

AGENTIC RELEVANCE TUNING

STAVROS MACRAKIS, OPENSEARCH @ AWS
DANIEL WRIGLEY, OPENSOURCE CONNECTIONS



Search is perfect!

- no CEO nowhere nohow

Relevance

- Everyone wants it
- No one knows how to achieve it
- It's too much work

More specifically

Customers can't find iPhones

Revenue is down

Boss says "just use AI"

THE VIRTUOUS CIRCLE OF SEARCH IMPROVEMENT

Search processing

- Lexical search
- Semantic search
- Sparse neural

Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

Search tuning

- Manual tuning
- Hybrid search
- Semantic reranking
- Learning to Rank (LTR)

LOTS OF PROGRESS

Search processing

- Lexical search
- Semantic search
- Sparse neural

Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

Search tuning

- Manual tuning
- Hybrid search
- Semantic reranking
- Learning to Rank (LTR)

LOTS OF PROGRESS

Search processing

- Lexical search
- Semantic search
- Sparse neural

- Embeddings
- Fast aNN
- Compression
- Caching strategies
- GPU Acceleration

- Faster lexical search

Search tuning

- Manual tuning
- Hybrid search
- Semantic reranking
- Learning to Rank (LTR)
- LLM-based query understanding and rewriting
- LLM-based result synthesis (RAG)

Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

SOME PROGRESS

Search processing

- Lexical search
- Semantic search
- Sparse neural

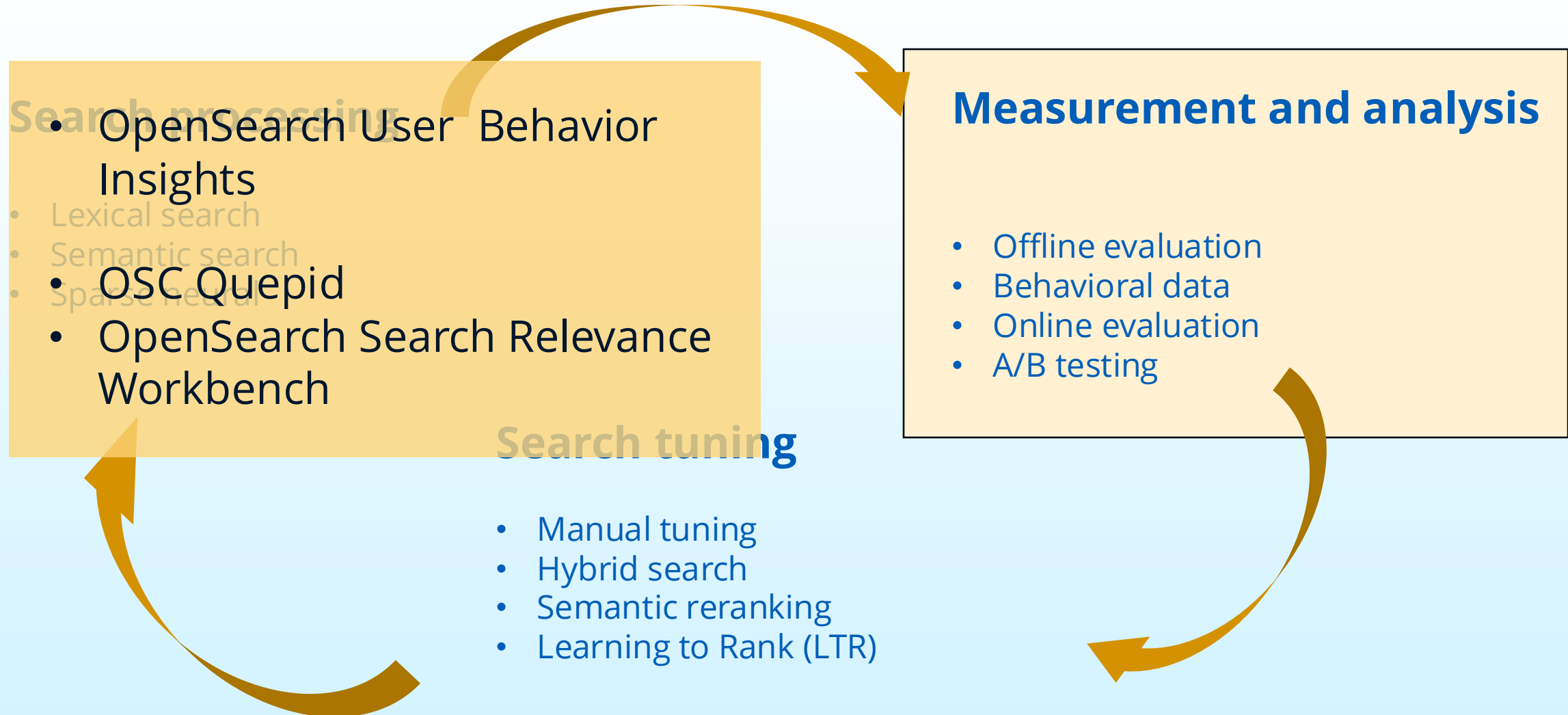
Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

Search tuning

- Manual tuning
- Hybrid search
- Semantic reranking
- Learning to Rank (LTR)

SOME PROGRESS



NOT MUCH PROGRESS

Search processing

- Lexical search
- Semantic search
- Sparse neural

Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

Search tuning

- Manual tuning
 - Field boosts
 - Lexical tricks
 - Hybrid search
- Re-ranking
 - Semantic reranking
 - Learning to Rank (LTR)
- LLM presentation (RAG)

NOT MUCH PROGRESS

Search processing

- Reranking for precision
 - Learning to Rank
 - Cross-encoders

- Lexical search
- Semantic search
- Sparse neural

Doesn't address recall

Measurement and analysis

- Offline evaluation
- Behavioral data
- Online evaluation
- A/B testing

Search tuning

- Manual tuning
 - Field boosts
 - Lexical tricks
 - Hybrid search
- Re-ranking
 - Semantic reranking
 - Learning to Rank (LTR)
- LLM presentation (RAG)

