# Building RAG Q&A Bots for I/O Psychologists

A step-by-step tutorial

# Introductions

**Yuyun Huang**

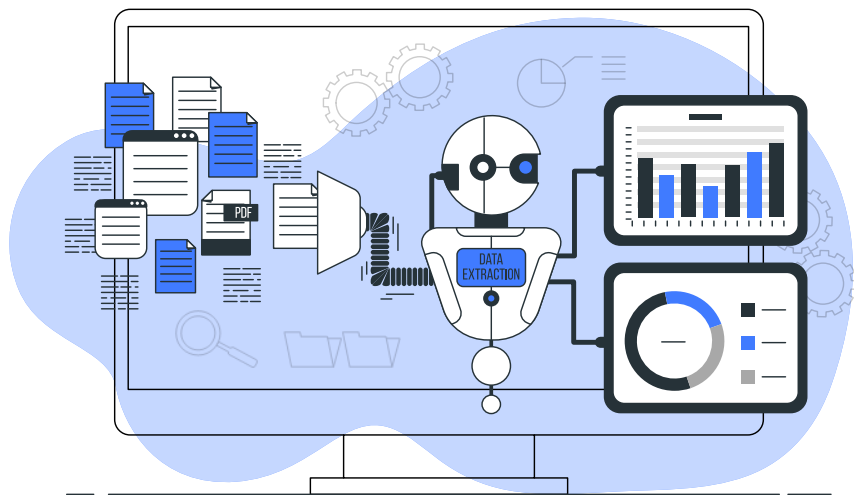*ML / NLP Expert*

**Karim Badr**

*Senior Research Scientist*

**Paul Agnello**

*Senior Associate Consultant*

# Agenda



**The Basics – What is a RAG Chatbot?**

**The Technicals – How to Build your own Chatbot?**

**The Practicals – Hands-On Demos**

# Chatbots for I/O Psychology

- ***Industry/Practitioner:***
  - Employee support and HR assistance
  - Policy and compliance queries
  - You can upload documents (e.g., policies or employee manuals) and ask questions about them.
- ***Academia/Student/Researcher:***
  - Asking questions about research articles
- ***Industry/Practitioner:***
  - Analyze employee surveys
  - Extract sentiments from comments
  - Analyze data
- ***Other:***
  - Real-time Coaching
  - Simulations of interactions with clients/employees/etc.
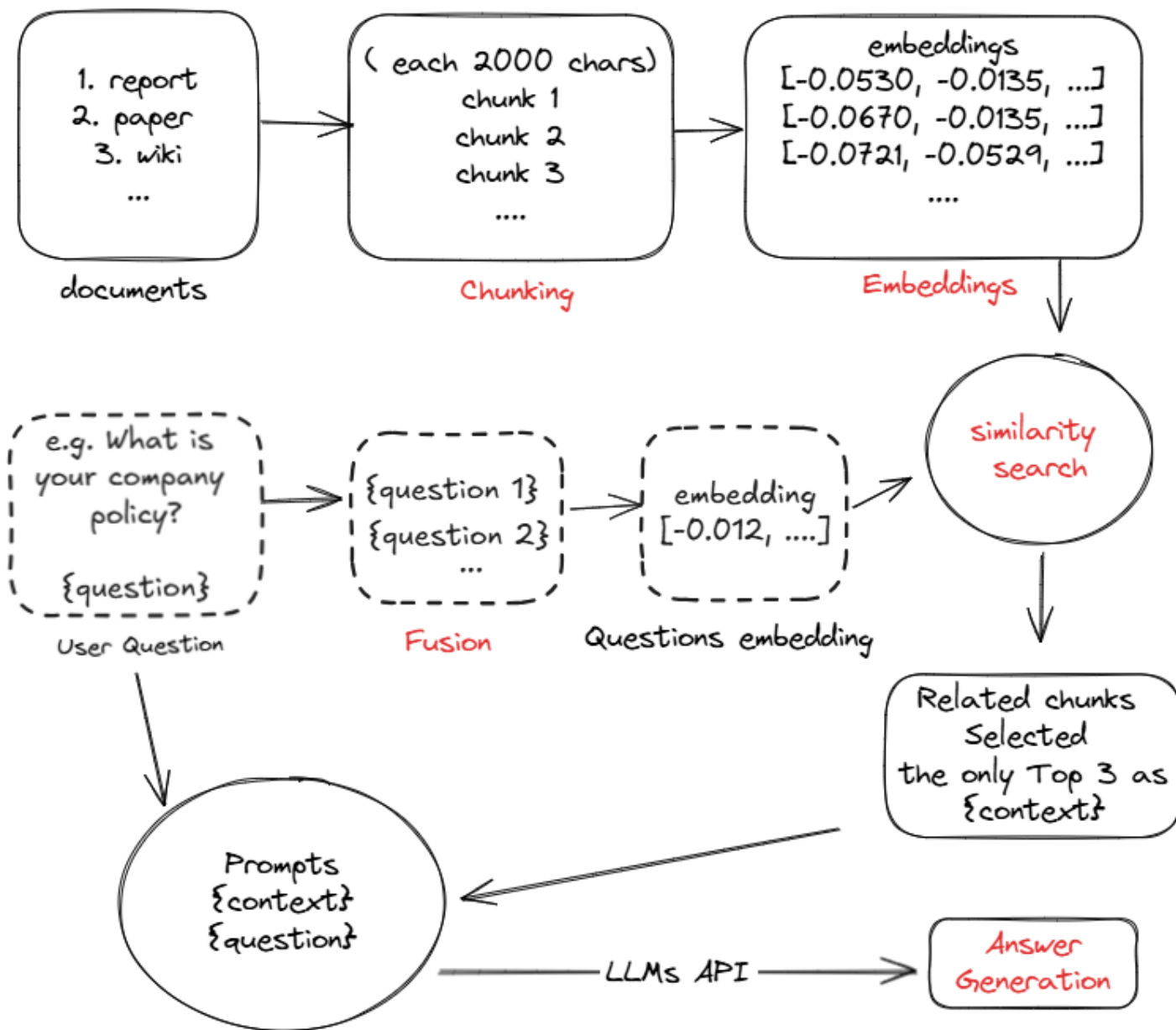
# RAG: An Overview

