

✓ This Managed RAG Tool Just Made AI Agents Stupid Simple

[Download All](#)

Duration: 5:35

[Notes](#) [Quiz](#) [Flashcards](#) [Transcript](#)

Full Transcript

I put this rack AI agent in less than 5 minutes. If you have built a rack before, you will hate how easy this new method is. Let's open up to chat and send off this request. It has three different queries. How many employees and children were at the big island of Hawaii event in 2012? Where can find the internet and give me the definition of a drop-up. In 2012, the big island of Hawaii event had 293 employees and 185 children participating. And it says it was in this handbook on page 1620. Let's check this. When we opened up the handbook, we see it was 293 employees and 185 children. And it was on page 18 like the script here. The internet can be accessed on this URL. And it says it was on page 327. Let's check this. MSC on page 5. Here's the internet. And last but not least, a drop-up is an effort to restart the play and football. When the referee stops the game of outfall or offense. And the information is in this period of on page 15 to 19. And on page 17, we find if the referee stops the game and no fall or other offenses has secured, the game restarts with a drop-up ball. Click you can see the RKA agent does not only give us the correct answers, but also sites the correct source. So you can easily verify the information. For those of you who built RKA agents before, you know it's not so easy to get your agent to site the sources where it's pulling the information from. It's doable, but you have to set up a complex pipeline for that. I did all this by simply dragging and dropping my files. I didn't have to worry about rack implementation details like extracting that, chunking it, embedding it, or the retrieval of the information. This way you can easily give your AI agent access to your knowledge base even when you don't know exactly how it works. And the best part is I will explain exactly how and even provide you this template for free. That way you can simply share it by activating it and copy and pasting this URL here into your browser and then you can give anybody access to it. Like for example to your employees and without further ado, let's get started. First make sure to open up the Google Drive folder provided in the description box below and download the example PDF and also the RKA template so you can easily follow along. Next up we open pinkon.io and sign up if we don't have an account yet otherwise look in. Then we select on the left menu assistant and click create assistant and give it name like fastest, rack and click create assistant. Then you click on the folder here and simply drag and drop your PDF side here and click import. No pinkon takes care of the heavy lifting. It's already finished and now let's try it with our initial questions and like you can see we received the same answers like an all our initial example. Now let's connect it with an at n. For that you go to nad.io and click as start if you don't have an account yet otherwise sign in. Then you click create workflow then on the three dots and click import from file and select the rack template provided and now we simply need to adjust this node. First we need to change this URL. For it to go to pinkon, click on connect then select here shell and then copy this URL that you find here and paste it in here. Now we simply need to replace your API key. For that you go back to pinkon, click here on API keys, click create API key, give it a name like nad and test, click create key and then simply copy it and paste it here. Now we can click go back to canvas, click save, open chat and provide questions from the beginning. And like you can see it worked like in our initial example and whole works. First our AI syntax in our requests here and then splits it in four pinkon tasks. The search queries are big island of how I event employee account 2012, big island how I event children count 2012,

Flashcards (10 Cards)

Card 1

Q: What is the main benefit of using the managed RAG tool described in the lecture?

 Show Answer

A: The main benefit is the ease of implementation. The user can quickly create a RAG agent by simply dragging and dropping files, without needing to manage complex RAG implementation details like extracting, chunking, embedding, or retrieval. This simplifies the process of giving an AI agent access to a knowledge base.

Card 2

Q: What steps are involved in setting up the RAG agent using Pinkon.io?

 Show Answer

Card 4

Q: What is the role of n8n in the described RAG setup?

 Show Answer

Card 3

Q: How does the lecture showcase the verification of the RAG agent's responses?

 Show Answer

Card 6

Q: What is OpenRouter and how is it integrated into the setup?

 Show Answer

Card 5

Q: How are different AI models used within the Pinkon and n8n setup?

 Show Answer

Card 8

Q: How can you give your employees access to the RAG agent created with Pinkon, according to the lecture?