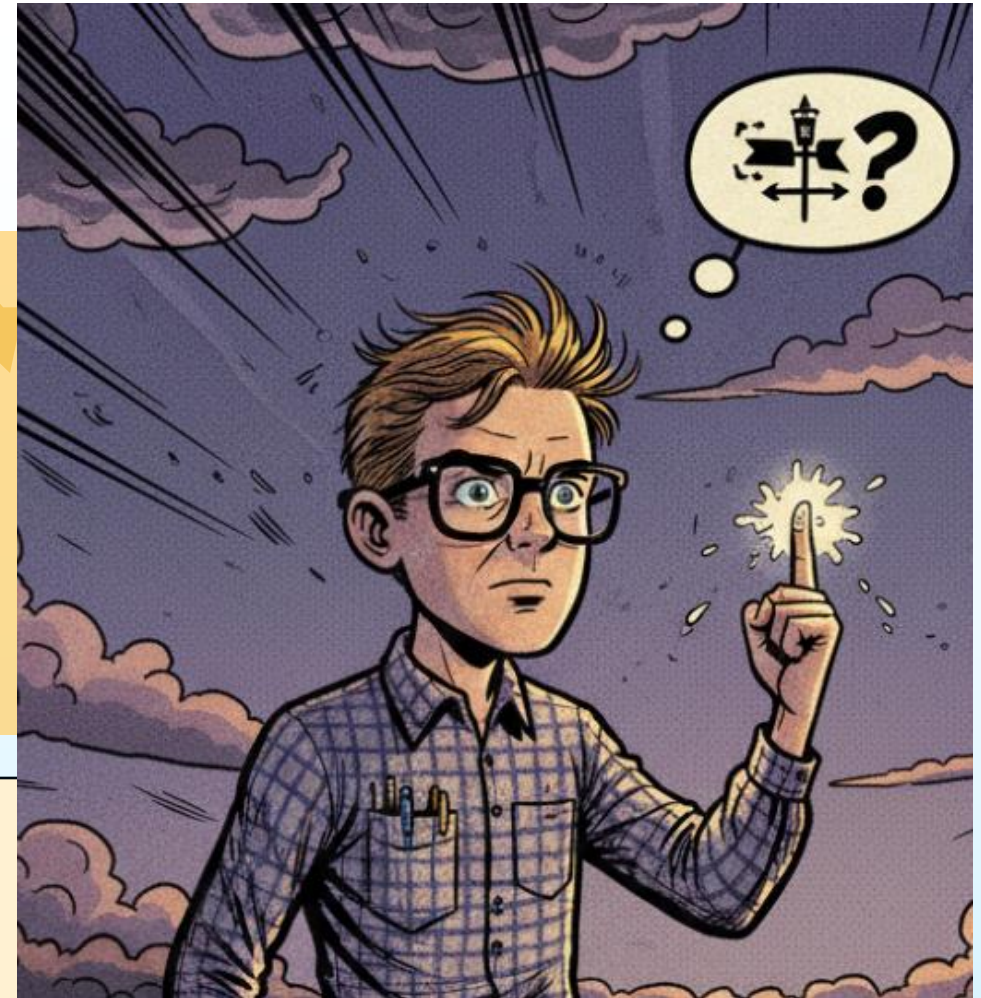
NOT MUCH PROGRESS

Search processing

- Reranking for precision
 - Learning to Rank
 - Cross-encoders
- Lexical search
- Semantic search
- Sparse neural
- Finger in the air?

Search tuning

- Manual tuning
 - Field boosts
 - Lexical tricks
 - Hybrid search
- Re-ranking
 - Semantic reranking
 - Learning to Rank (LTR)
- LLM presentation (RAG)



THE MANUAL PROCESS

Identify a problem

Generate hypotheses

Create search configurations

Run off-line A/B tests

Run on-line A/B tests (months)

Deploy

Iterate on a scale of months

CAN WE DO BETTER?

THE AUTOMATIC PROCESS

Identify a problem — library of business cases

Generate a hypothesis — library of proposed fixes

Create search configurations — rewrite queries

Run off-line A/B tests — use SRW tools,
with and without judgments

Run on-line interleaved A/B test — with real-time feedback

Deploy



Iterate on a scale of days

Relax must match

Add click data

Add phrase matching

Collect data (UBI)

Improve synonyms

Change field weights

Add ngrams

Add/remove fields

Change field mapping

Rewrite query (boost, filter, ...)

Add lemmatization

Boost exact matches v. keyword

Map to structure of data

Add entity extraction for entity products

(query)

Re-rank results

Semantic

Chunking strategy

fine-tune model; finetune

date/freshness)

Add taxonomy; engage external taxonomist

Retrieval (hybrid; sparse)

Faceted search

Add multimedia data

Personalize results

Stock status, shipping, tax (post filtering)

Use better images

Add entities to corpus

Add query terms to a field on clicked X

Use machine learning to figure out the value of terms

Spell checking

Compounding/decompounding

Security trimming

Add geo

products v. accessories

hyperparameters

Map to structure of data

Change field mapping

Add lemmatization

Add entities to corpus

Re-rank results

Use machine learning to

the value of terms

Semantic

Faceted search

Remove query clauses

Add entity extraction for

(query)

Chunking strategy

Retrieval (hybrid; sparse)

Add geo

Add click data

Rewrite query (boost, fi



A QUICK WORD ON INTERLEAVED A/B TESTING

Classic A/B test:

partition *users* into A group and B group

Problem: each user issues different queries

Very slow statistical significance

Interleaved A/B test:

Combine A and B results for *each user*

Faster statistical significance

Much much faster feedback

Accelerating iteration speed:
ranker evaluation with
debiased interleaving

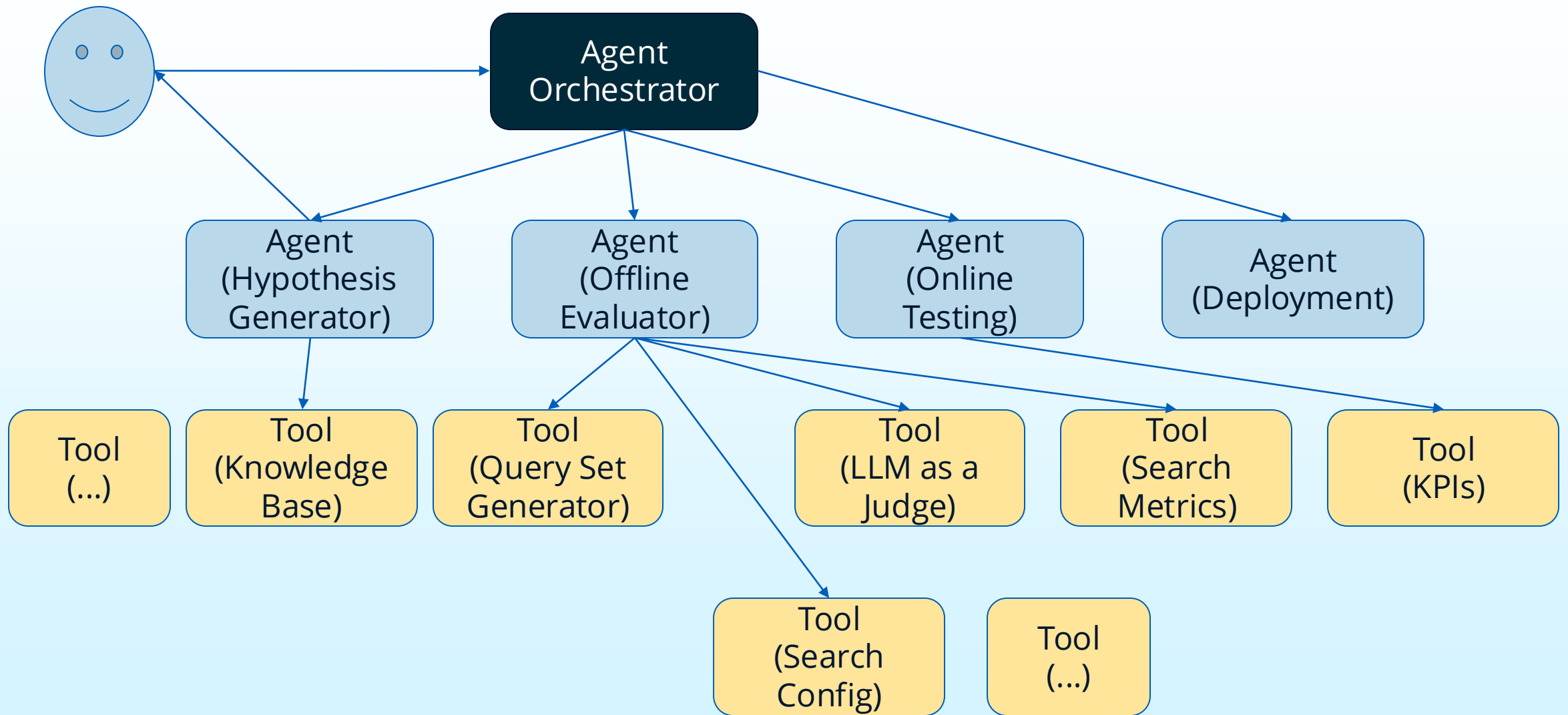
MICES 2025

