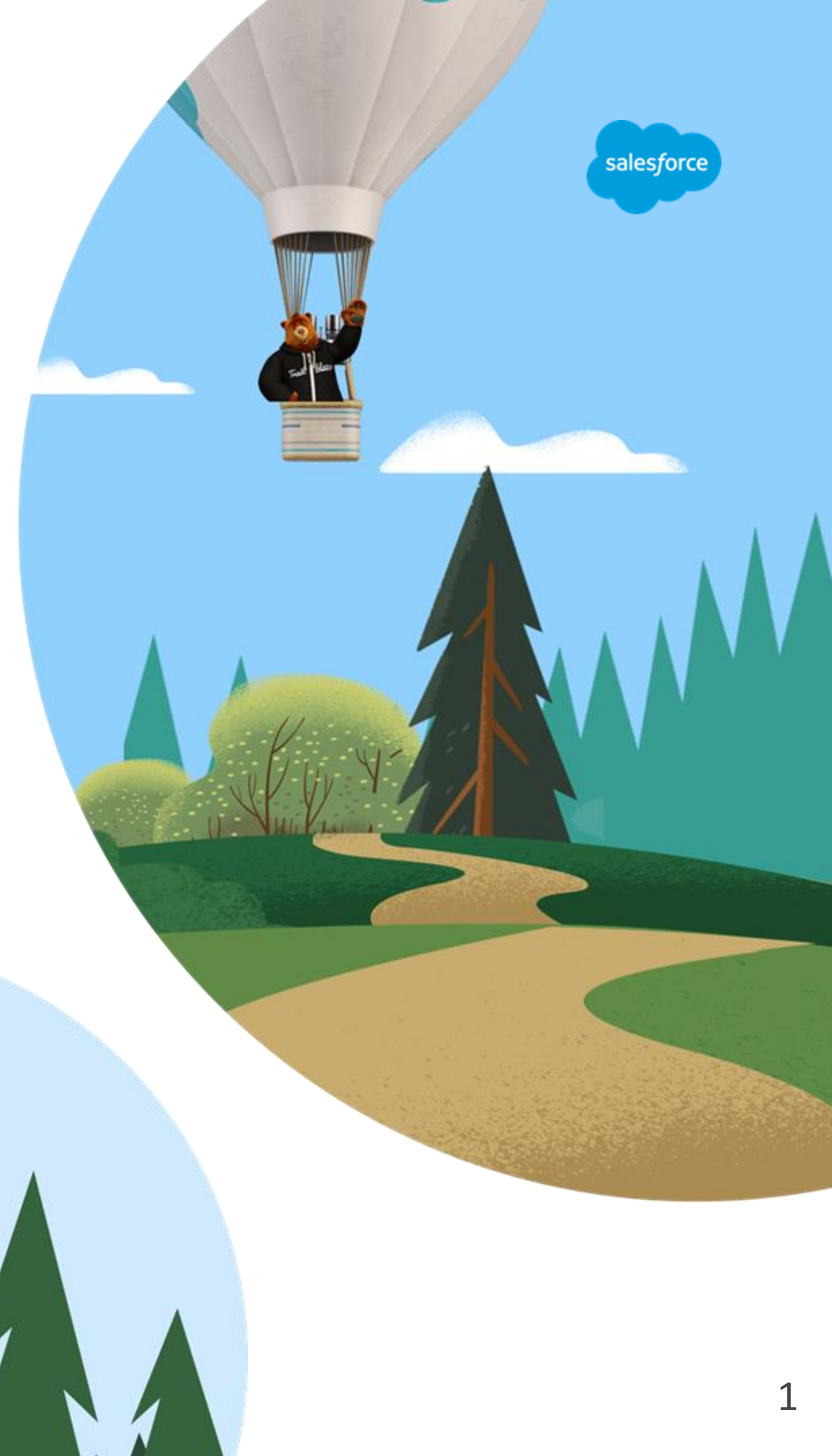# Agenda

Evaluation and Benchmark

Parametric Knowledge Adaptation
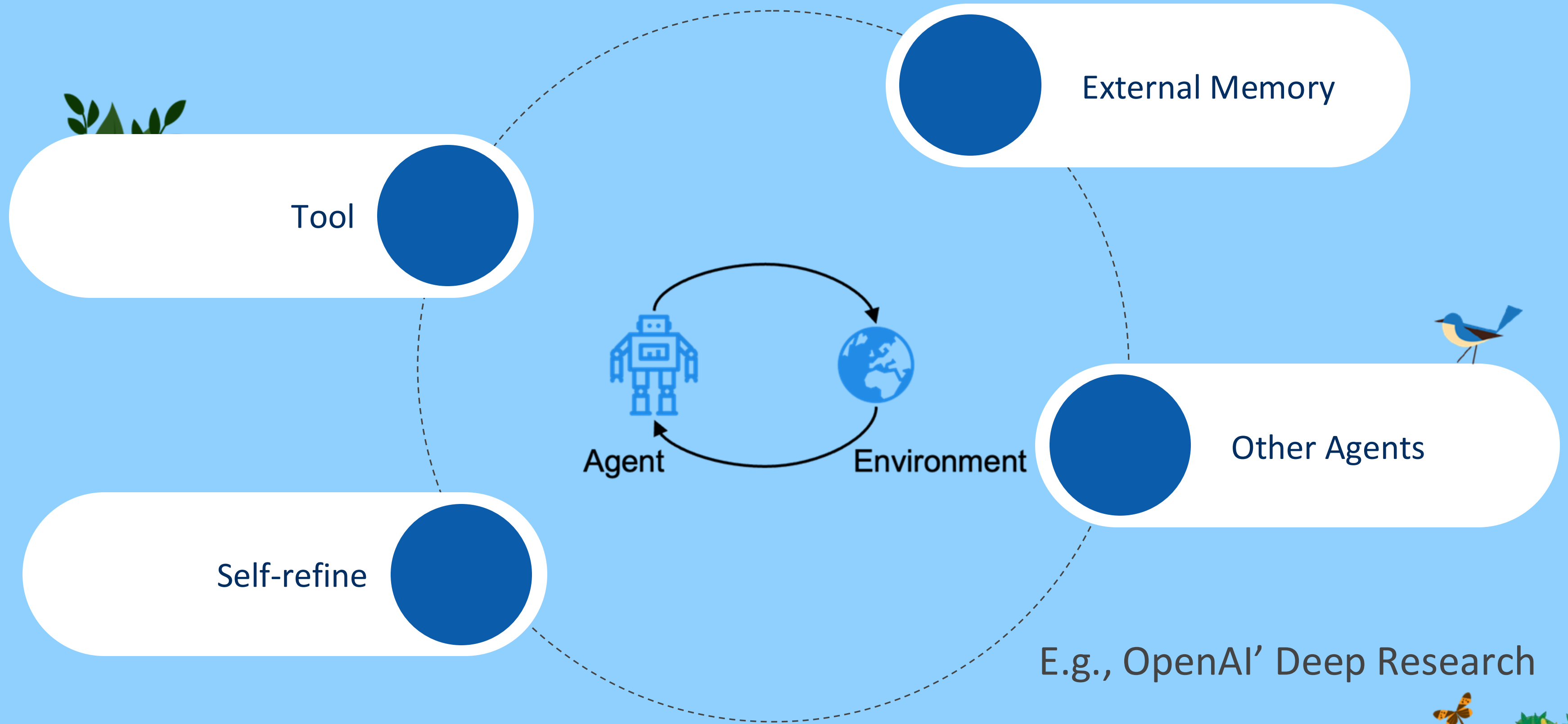
Semi-Parametric Knowledge Adaptation ~30min

Summary, Discussion, QAs

# Semi-Parametric Knowledge



Tool

External Memory

Self-refine

Other Agents

Agent

Environment

E.g., OpenAI' Deep Research

# RAG – Role

## Bridge Gap

Off-the-shelf LLMs may not have