Question 1

What problem does the managed RAG tool discussed in the lecture primarily aim to solve?

A) Simplifying the process of building and deploying AI agents.

B) Increasing the processing speed of AI agents.

C) Enhancing the visual interface of AI applications.

D) Improving the accuracy of language translation models.

 **Explanation:** The lecture focuses on how the tool simplifies the process of creating RAG (Retrieval-Augmented Generation) AI agents, emphasizing the ease of use and lack of complex setup compared to traditional methods.

Question 2

According to the lecture, what is a key feature of the discussed RAG tool that makes it beneficial for users who have previously built RAG agents?

A) Automatic language translation capabilities.

B) The ability to cite the source of information used in its responses.

C) Integration with blockchain technology for data security.

D) Generation of marketing copy based on user input.

the response back to the user.

14. Tokens: Units of text used by AI models to process and generate language. Token count is often used to measure the usage of AI models and determine pricing.

Examples

- 1. Querying about the Hawaii event:** The lecturer demonstrates how to ask the RAG agent about the number of employees and children at an event in Hawaii in 2012. The agent correctly retrieves the information and cites the specific page in the handbook where the information is found.
 - 2. Defining "drop-up":** The lecturer asks the RAG agent for the definition of "drop-up" in football. The agent provides the correct definition and cites the relevant pages in the document.
 - 3. Integrating Pinkon with n8n:** The lecturer shows how to create a workflow in n8n that sends queries to Pinkon, retrieves the results, and then uses another AI model to generate a response. This example demonstrates how to chain together different operations to create a more complex and customized RAG system.
 - 4. Changing the AI Model in Pinkon:** The lecturer shows how to change the model to better suite the needs. The agent will utilize a different model for information retrieval based on the configurations.
-

Study Tips

- 1. Follow the Tutorial Step-by-Step:** Recreate the RAG agent using Pinkon and n8n by following the steps outlined in the lecture. This hands-on experience will solidify your understanding of the process.
 - 2. Experiment with Different AI Models:** Explore different AI models in Pinkon and n8n to see how they affect the performance and quality of the RAG agent. Pay attention to the trade-offs between different models in terms of accuracy, speed, and cost.
 - 3. Explore Pinkon's Documentation and Community:** Dive deeper into Pinkon's documentation and community forums to learn more about its features and capabilities. Engage with other users to ask questions and share your experiences.
-

Key Concepts

1. RAG (Retrieval Augmented Generation): A technique for improving the accuracy and reliability of AI language models by grounding them in external knowledge. The model first retrieves relevant information from a knowledge base and then uses that information to generate a response. This prevents the model from relying solely on its pre-trained knowledge, which might be outdated or inaccurate.

2. AI Agent: In this context, an AI program designed to answer questions and provide information based on a specific knowledge base. It acts as an interface to the RAG system, taking user queries and delivering informative responses.

3. Pinkon.io: A platform that simplifies the creation of RAG agents by handling the complex technical details of data processing and retrieval. It automates tasks such as extracting text from documents, chunking the text into smaller segments, creating embeddings (vector representations) of the text, and retrieving the most relevant chunks in response to a query.

4. Data Extraction: The process of retrieving text and information from various document formats (e.g., PDFs). This involves converting the document into a machine-readable format and identifying the relevant text content.

5. Chunking: Dividing large documents into smaller, manageable segments of text. This is important for efficient retrieval and processing of information. Smaller chunks allow for more precise matching of user queries.

6. Embeddings: Vector representations of text that capture the semantic meaning of words, phrases, and documents. These vector representations allow the system to compare the similarity between different pieces of text and identify the most relevant chunks for a given query.

7. Retrieval: The process of identifying and extracting the most relevant chunks of text from the knowledge base in response to a user query. This is typically done by comparing the embedding of the query to the embeddings of the chunks.

8. n8n.io: A workflow automation platform that allows users to connect different applications and services together to create automated workflows. In this context, it is used to integrate Pinkon with other AI models and services, allowing for more complex processing and routing of queries and responses.