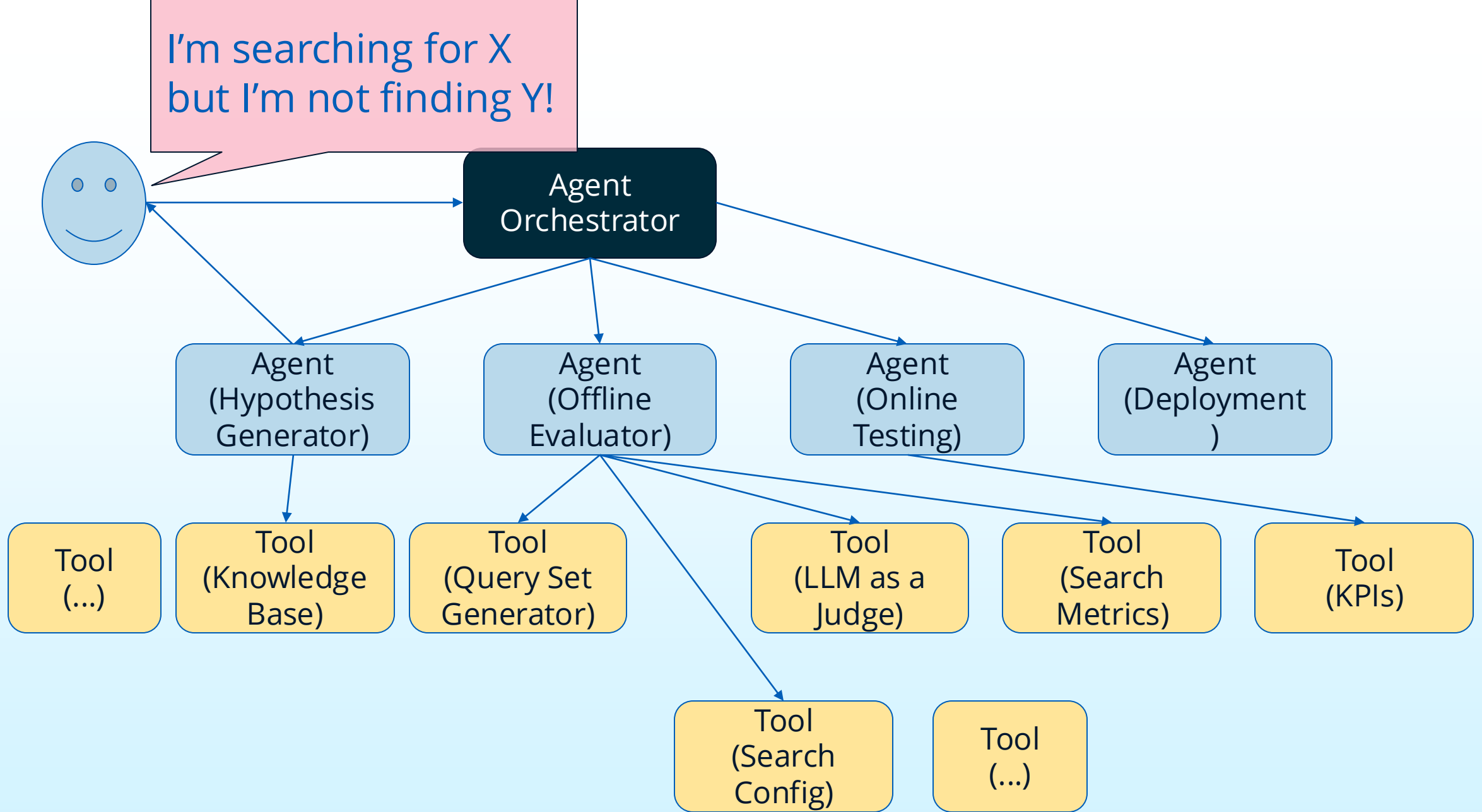
OTTO

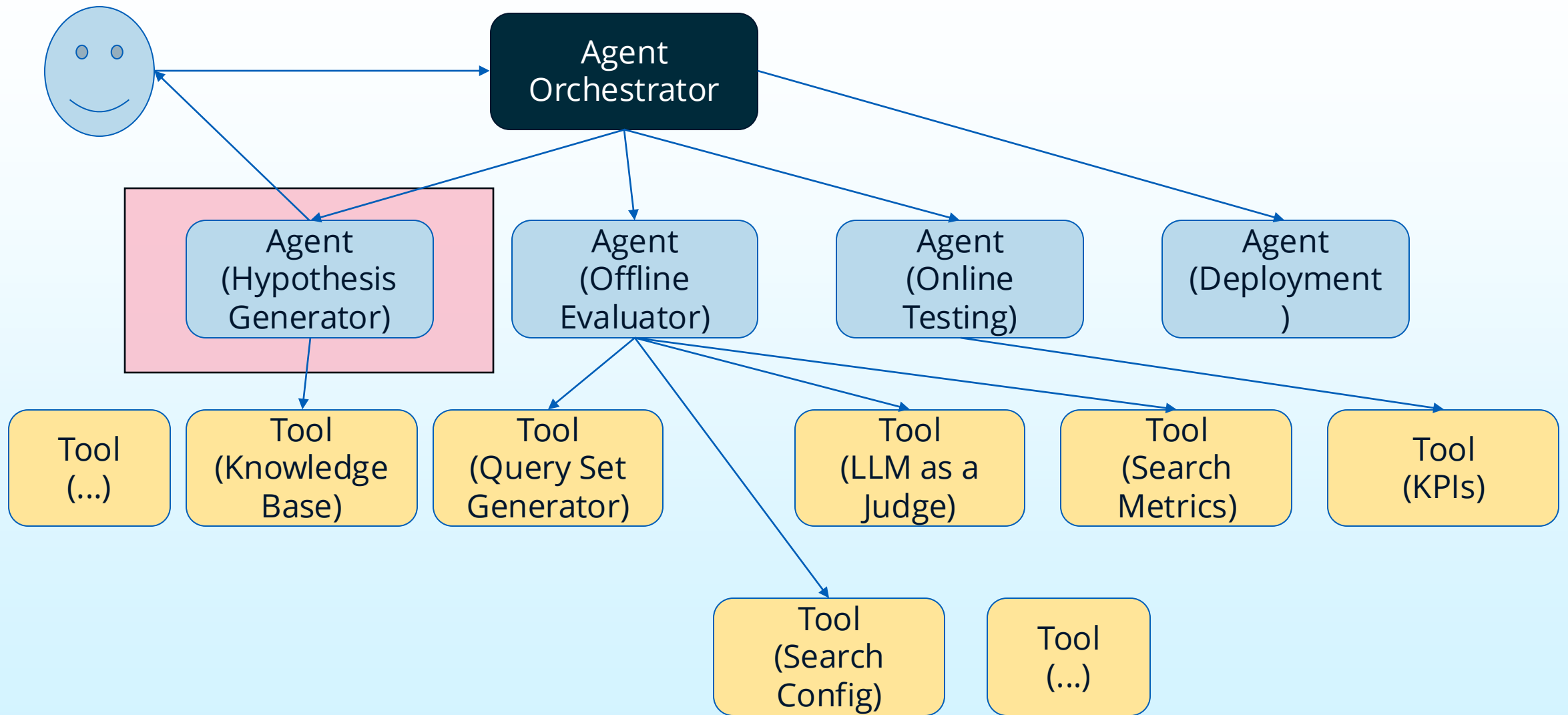
USING AGENTS

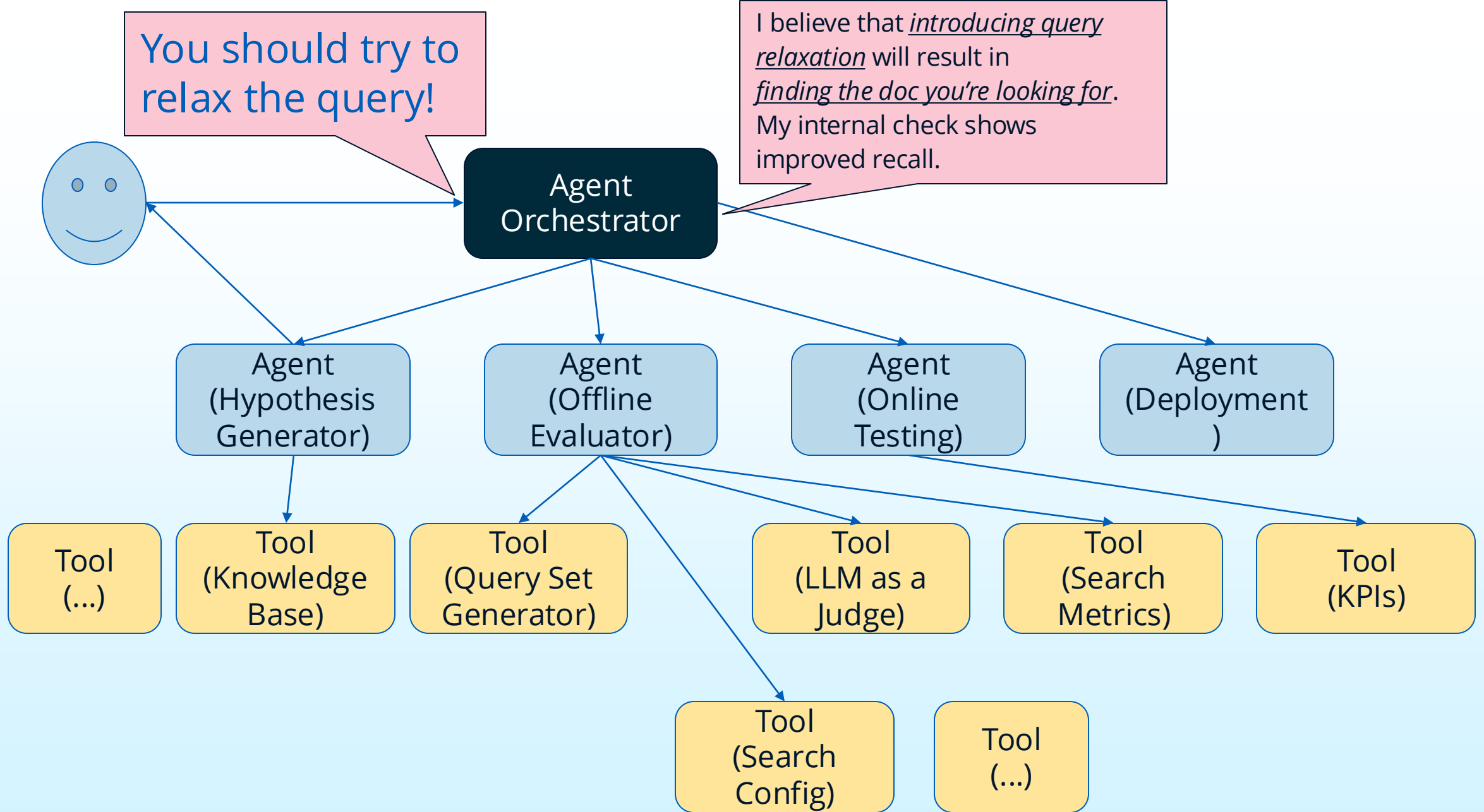
Orchestration and user interaction
Business Problem
Search Strategy
Offline Experiment
Online Experiment runner

