT FACTORS IN A SHARED LOW-DIMENSIONAL SPACE.

SCORING (MF)

SIMILARITY OF LATENT VECTORS DRIVES RATINGS.

TRAINING (EXPLICIT)

MINIMIZE MSE WITH L2 REGULARIZATION; OPTIMIZE VIA SGD OR ALS.

STRENGTHS

HIGH THROUGHPUT, STABLE OPTIMIZATION, SIMPLE TO SCALE AND COMBINE WITH SIDE FEATURES.

LIMITS

COLD-START PAIN, WEAK SEMANTICS, LONG-TAIL UNDEREXPOSURE, LIMITED INTERPRETABILITY.

TAKEAWAY

MF LEARNS CO-OCCURRENCE STRUCTURE; PAIR WITH LLMS TO INJECT SEMANTICS, CONSTRAINTS, AND EXPLANATIONS.

$$\hat{r}_{ui} = \mu + b_u + b_i + \mathbf{p}_u^\top \mathbf{q}_i$$

知恵で薦める

BPR

$$\max_{\Theta} \sum_{(u, i, j)} \log \sigma(\hat{y}_{ui} - \hat{y}_{uj}) - \lambda \|\Theta\|_2^2$$

Graph (LightGCN)

$$\mathbf{E}^{(k+1)} = \tilde{\mathbf{A}} \mathbf{E}^{(k)} \quad (\text{propagation})$$

$$\mathbf{e}_u = \sum_{k=0}^K \alpha_k \mathbf{E}_u^{(k)}, \quad \mathbf{e}_i = \sum_{k=0}^K \alpha_k \mathbf{E}_i^{(k)} \quad (\text{layer aggregation})$$

$$\hat{y}_{ui} = \mathbf{e}_u^\top \mathbf{e}_i \quad (\text{prediction})$$

EFFICIENT, CAPTURES HIGH-ORDER NEIGHBORS.

SEQUENCE (SASREC/GRU4REC)

NEXT-ITEM FROM ORDERED CLICKS; SELF-ATTENTION/RNN MODELS SHORT-TERM INTENT.

BRIDGE

STRONG STRUCTURE, WEAK SEMANTICS ☐ PAIR WITH LLMS FOR INTENT/CONSTRAINTS/EXPLANATIONS.

SEMANTIC BLINDNESS: CF/GRAFH/SEQUENCE RELY ON CO-OCCURRENCE; WEAK AT UNDERSTANDING TEXT, IMAGES, ATTRIBUTES, AND NUANCED INTENT.

COLD-START & OOD: NEW USERS/ITEMS, UNSEEN VOCABULARY, CROSS-DOMAIN TRANSFER ALL DEGRADE.

COMPLEX CONSTRAINTS: BUDGETS, STYLES, COMPATIBILITY, POLICIES—HARD TO ENCODE AND ENFORCE IN TRADITIONAL SCORERS.