been optimized for leveraging external information in its context

Additional adaptation is required for better performance

## Autonomous Decision Making

A RAG system needs to decide whether it needs external information or it can respond directly
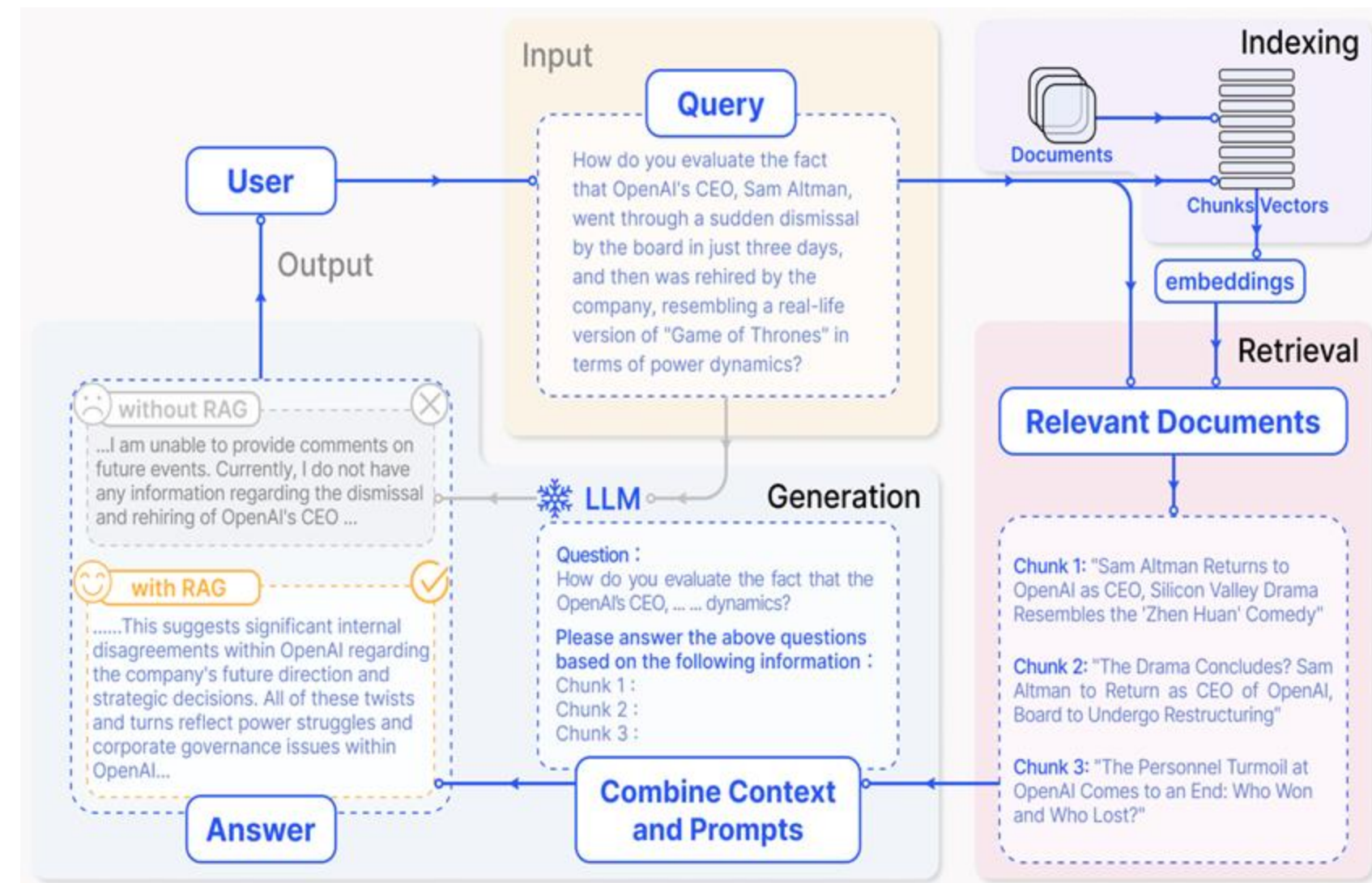
It may need to ask for clarification to the user, do multiple searches via retrieval and aggregate results across documents