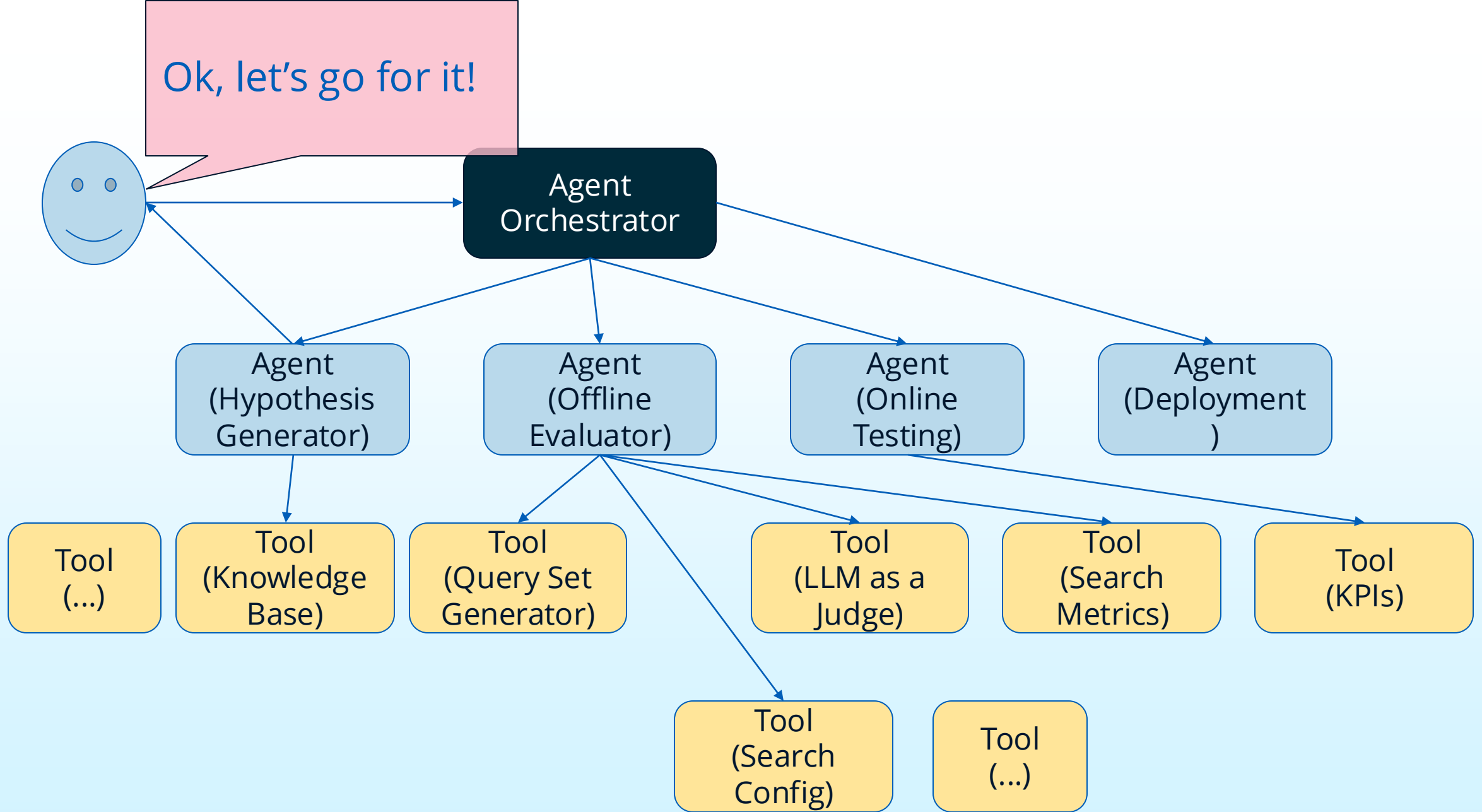
Agents, tools, and modules

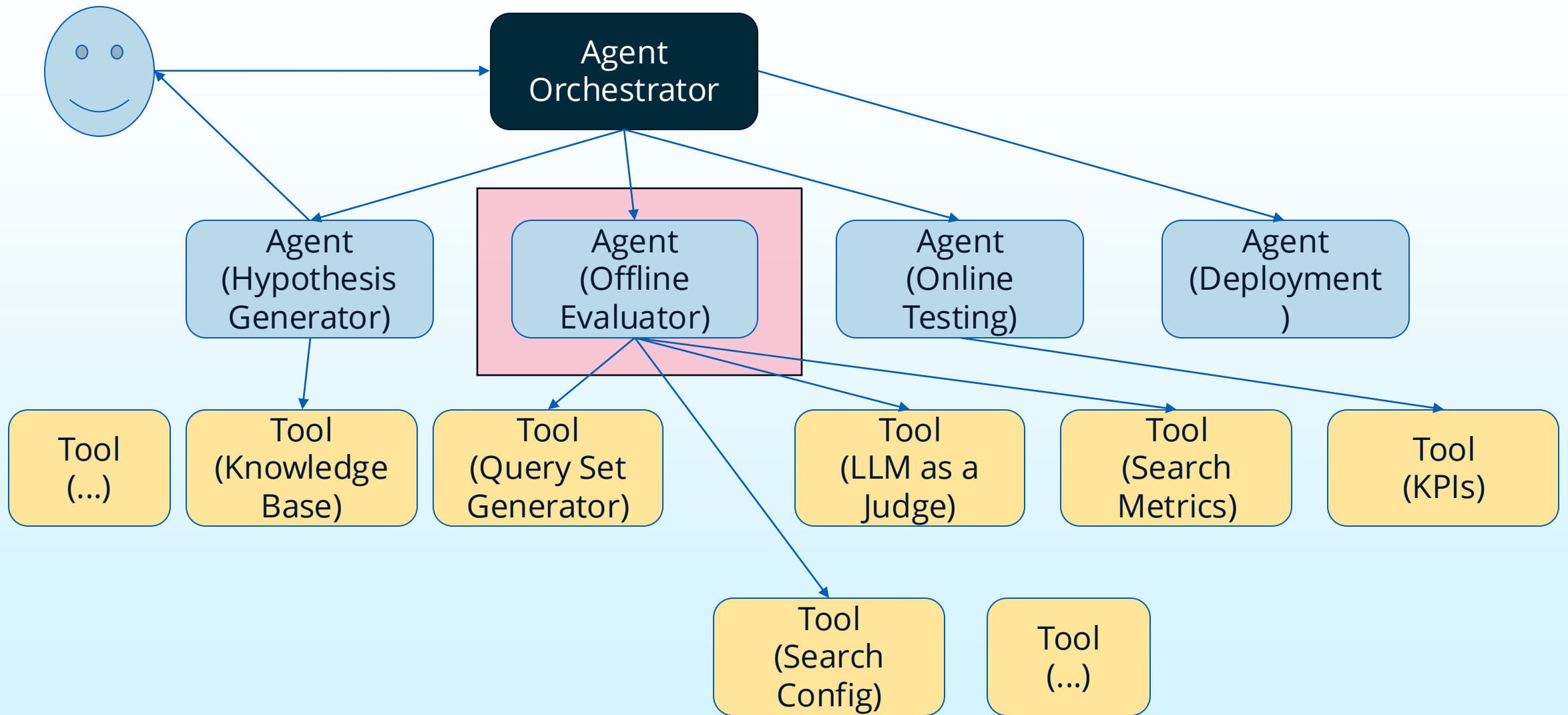