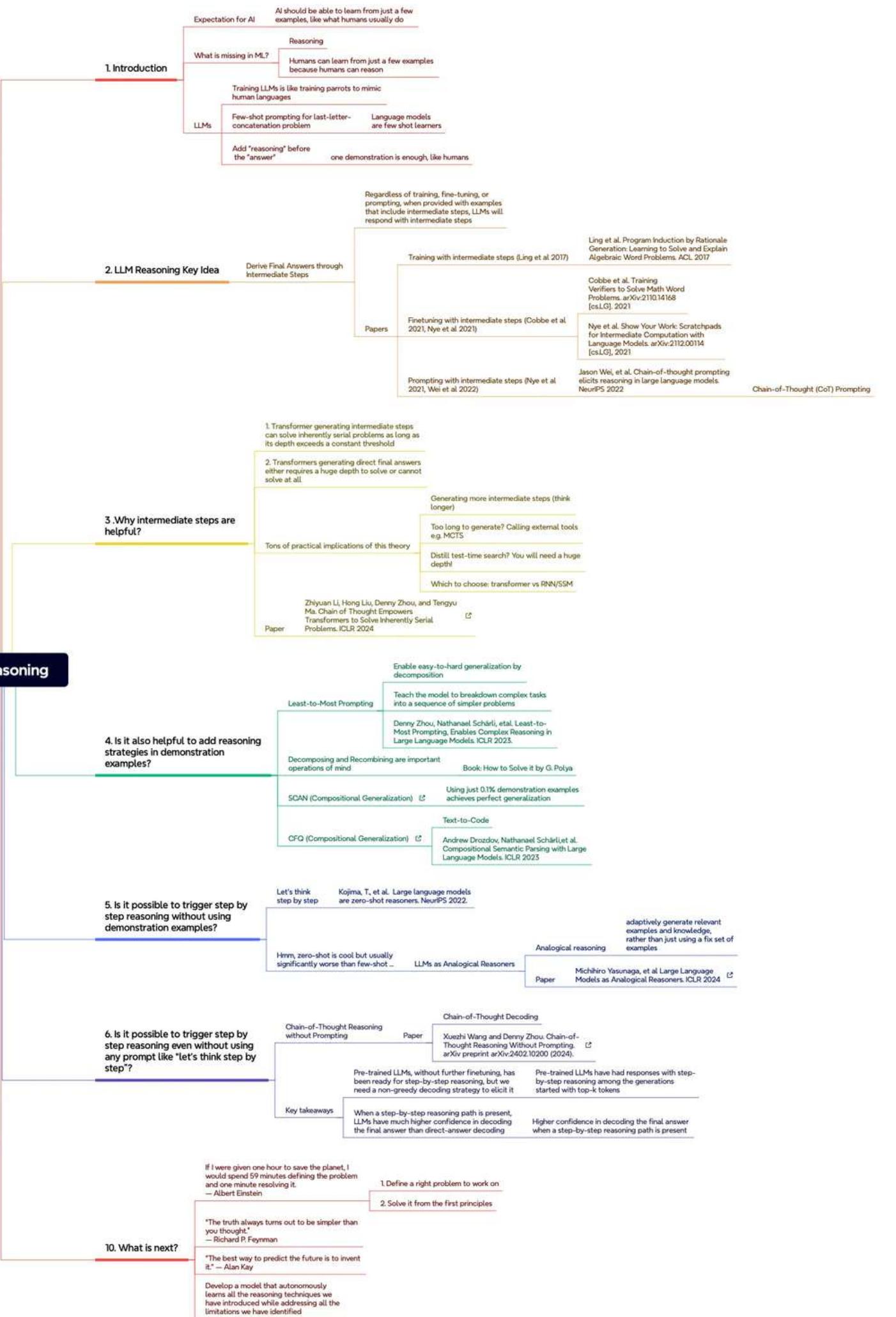
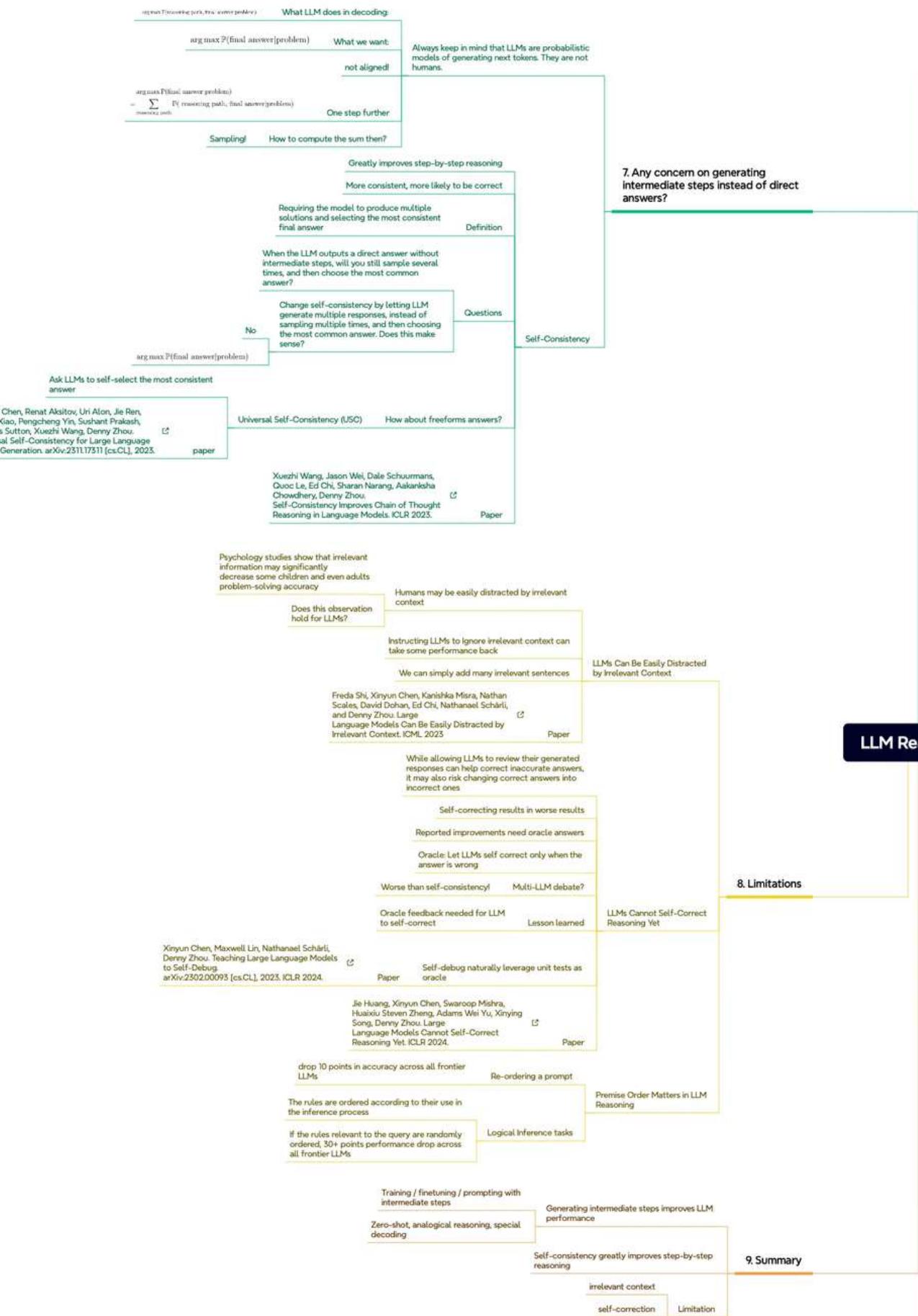
EXPLAINABILITY: HARD TO PRODUCE FAITHFUL, USER-FACING RATIONALES TIED TO EVIDENCE.