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RAG - Key Ideas

## Example Workflow

### Three Main Components

- LLM
- Retriever
- LLM-Retriever Interaction



## Minimalist RAG System

# RAG – Key Considerations

## Training Recipe

**Data Recipe:**
- Hard to obtain ground truth decision-making trajectory data.
- Model should be robust to potentially noisy context.

**Model Recipe:**

**Algorithm**: How to optimize the LLM for search-based interactions?

**Training Workflow**: What kind of workflow we should use?

## Seed Data

**Data Source:** Where to get the data?

**Data Mixture:** What should be included in the RAG data?

**Data Budget:** How much data we need?

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RAG – Key Ideas

## LLM and Decision Making

**Post-train LLMs for contextual usage**

Deal with:

- Noisy context (passages from same document and different documents)
- Conflicting evidence
- Counterfactual evidence
- Absence of knowledge

E.g., SFR-RAG (Salesforce), RAG 2.0 (Contextual AI)

**LLMs with agentic workflow**

- Predefined or autonomous workflow.
- Single agent vs. multi-agent system
- Planner and worker agents

E.g., Infogent, Manus Agent, Deep Research (OpenAI)

INFOGENT: An Agent-Based Framework for Web Information Aggregation, Reddy, et al., 2024
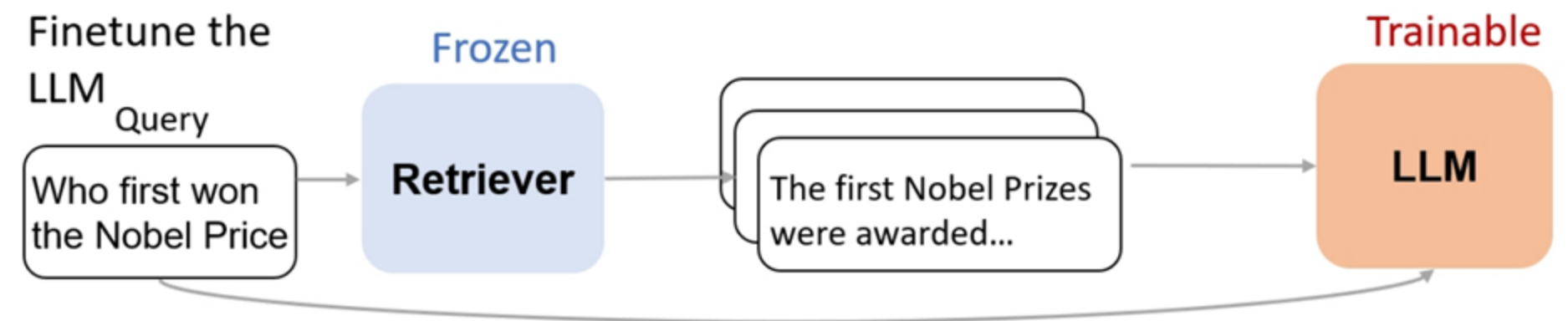
# RAG – Key Ideas

## Train LLMs for Contextual Use

**Post-train LLMs for RAG scenarios:**

Create contextual fine-tuning data to deal with noisy contexts, counterfactual contexts, no-answer contexts and conflicting