**RAG:**

Retrieval-Augmented Generation

It's a way for AI to give better answers by first searching for useful info in a database (e.g., books, websites) and then using that info to make a clear and smart response.

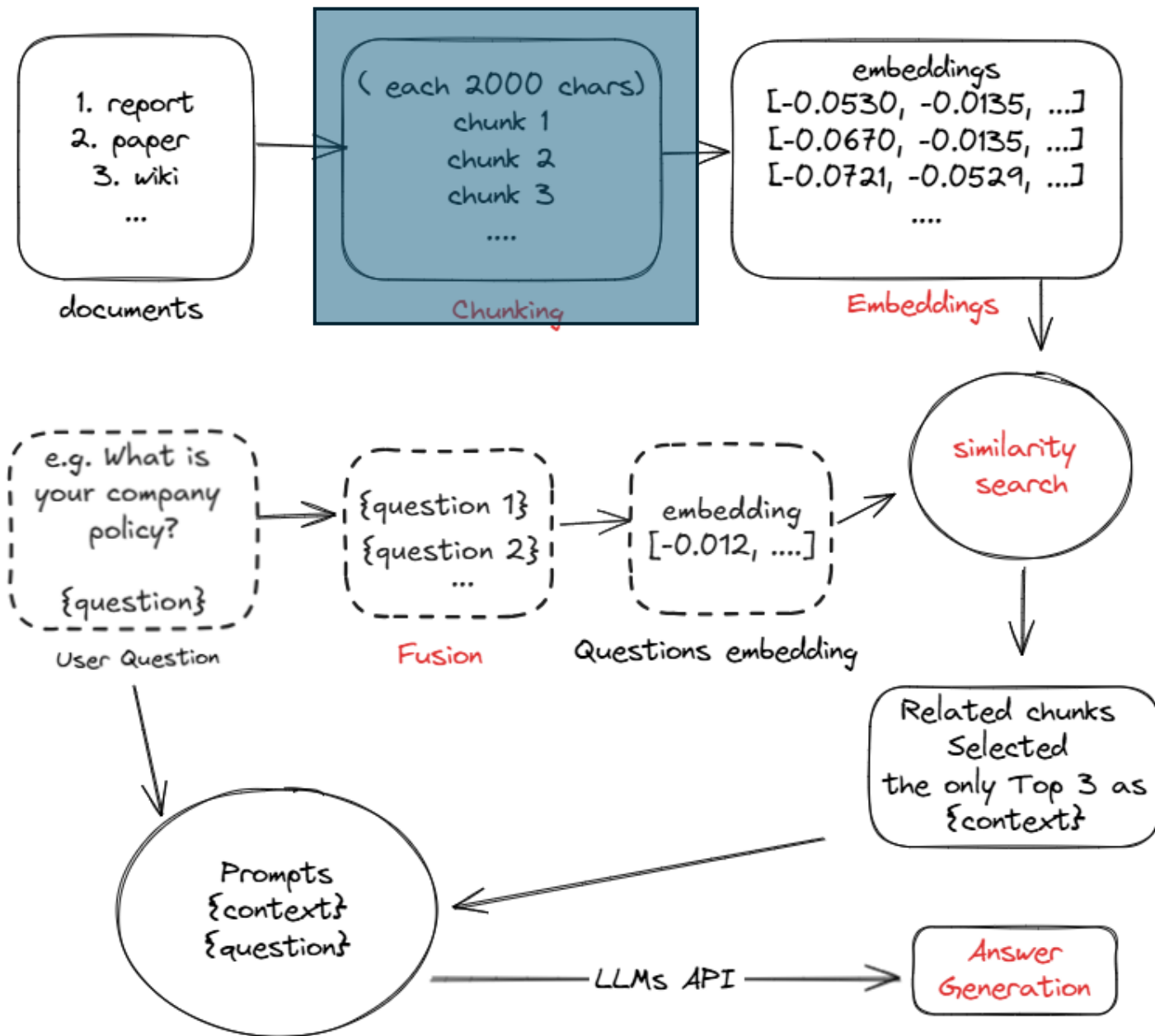It's like an assistant who checks the facts before talking!