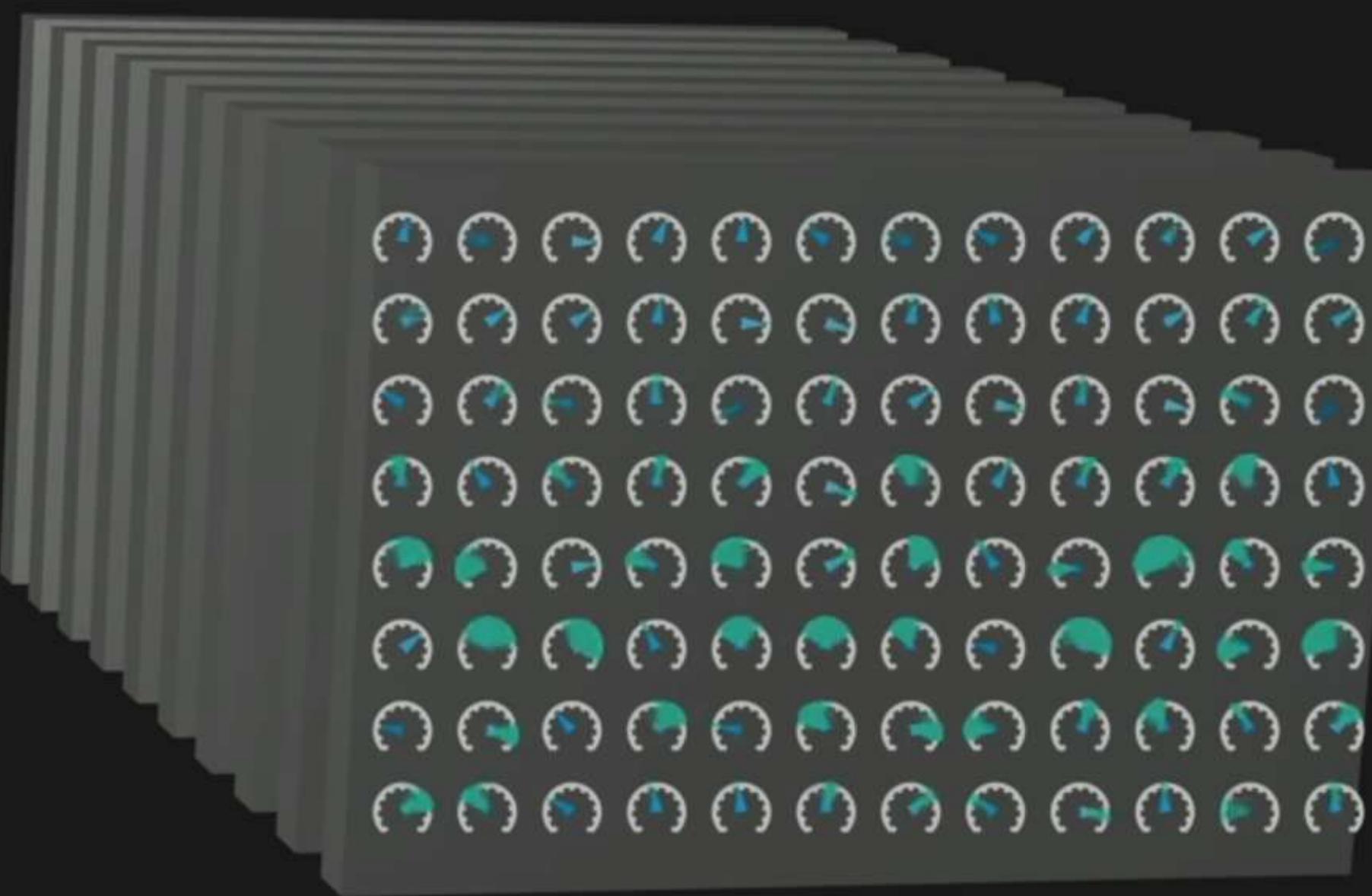
LONG-TAIL BIAS: POPULARITY AMPLIFICATION; LIMITED NOVELTY/DIVERSITY CONTROL.