Examples: SFR-RAG, RAG 2.0



1. **Fix the retriever**
2. **Train the LLM for contextual usage**

SFR-RAG: Towards Contextually Faithful LLMs, Nguyen et al., 2024
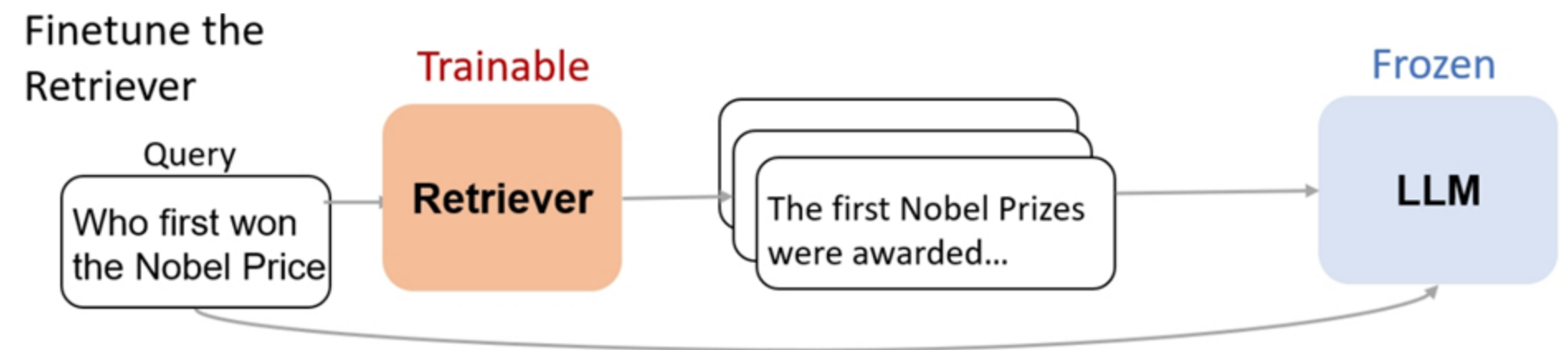RAG2.0: https://contextual.ai/introducing-rag2/

# RAG – Key Ideas

## Align Retriever to LLM

The output of a frozen LLM is used as supervision signals to train the retriver

Examples: REPLUG, Atlas

1. **Fix the LLM**
2. **Align the retriever to LLM**



REPLUG: Retrieval-Augmented Black-Box Language Models, Shi et al., 2023
Atlas: Few-shot Learning with Retrieval Augmented Language Models, Izacard, 2022
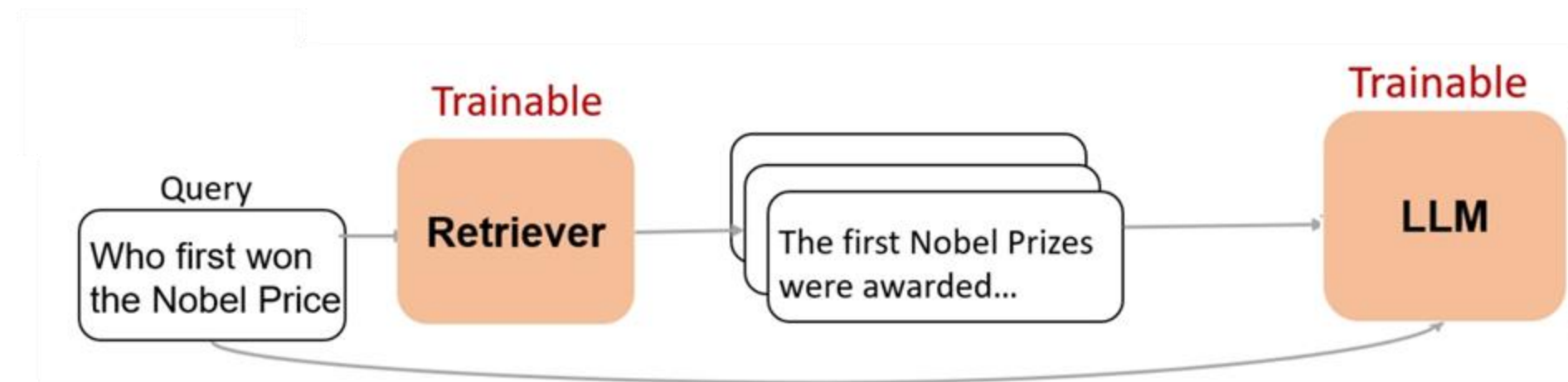
# RAG – Key Ideas