9. API Key: A unique identifier that allows a program or service to access an API (Application Programming Interface). API keys are used to authenticate requests and track usage.

10. AI Model Selection: The process of choosing the appropriate AI model for different tasks within the RAG system. The lecturer highlights the importance of using different models for information retrieval (Pinkon's model) and response generation (the AI agent's model).

11. Open Router: A service that allows access to multiple AI models through a single API. This can be useful for selecting the best model for a given task or for comparing the performance of different models.

12. Citations/Source Attribution: The ability of the RAG agent to provide the specific source and location (e.g., page number) within the knowledge base from which the information was retrieved. This is crucial for verifying the accuracy and reliability of the information.

13. Workflow Automation: Automating repetitive tasks by creating a sequence of actions that are executed automatically. Using n8n is an example of workflow automation as it takes the API call and automatically gets the response back to the user.

14. Tokens: Units of text used by AI models to process and generate language. Token count is often used to measure the usage of AI models and determine pricing.



Lecture Intelligence

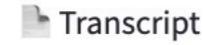
Transform lectures into AI-powered study materials



This Managed RAG Tool Just Made AI Agents Stupid Simple

[!\[\]\(950a62bbddad88d64435fd35607dfc42_img.jpg\) Download All](#)

Duration: 5:35

 Notes  Quiz  Flashcards  Transcript

Summary

Here's an analysis of the provided lecture transcript, broken down into a summary, key concepts, examples, and study tips:

The lecture introduces a simplified method for building Retrieval Augmented Generation (RAG) AI agents using a tool called Pinkon.io. The demonstrator emphasizes the ease of use, stating that a functional RAG agent can be created in under 5 minutes. The process involves simply uploading documents to Pinkon, which handles the complexities of data extraction, chunking, embedding, and retrieval. The lecture highlights the ability to provide the AI agent with a knowledge base without needing to understand the underlying technical implementation details.

The lecture then walks through a practical demonstration of creating a RAG agent. It begins by showcasing the agent answering complex queries and accurately citing the sources within uploaded documents. The tutorial then proceeds to demonstrate how to integrate Pinkon with n8n, a workflow automation platform, allowing for more complex processing and routing of queries and responses. This integration allows users to chain together various operations, such as sending queries to Pinkon, processing the results, and then formulating a final answer using a separate AI model.

The lecturer also discusses customization options within Pinkon and n8n, including the ability to select different AI models for both information retrieval and answer generation. They highlight the importance of using different models for these two tasks: one for accurately finding the correct information within the knowledge base (Pinkon's model) and another for generating a coherent and user-friendly response based on that information (the AI agent's model). The lecture concludes with a brief overview of Pinkon's pricing and a call to action to join a community for more structured learning and support.

⚙️ Configuration

Backend URL 

`https://subclavate-hypatia-squashily.ng`

Gemini API Key 

..... 

📊 Options

Quiz Questions 

10

3 20

Flashcards 

15

5 30

 New Lecture

 Tip: Keep your Colab notebook running!

Lecture Intelligence

Transform lectures into AI-powered study materials

 YouTube Lecture URL 

`https://youtu.be/-OXHWET_RFI?si=l5X2Nuy-U2f8zDQM`

 Generate Study Materials

Made with ❤️ using Streamlit | Powered by OpenAI Whisper & Google Gemini

Process lectures • Generate notes • Create quizzes • Make flashcards



Configuration

Backend URL ?

`https://subclavate-hypatia-squashily.ng`

Gemini API Key ?

..... eye icon

Options

Quiz Questions 10



Flashcards 15



New Lecture

Tip: Keep your Colab notebook running!

Lecture Intelligence

Transform lectures into AI-powered study materials

YouTube Lecture URL

`https://youtu.be/-OXHWET_RFI?si=l5X2Nuy-U2f8zDQM`

Generate Study Materials

Job started! ID: 32af965f-4980-4250-8864-daeddc9ffad5

Status: transcribing | Progress: 30%