SPARSE/TEMPORAL DRIFT: DATA SPARSITY, EVOLVING TASTES, SEASONALITY HURT STABILITY AND CALIBRATION.

MULTI-OBJECTIVE TRADE-OFFS: RELEVANCE VS. DIVERSITY, BUSINESS RULES, SAFETY—DIFFICULT TO BALANCE DECLARATIVELY.

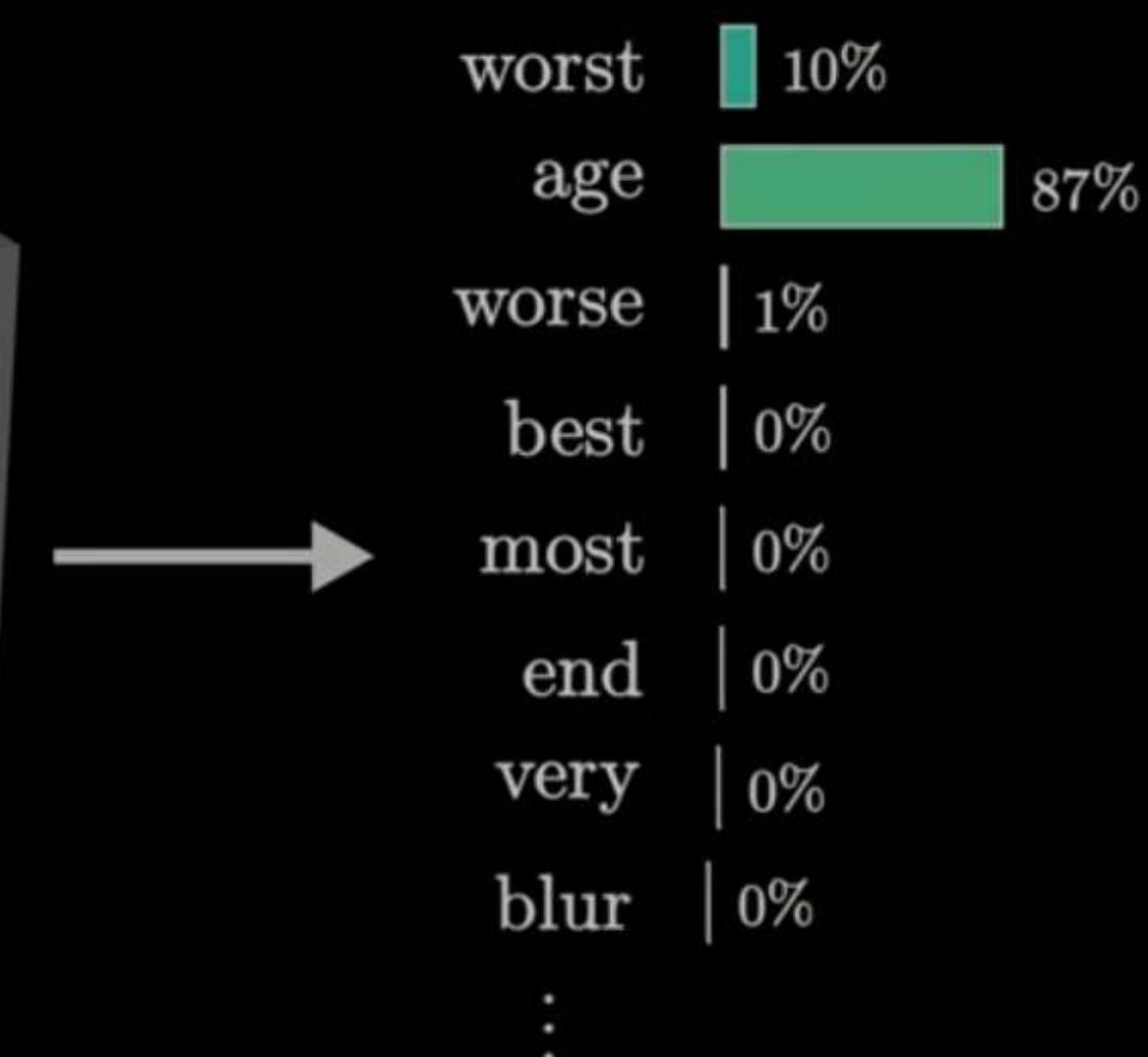
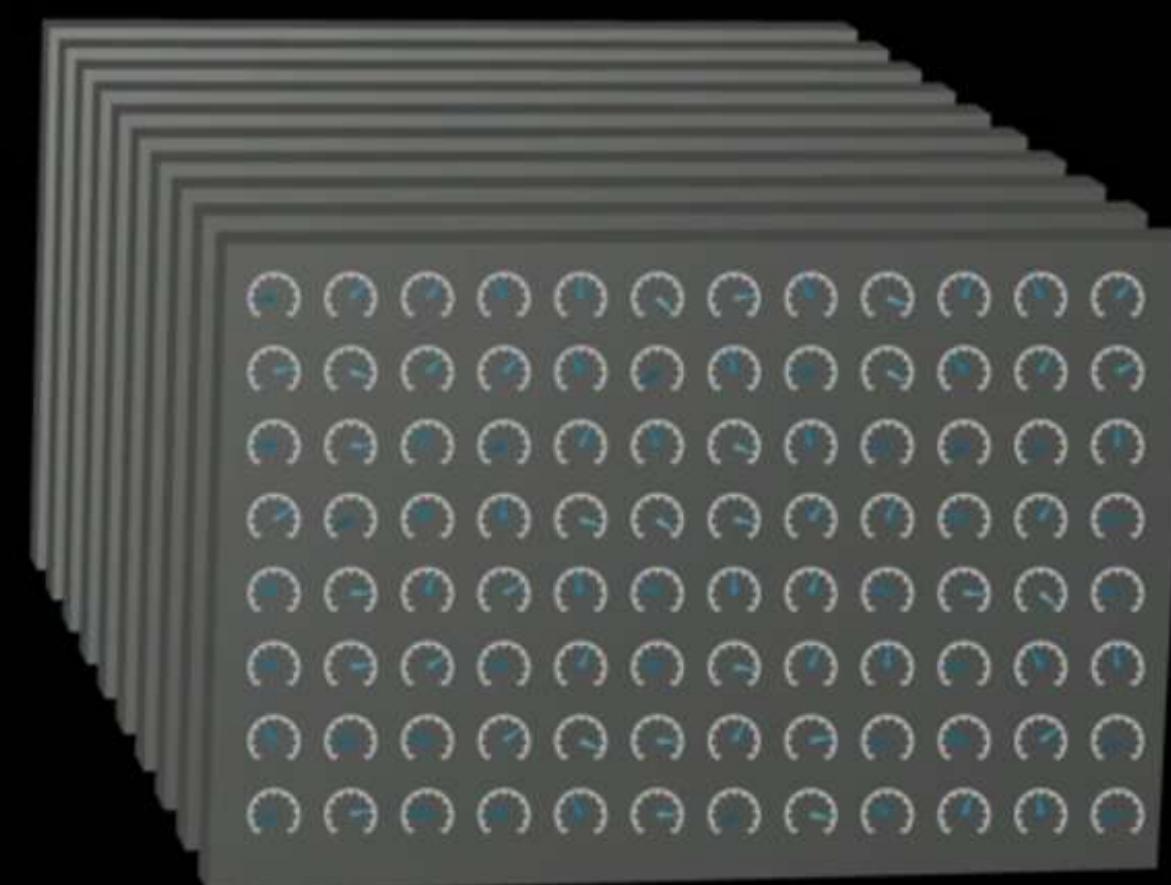
PRIVACY & COMPLIANCE: SENSITIVE SIGNALS REQUIRE SELECTIVE USE, REDACTION, AND AUDITABLE DECISIONS.