## Train both the LLM and Retriver

Jointly or sequentially train the retriever and LLMs so that they are aligned

Examples: RA-DIT

1. **Train both the LLM and the retriever**



RA-DIT: Retrieval-Augmented Dual Instruction Tuning, Lin et al, 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025
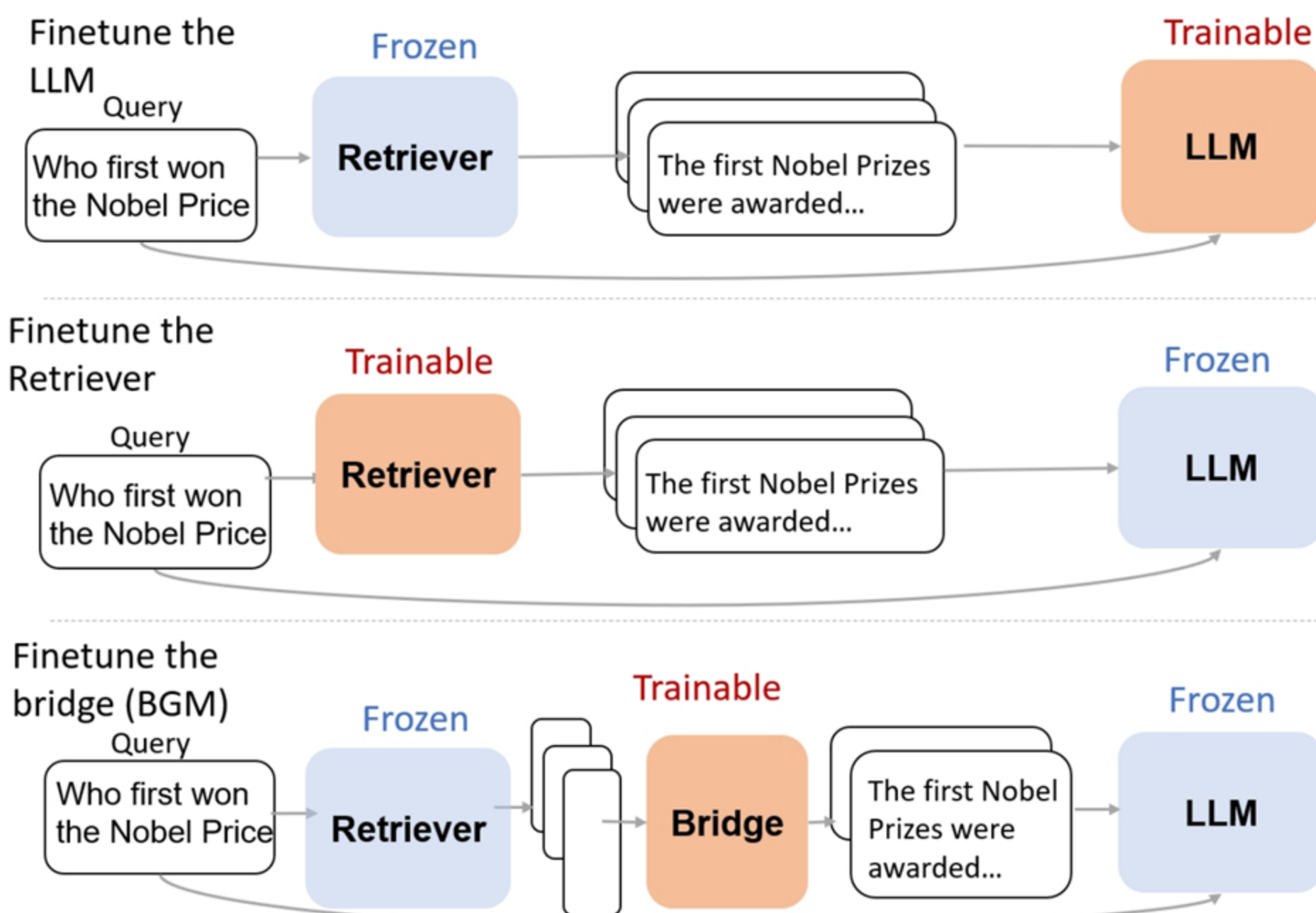
# RAG – Key Ideas