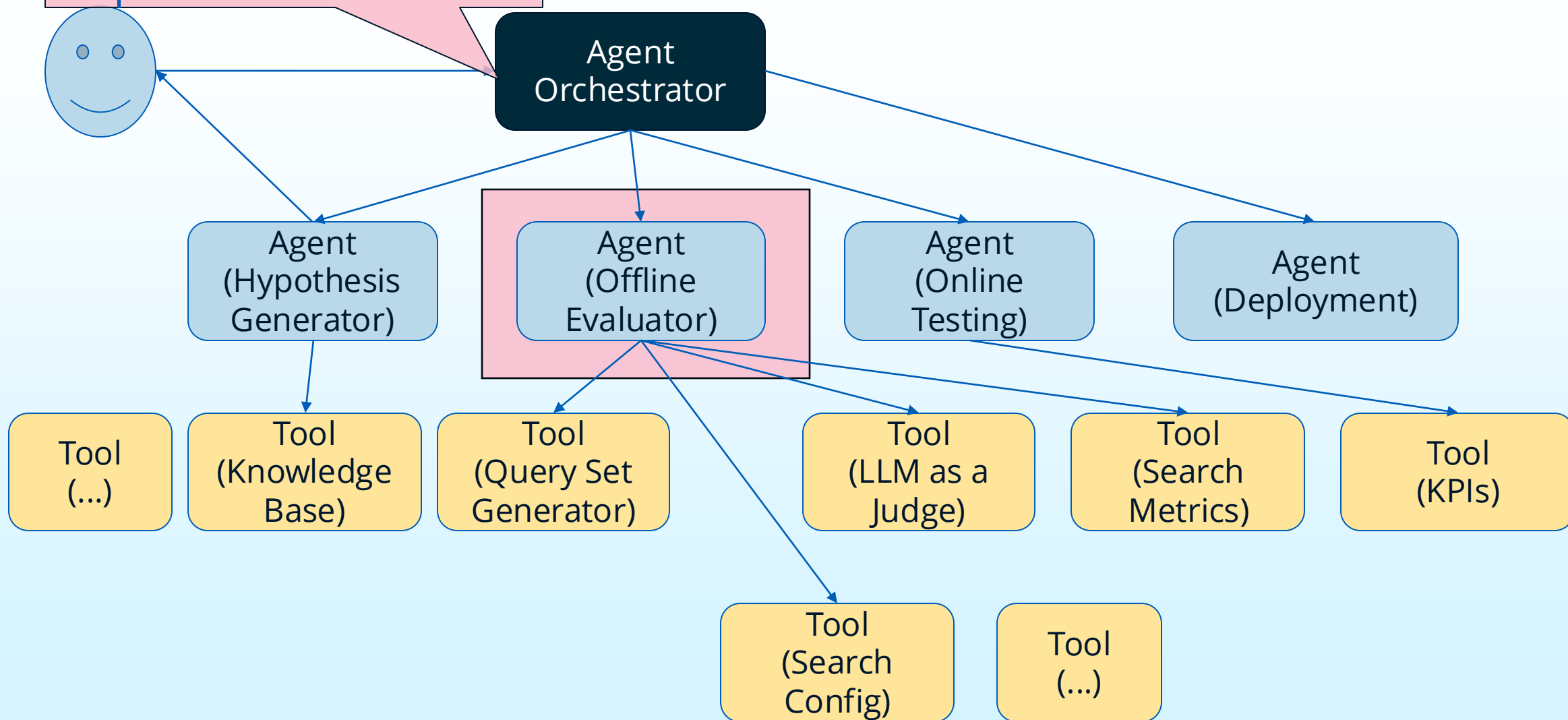


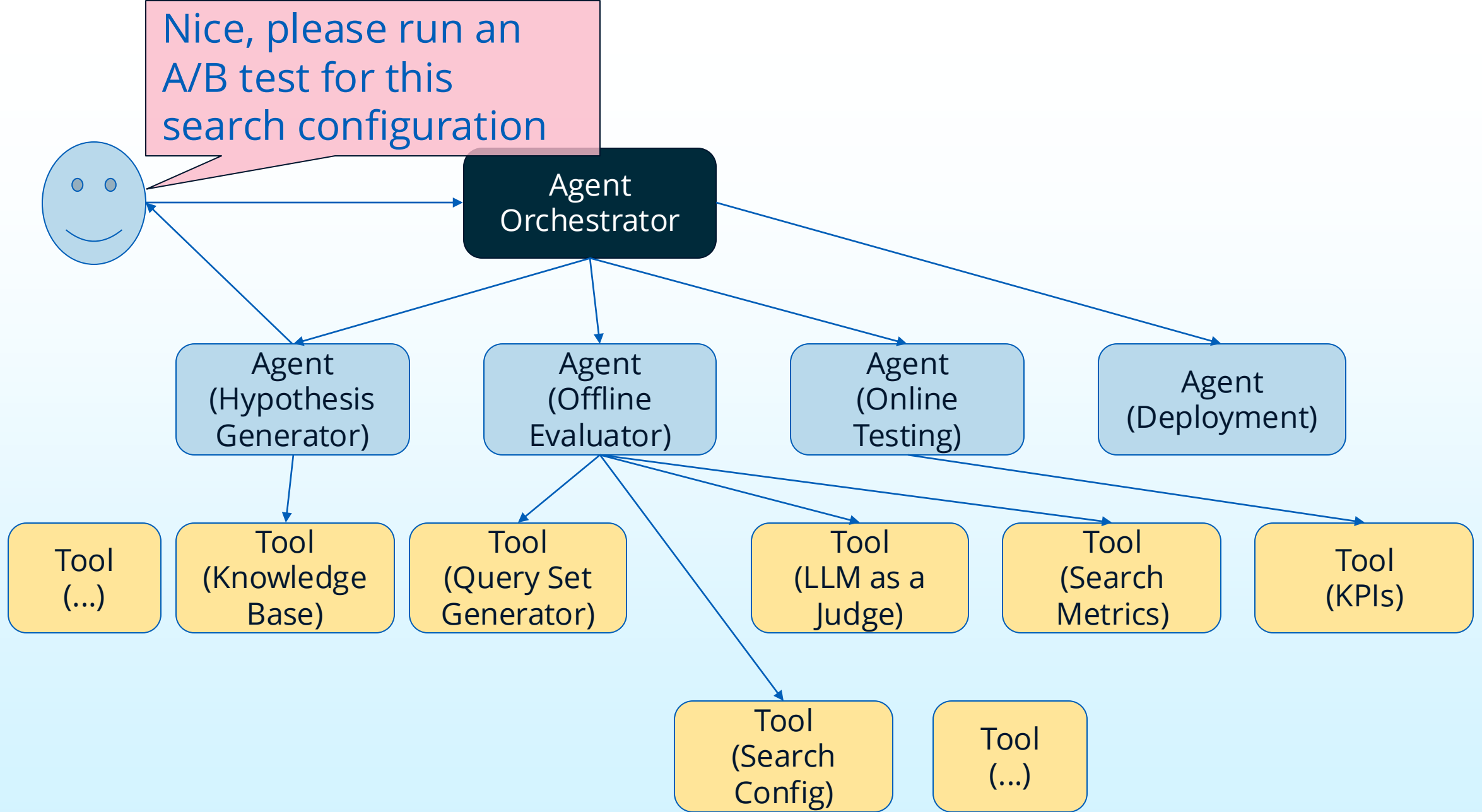


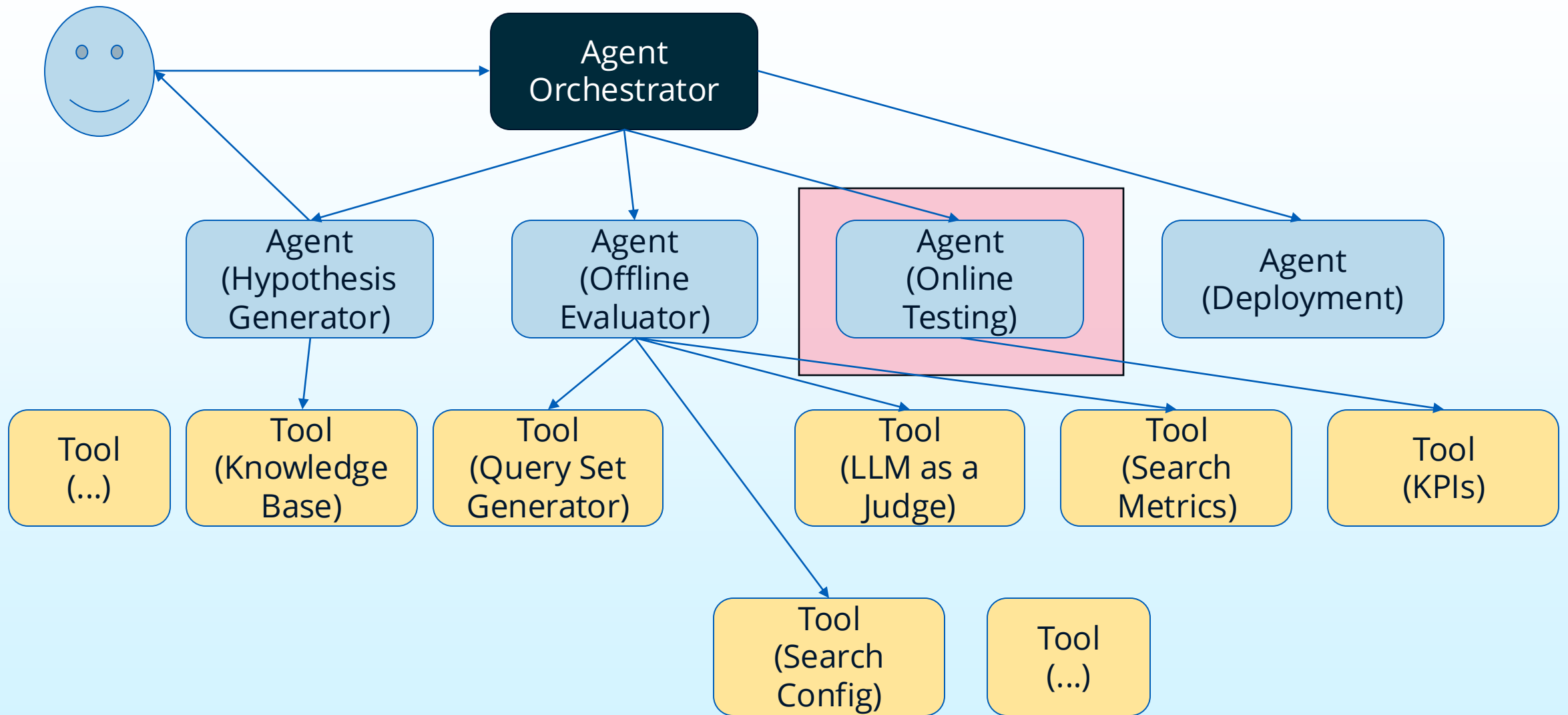