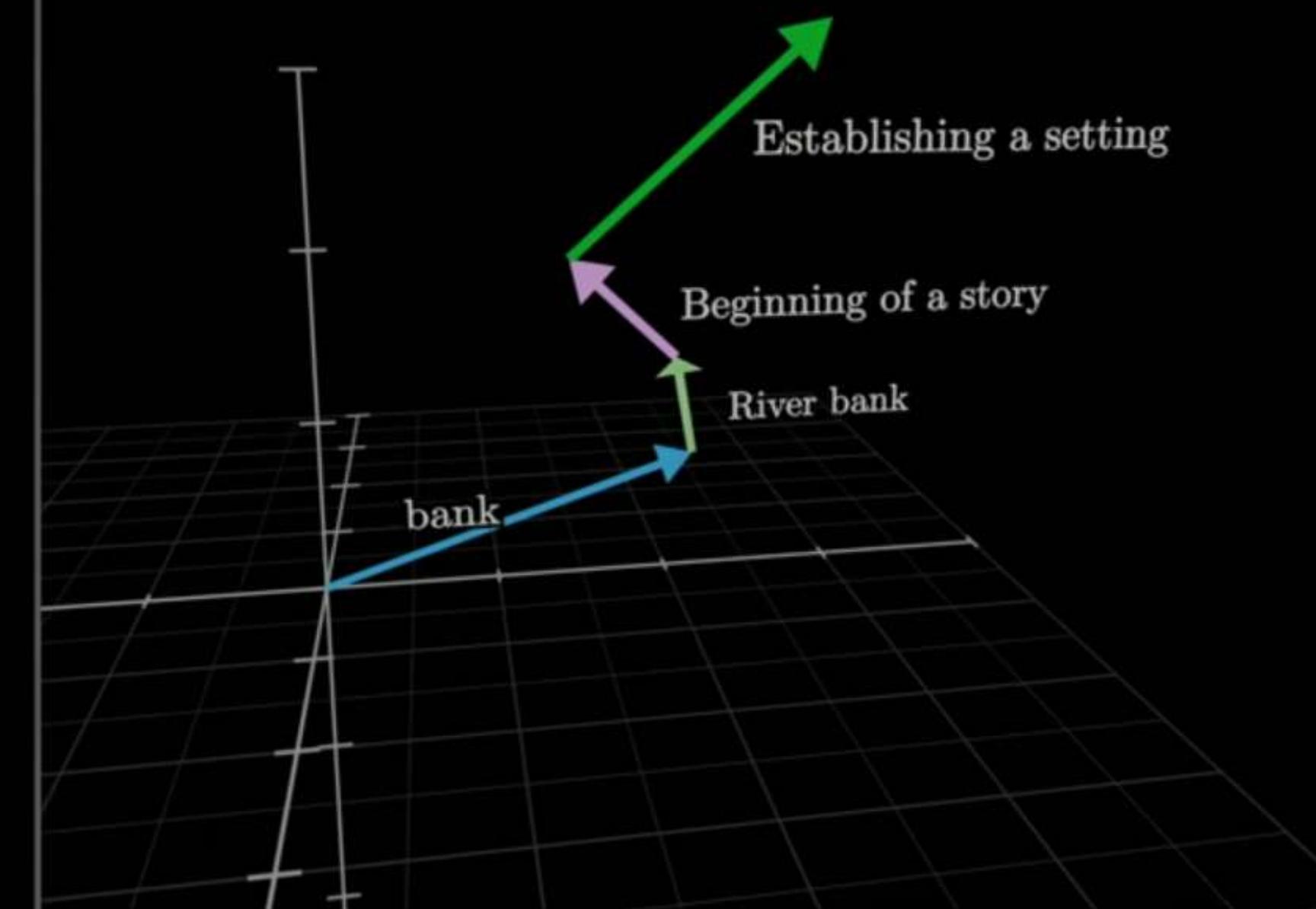


It was the best of times it was the **worst**

It was the best
of times it was →
the _____



Down by the river bank



Theoretical Summary

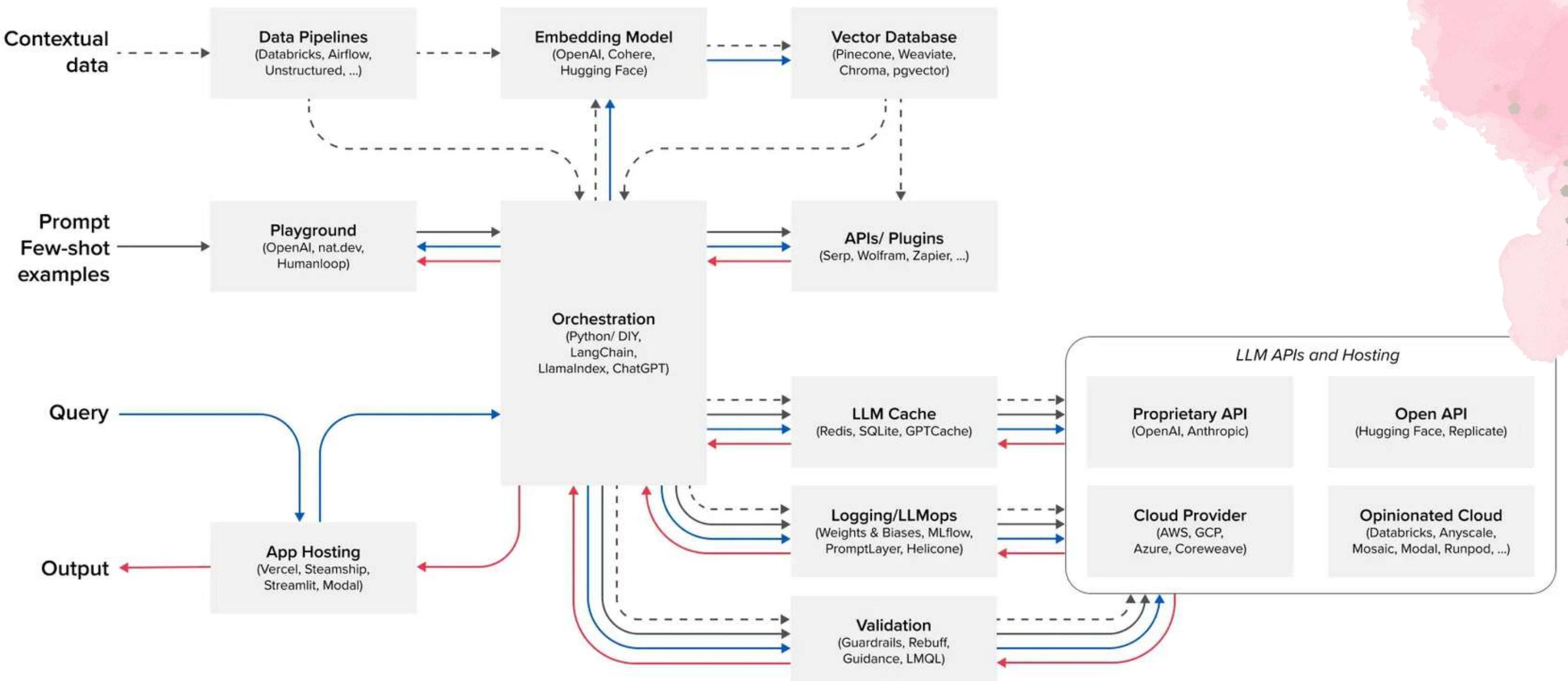
LLMs enhance recommendation pipelines by unifying semantics across recall, modeling, and ranking.

Candidate Generation: Texts (titles, reviews) are embedded into a shared vector space for ANN retrieval, combining sparse BM25 and dense vectors to improve findability.

User Modeling: LLMs summarize history and dialogue into structured JSON constraints (e.g., `price_range`, `must_have`).

Re-ranking: Using these constraints, LLMs re-rank candidates with interpretable scores and rationales, balancing relevance, diversity, and controllability.

Emerging LLM App Stack



LEGEND

- Gray boxes show key components of the stack, with leading tools/systems listed
- Arrows show the flow of data through the stack
 - - - → Contextual data provided by app developers to condition LLM outputs
 - Prompts and few-shot examples that are sent to the LLM
 - Queries submitted by users
 - Output returned to users