## LLM-Retriever Interaction

**Fix the LLM and Retriver**

**Train a "bridge" (a LLM) to connect their preference**

**Main innovation:** There is preference gap between **retriever** (built for human) and **LLM** (can prefer different order, selection..). One alternative way besides training LLM or retriever is to train an intermediate bridge
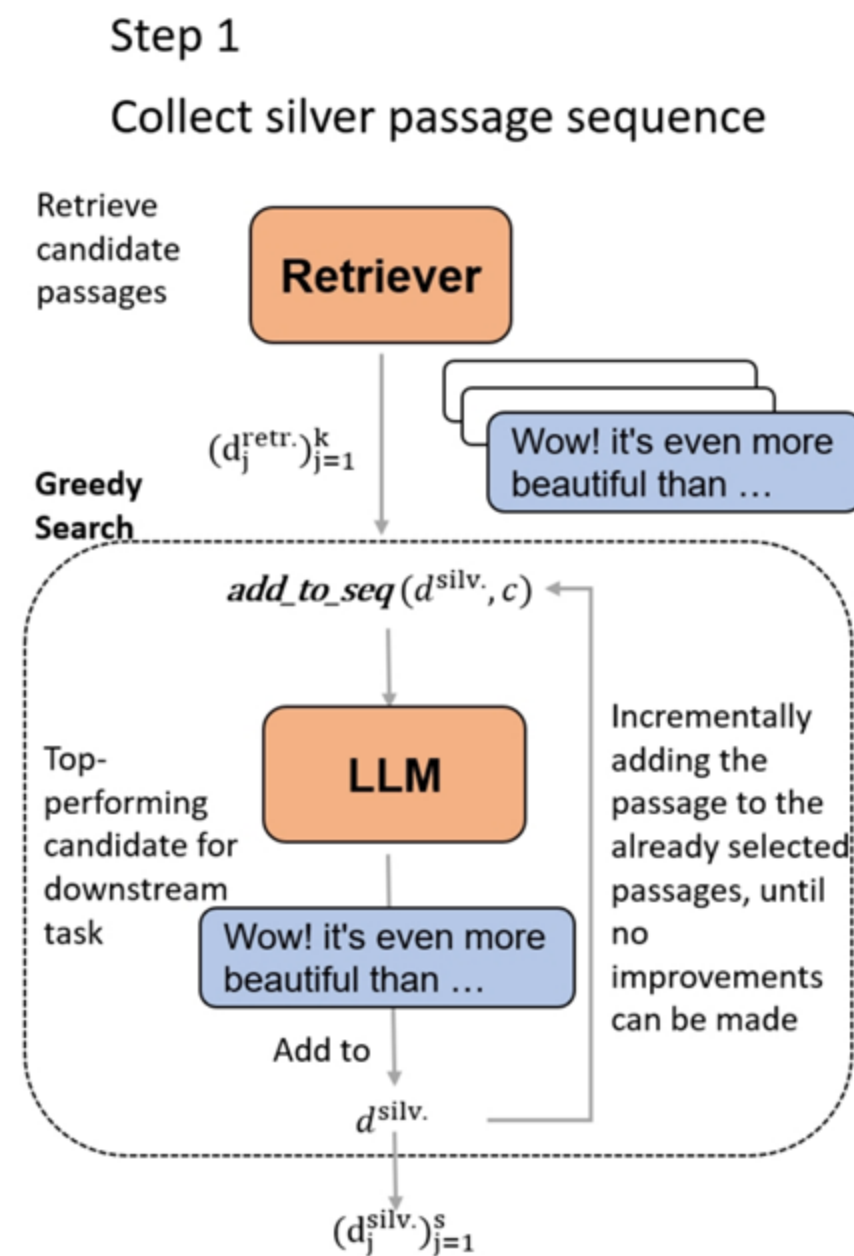


Bridging the Preference Gap between Retrievers and LLMs, Ke et al., 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RAG – Key Ideas

## LLM-Retriever Interaction

**Ground Truth Data:** Use greedy search to find the silver passage



Step 1

Collect silver passage sequence

Bridging the Preference Gap between Retrievers and LLMs, Ke et al., 2024

# RAG – Key Ideas

## LLM-Retriever Interaction

**Ground Truth Data:** Use greedy search to find the silver passage

**Workflow:** IT → RL



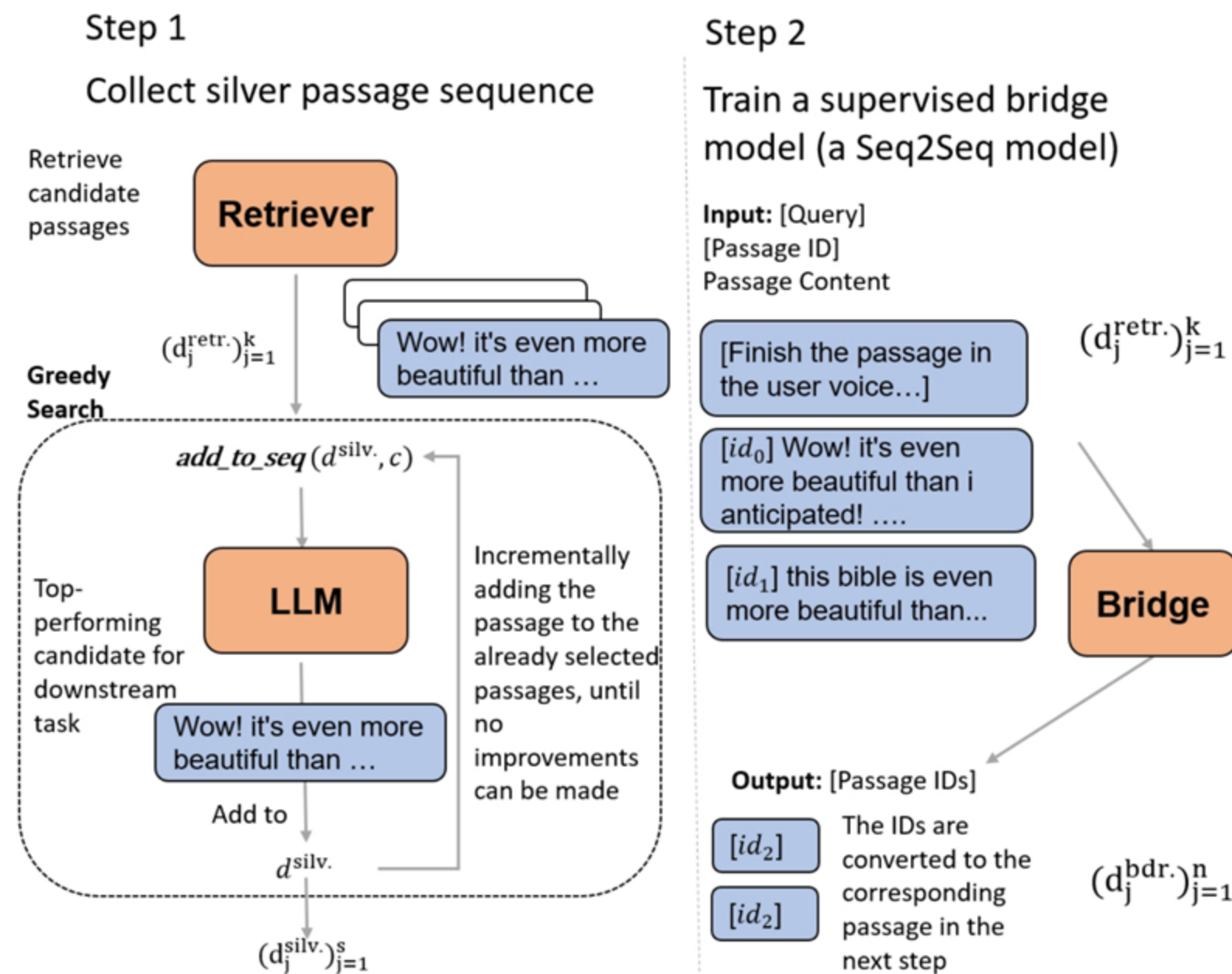Bridging the Preference Gap between Retrievers and LLMs, Ke et al., 2024