# RAG: The Process
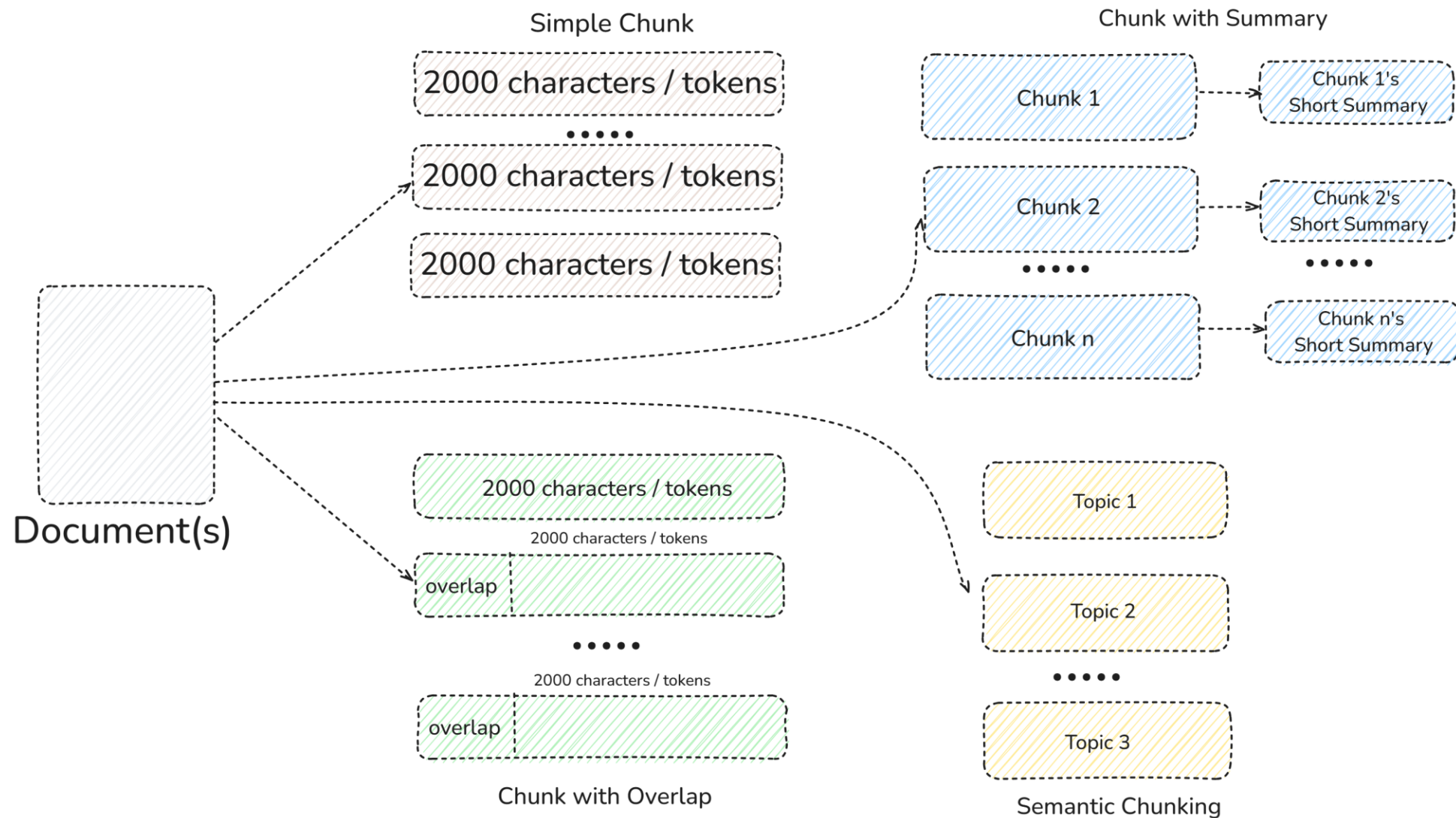
# RAG:
# The Process

# **Chunking**

# *Chunking*

- Definition
  - The process of breaking down large documents or datasets into smaller, manageable pieces (chunks) for the retriever to process.

- Types
  1. **Simple Chunk** (Character / Token cut, Sentence Cut)
  2. **Overlap Chunking**: Add overlapping text between chunks to preserve context (e.g., 20% overlap).
  3. **Chunk with Summary**: Creates chunks and generates a summary for each chunk to preserve context and aid retrieval.
  4. **Semantic Chunking**: Use NLP to split based on meaning (e.g., topic shifts).

**Tip**: Experiment with chunking strategies based on your data and use case!

# *Chunking*

### Simple Chunk

2000 characters / tokens

• • • • •

2000 characters / tokens

2000 characters / tokens

### Chunk with Summary

Chunk 1 ⇢ Chunk 1's Short Summary

Chunk 2 ⇢ Chunk 2's Short Summary

• • • • •                    • • • • •

Chunk n ⇢ Chunk n's Short Summary

**Document(s)**

### Chunk with Overlap

2000 characters / tokens