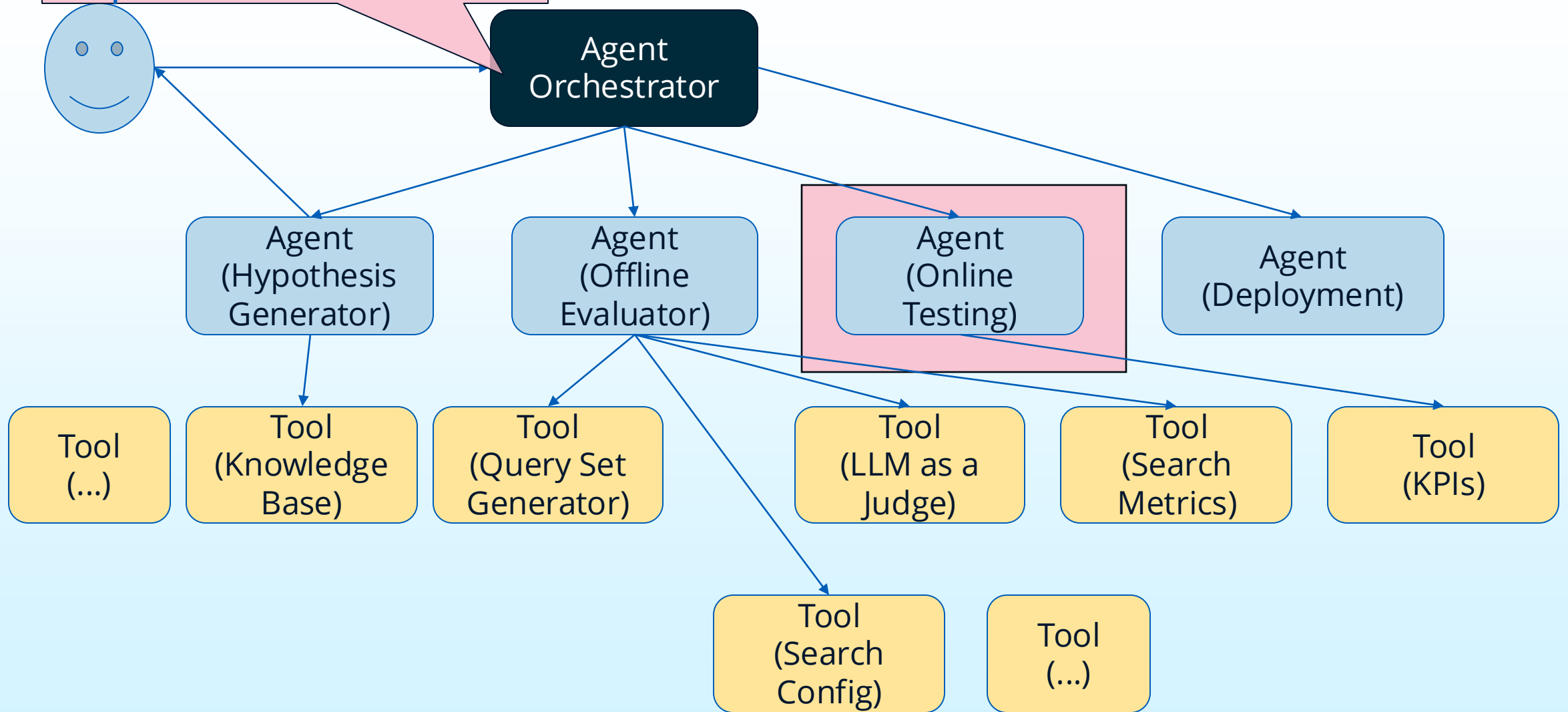
Your search metrics show the following improvements ...