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# RAG – Key Ideas

## LLM-Retriever Interaction

**Ground Truth Data:** Use greedy search to find the silver passage

**Workflow:** IT → RL



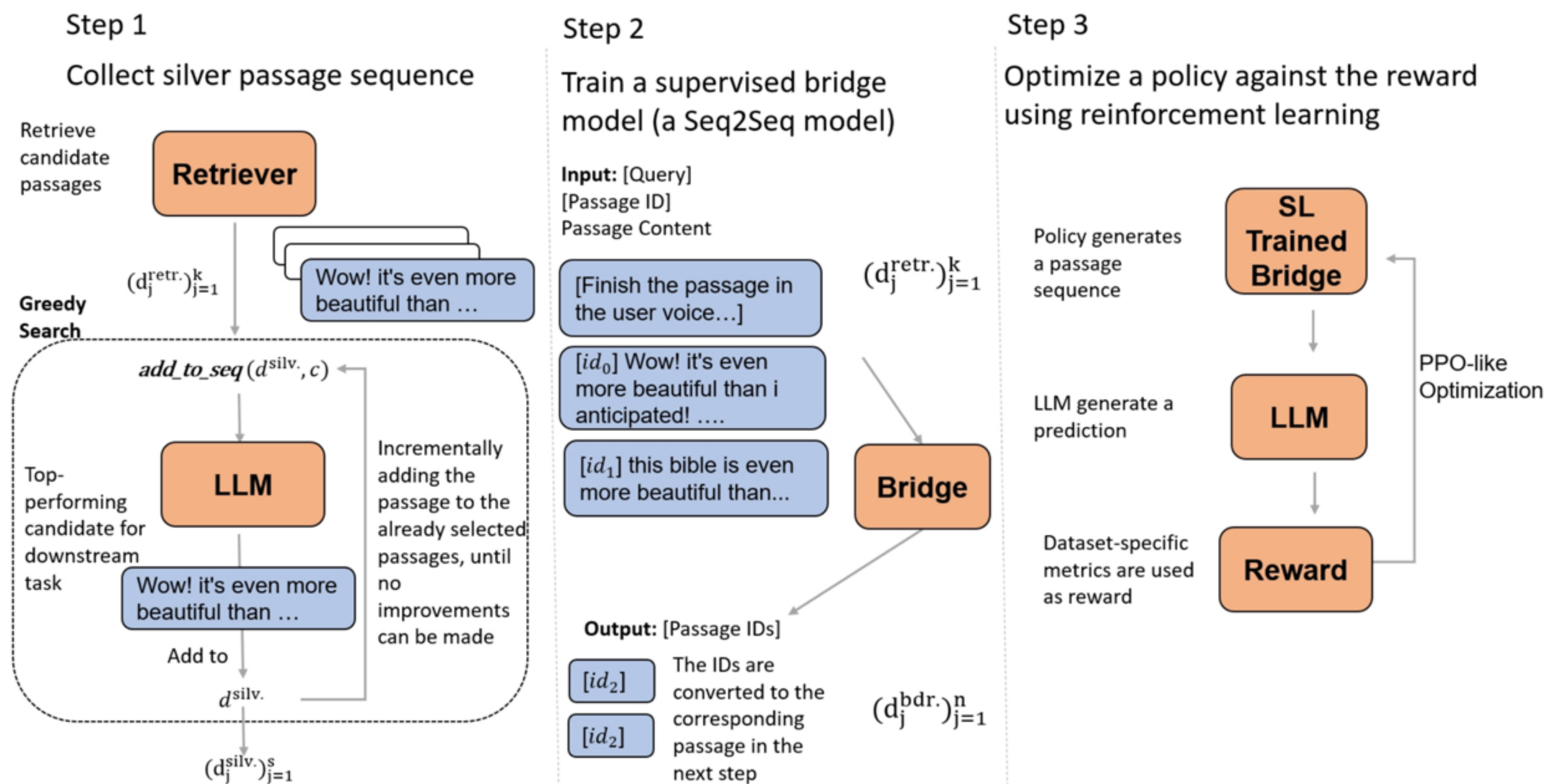Bridging the Preference Gap between Retrievers and LLMs, Ke et al., 2024
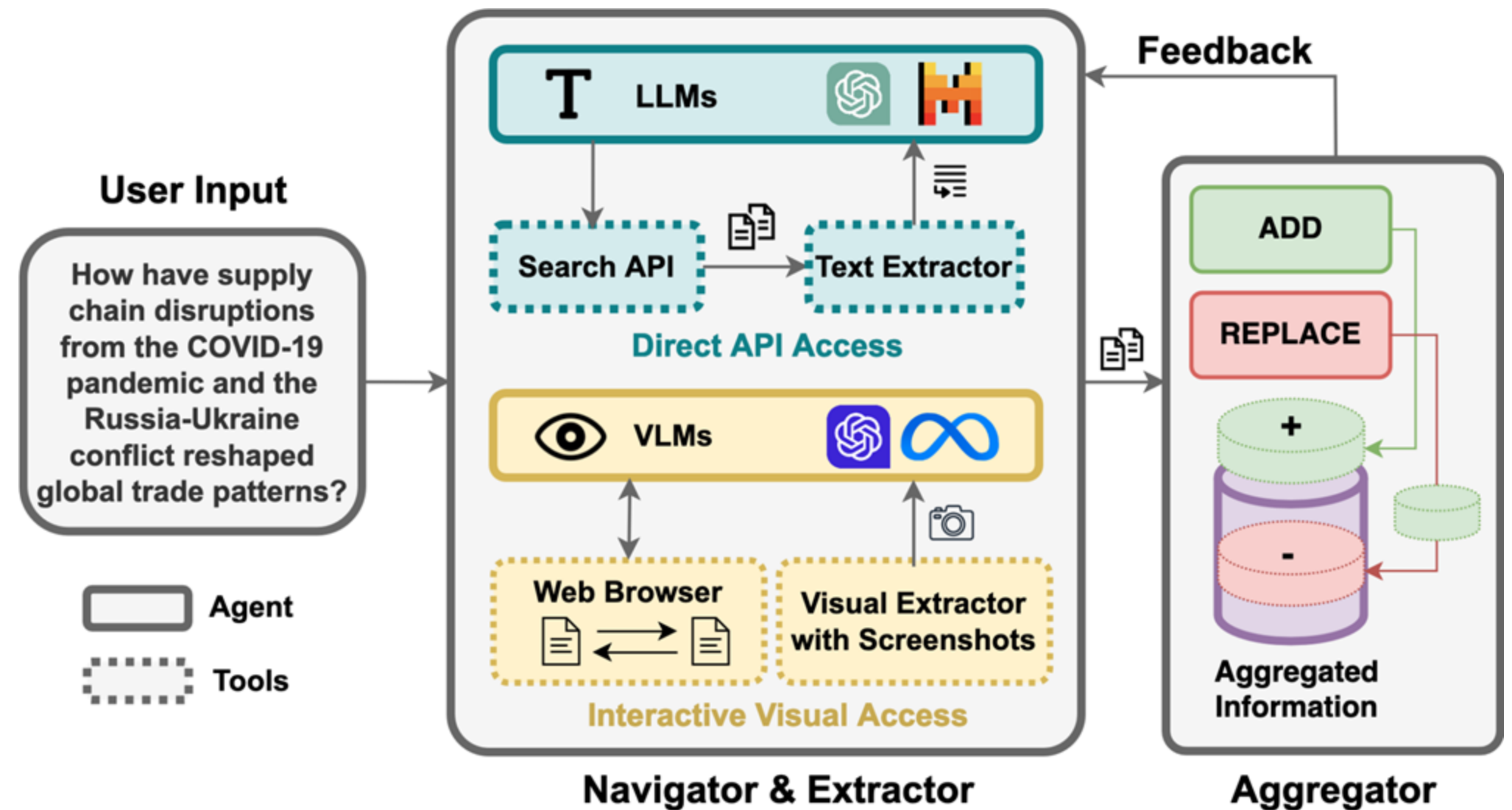
Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025

# Agentic RAG

## RAG with Predefined Workflow

**Main innovation:** RAG can be performed in multiple predefined steps (workflow) to approach the final goal. Those steps usually involve API call, web browser, planner, etc.

Examples: Infogent, MindSearch



INFOGENT: An Agent-Based Framework for Web Information Aggregation, Reddy, et al., 2024
MindSearch: Mimicking Human Minds Elicits Deep AI Searcher, Chen et al., 2024

# RAG – Key Ideas Summary

## Training Recipe

**Data Recipe:**
   often use heuristic way to construct the ground truth

**Model Recipe:**
   **Algorithm and Workflow**: so far, it is largely follows the parametric knowledge adaptation

## Seed Data

**Data Source:** Knowledge-extensive tasks

**Data Mixture:** Can be large scale (e.g., Math, Logic, Code, Science, Reasoning..)

**Data Budget:** Follow the budget required in the specific method

Ke, Ming, Joty - Adaptation of LLMs Tutorial, NAACL 2025