2000 characters / tokens
overlap | 2000 characters / tokens

• • • • •

2000 characters / tokens
overlap | 2000 characters / tokens

Chunk with Overlap

### Semantic Chunking

Topic 1

Topic 2

• • • • •

Topic 3

Semantic Chunking

# *Creating Chunks*

## Step 1

Identify source / data e.g.,
- SIOP Principles
  (Fifth Edition, 2018).
- 60-page PDF on
  personnel selection validation

## Step 2

Chunking Strategy  e.g.,
Character Cut:
- Split by 500 chars & 20% overlap
- Split by sections, e.g.,
  "Validity" vs. "Fairness"

## Step 3

Automate Tools  e.g.,
Character Cut:
- NLTK, langchain etc.
- Store Chunks (vector DB)

**Principles for the
Validation and Use of
Personnel Selection Procedures**

FIFTH EDITION
AUGUST 2018

AMERICAN
PSYCHOLOGICAL
ASSOCIATION

Test:
Can a chunk from

"Fairness and Bias" (p. 27)

answer

"What does SIOP say about predictive
bias?

# RAG: The Process

# Embeddings

# *Embeddings*

- **Definition**
  - Embeddings are numerical representations (vectors) of text, images, or other data in a high-dimensional space, capturing their meaning or context.