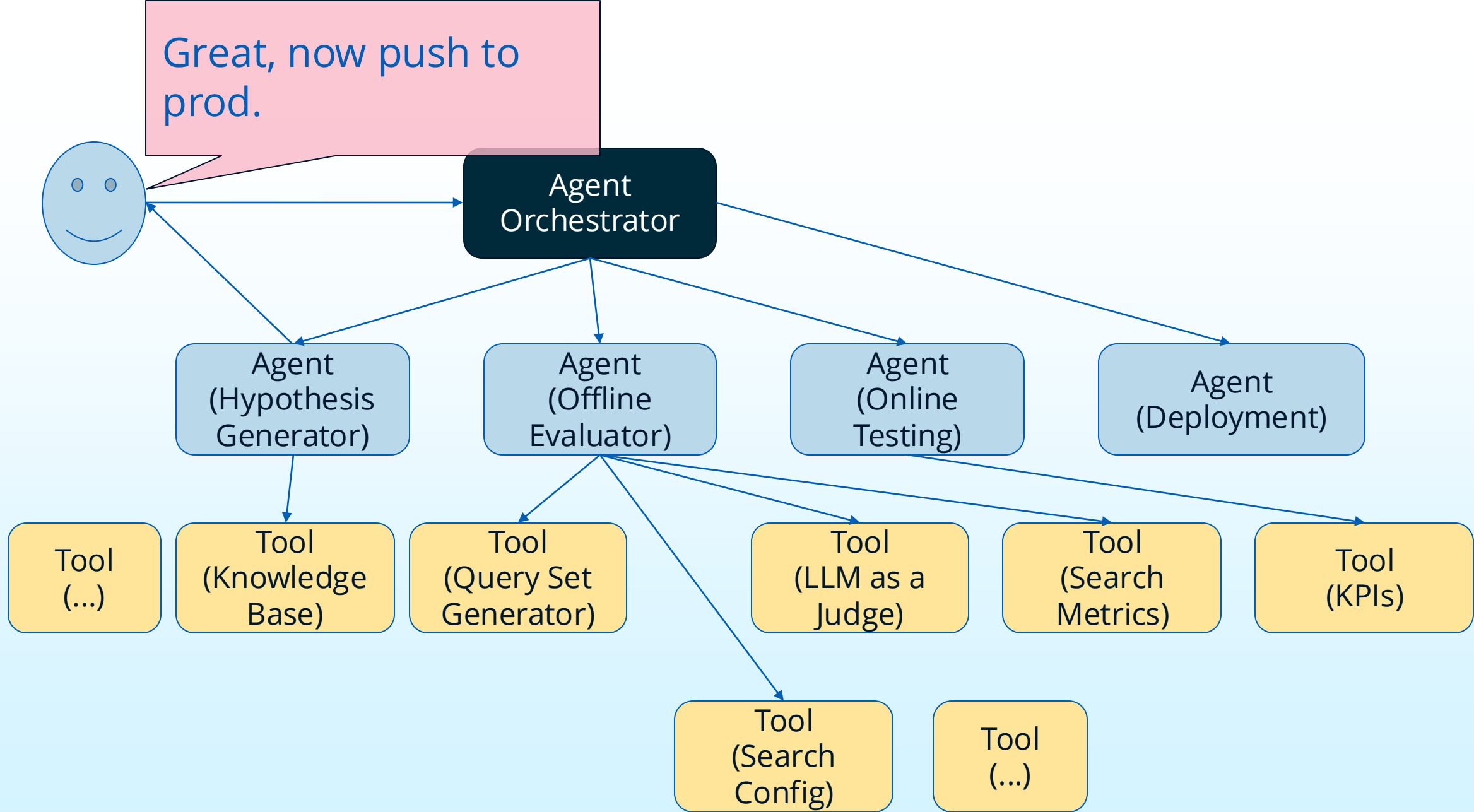




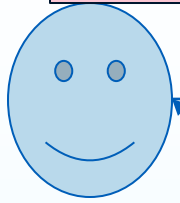


After 38 hours, your
KPIs are significantly
improved ...





Done. Next problem?



Agent
Orchestrator

Agent
(Hypothesis
Generator)

Agent
(Offline
Evaluator)

Agent
(Online
Testing)

Agent
(Deployment)

Tool
(...)

Tool
(Knowledge
Base)

Tool
(Query Set
Generator)

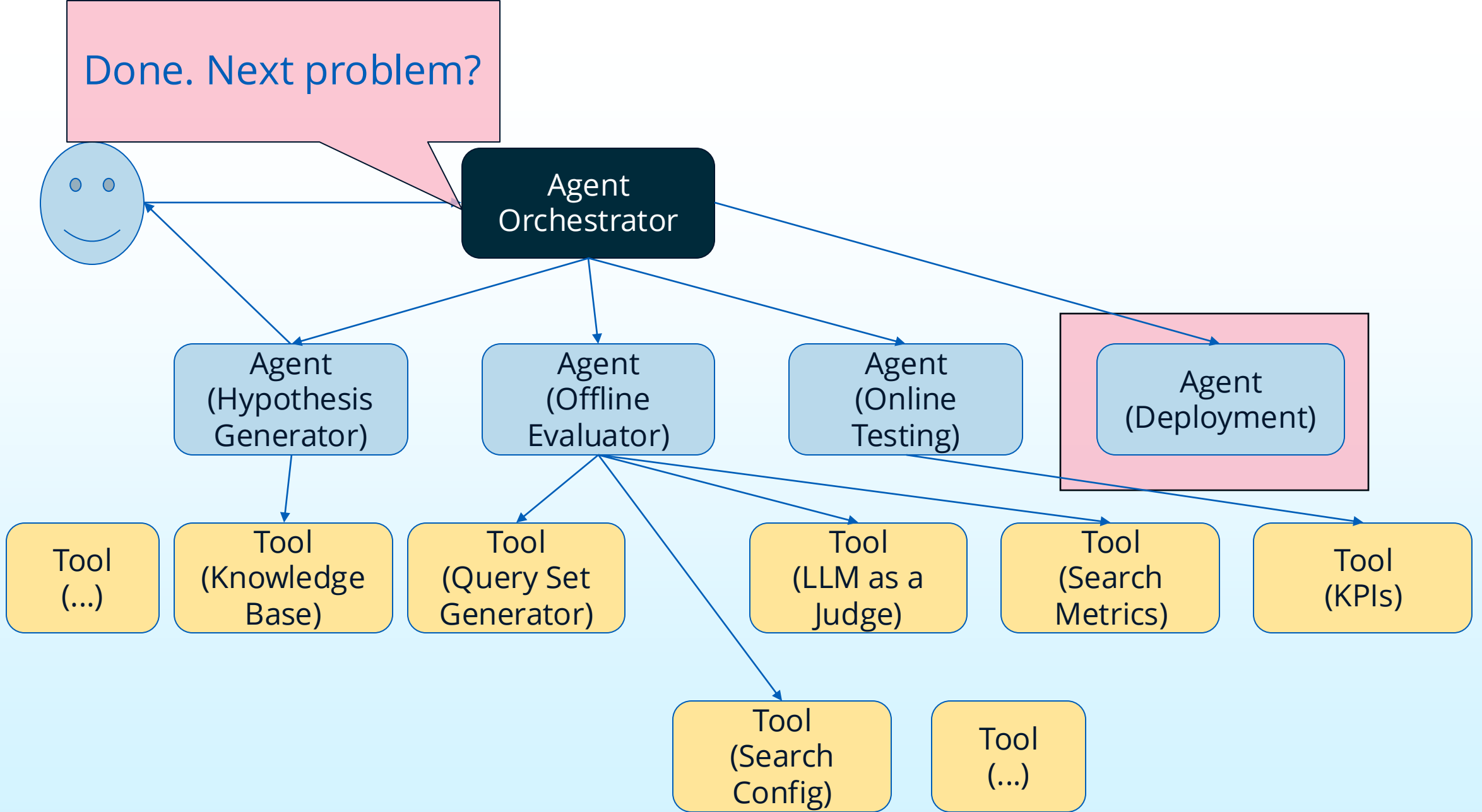
Tool
(LLM as a
Judge)

Tool
(Search
Metrics)

Tool
(KPIs)

Tool
(Search
Config)

Tool
(...)



AGENTIC RELEVANCE TUNING

Agents to automate tuning

+

Interleaving for fast feedback

=

Fast, goal-driven relevance improvements