- **Usage in RAG Chatbots**
  - Enables the retriever to understand and compare the similarity between a user's query and stored chunks.
  - Enables semantic search: e.g., find SIOP sections on "predictive bias" without exact keywords).

# *Embeddings*



**Chunks**

| 0.15 | 0.56 | -0.21 | ...... |

| 0.25 | 0.66 | -0.81 | ...... |

| 0.45 | 0.26 | -0.11 | ...... |

...

**User Query**

"What SIOP principles said about Criterion-Related Validity ?"

**768-D space**

| 0.15 | 0.56 | -0.21 | ...... |

Search

**Plot**

**Text**     **Embedding**

"Criterion-Related Validity"

| 0.15 | 0.56 | -0.21 | ...... |

"Predictor-Criterion Relationship"

| 0.12 | 0.38 | -0.18 | ...... |

*I/O psychology topics*

"Turnover"

| 0.9 | 0.8 | -0.7 | ...... |

*A contrasting concept*

"Criterion-Related Validity"     "Turnover"

"Predictor-Criterion Relationship"

# *Creating Embeddings*

- Pre-trained models (e.g., all-MiniLM-L6-v2)
- APIs (OpenAI's Text-embedding models)

turn text into numbers that preserve meaning

**Pick Models and Tools**

e.g.,
- Python Library
- text-embedding-3-large

Text Chunk

Text Chunk

Text Chunk

Text Chunk

Chunks

vectors

vectors

vectors

vectors

embeddings

**Step 1: Pick a Model**          **Step 2: Process Chunks**          **Step 3: Generate Vectors**

# RAG:
# The Process

# Searching

# *Searching: Topic Fusion*

- **Fuse Topics**: Split query into two richer sub-queries:
  - "What's criterion validity?" →
    - "How is criterion validity defined?" (targets definitions).
    - "What evidence supports criterion validity?" (targets examples).
- **Embed Queries**: Convert both to vectors
- **Compare Vectors**: Match each to SIOP chunks (e.g., p. 18's "Evidence for criterion-related validity…").
- **Rank Results**: Pick top 2 chunks per sub-query (4 total) using cosine similarity.
- **Retrieve**: Return chunks with metadata (e.g., "SIOP Principles," "p. 18").

# *Searching: Topic Fusion*

# Searching: Cosine Similarity & Vector Search

- **Cosine Similarity**: a measure of how close 2 vectors are
  - 1 = identical
  - 0 = unrelated
- **Example:** Query "What's criterion validity?" vs. p. 18 chunk "Evidence for criterion-related validity…" → high score (~0.9).
- **Tools**:
  - FAISS (Local) : Fast vector search library
  - Pinecone (Cloud): Scalable vector DB

# RAG:
# The Process

## Prompt Forming

# Prompt Forming

- Goal: Combine query + retrieved chunk → generate coherent answer

- Example:

```
Query: {user_query}
Context:
{retrieved_chunk1}
{retrieved_chunk2}
{retrieved_chunk3}
{retrieved_chunk4}
Answer:"
```

# Build your own Chatbot – Code

*Code demo in Jupyter Notebooks*

# Build your own Chatbot: No Code via DIFY

- Many no code tools available

- Example: **DIFY**
  - Online cloud version (free & paid)
  - On-premises (Open-source, install on your server)

**Dify_** | **Define & Modify Do It For You**

https://dify.ai/

*DIFY Demo*



Principles for the Validation and Use of Personnel Selection Procedures

FIFTH EDITION
AUGUST 2018

AMERICAN PSYCHOLOGICAL ASSOCIATION

# DIFY: Landing Page

# DIFY: Chatbot

# DIFY: Create Knowledge

# DIFY: Upload Files

# DIFY: Text Chunking

# DIFY: Embedding

# DIFY: Debugging & Trialing the Chatbot

# DIFY: Debugging & Trialing the Chatbot

# DIFY: Using the Chatbot

# DIFY: Using the Chatbot

# Q&A Resources

- Q&A
- Slides available on **Whova**
- All materials (slides, code, and supplementary resources) found on **GitHub repo** here: https://github.com/karimhbadr1/SIOP_2025_ChatBot_Master_Tutorial