

# SAP Joule with RAG and Document Grounding: Next-Gen Document-Centric Processes

**Arun Chinnannan Balasubramanian**

Verizon Communications  
Basking Ridge, USA

## Abstract

The rapid evolution of artificial intelligence (AI) has introduced groundbreaking methodologies like Retrieval-Augmented Generation (RAG) that enable systems to provide accurate, real-time, and contextually aware responses by leveraging external knowledge repositories. SAP Joule, a business AI solution, can be integrated with RAG and document grounding techniques to enhance decision-making and operational efficiency across the SAP ecosystem. This paper delves into the concepts of RAG via document grounding, explores their synergies with SAP Joule, and provides a comprehensive implementation guide for activating these capabilities within the SAP ecosystem. Document grounding will aid the rich factual results in conversational AI setup of corporate systems. Usage can be endless, where documents can be read and give a factual output.

**Keywords:** Retrieval-Augmented Generation (RAG), Document Grounding, Large Language Models (LLMs), SAP Joule, Conversational AI.

## I. Introduction

Overview of SAP Joule: SAP Joule is an AI-powered solution designed to enhance enterprise applications by providing actionable insights and predictive capabilities. SAP Joule is one of the services that run on SAP Business Technology Platform(BTP) in a multi-cloud / Cloud Foundry environment (CF). SAP Joule leverages the SAP Business Technology Platform to simplify complex integration landscapes. To activate or set up the Joule service with the Line of Business (LoB), It needs certain services. SAP Build Work Zone that supports Joule interactions rely on navigational patterns. SAP Cloud Identity Service offers Identity Authentication Services (IAS) and Identity Provisioning Services (IPS), so customers can achieve a Consistent Identity lifecycle(CIL), as a Identity Management[1]. By leveraging AI models, businesses can improve their operational workflows and achieve greater productivity[2].

### Introduction to RAG and Document Grounding

RAG combines generative AI models with retrieval systems, enabling the use of knowledge repositories to produce precise and contextually enriched outputs. The retrieval mechanism employs a high-performance index that queries structured and unstructured data in real time, while the generative model integrates this retrieved data to construct coherent, domain-specific responses[4]. Advanced ranking algorithms prioritize relevance, ensuring responses are both accurate and aligned with the specific enterprise context. Document grounding ensures that AI responses are rooted in specific, authoritative documents, leveraging structured indexing and retrieval pipelines to enhance reliability and trustworthiness. This approach integrates contextually relevant documents into the AI processing workflow, ensuring outputs are not only accurate but also aligned with enterprise-grade compliance and operational standards.

## II. Literature review and Understanding the Core Concepts

### SAP Joule: Features and Capabilities

SAP Joule offers advanced analytics, predictive insights, and integration with SAP's suite of enterprise tools. Key features like:

- Predictive analytics for supply chain and financial operations leveraging SAP Analytics Cloud and SAP Integrated Business Planning (IBP). These tools use advanced machine learning algorithms to detect trends, forecast demand, and optimize inventory levels in real time (SAP, 2023).
- Seamless integration with SAP Business suite application and SAP Business Technology Platform through APIs and middleware such as SAP Cloud Integration, enabling the unification of disparate enterprise processes into a cohesive ecosystem.

### Retrieval-Augmented Generation (RAG): A Deep Dive

RAG is a hybrid AI approach combining large language models (LLMs) with a retrieval mechanism. Retrieval Mechanism fetches relevant documents or data from external sources. Generative Component, produces coherent responses using retrieved data as context.

This methodology overcomes the limitations of standalone generative models by providing grounded, accurate, and updated information [4].

### Document Grounding: Definition and Benefits

While Large Language Models (LLMs) excel at understanding and generating human-like text, they often lack the specificity and context of business data. To fully unlock the potential of LLMs, particularly in conversational use cases, providing the right context to ground LLMs in business reality is essential. This approach ensures that the generated answers are:

- Relevant (Factually Correct): Responses are aligned with specific enterprise scenarios.
- Reliable (Based on Up-to-Date Information): Outputs reference the most current and accurate data.
- Responsible (Trace Errors to Sources): Factual discrepancies can be traced back to the original data source.

SAP Joule leverages the Retrieval-Augmented Generation (RAG) technique to optimize the output of LLMs. This approach integrates external document repositories, dynamically retrieving information that is specific to use cases and not part of the LLM's original training dataset. Through this method, document grounding bridges the gap between general AI capabilities and enterprise-specific requirements, delivering precision and accountability in AI-driven workflows[3]. Key technical benefits include:

- Dynamic Querying: Real-time access to enterprise data ensures contextual relevance.
- Scalable Integration: Handles complex queries across structured and unstructured data sources.
- Enhanced Traceability: Links AI-generated outputs to authoritative documents, facilitating compliance and audit readiness.

This robust framework not only augments the capabilities of LLMs but also addresses the demands of enterprise-grade accuracy, reliability, and traceability in AI-driven decision-making.

## Enhancing Contextual Understanding

By integrating RAG with SAP Joule, businesses can provide AI models with real-time access to enterprise data, enhancing contextual awareness and decision-making. For example, a supply chain management system can leverage SAP Joule to retrieve inventory levels and supplier lead times dynamically. This real-time data integration allows decision-makers to identify bottlenecks and optimize workflows, a process supported by SAP's ability to harmonize structured data with unstructured records from external sources. Such use cases demonstrate the value of RAG in complex, data-driven environments. RAG enables SAP Joule to address knowledge gaps by retrieving information from sources, ensuring comprehensive and relevant outputs.

## Use Cases in the SAP Ecosystem

- **Finance:** To provide a SOX audit trail for an intercompany invoice that can go through specific documents and provide factual evidence.
- **Supply Chain:** Supplier invoice automation aggregating the multiple invoice documents against goods receipt and customs landing documents.
- **HR:** Enhancing employee engagement through personalized AI-driven interactions [3]. Like conversational AI can be set up to provide benefits for short term disability.

## III. Methodology and Implementation Framework

**Prerequisites:** Before initiating the implementation, ensure the following prerequisites are met:

**Active SAP Joule Setup:** SAP Joule must be operational within your line of business, fully configured with necessary enterprise scenarios and models aligned with organizational requirements. This setup includes ensuring connectivity with SAP Business Technology Platform (BTP) and alignment with SAP AI Core for streamlined workflow integration.. Ensure proper activation as per SAP documentation.

**AI Units SKU:** Procure the AI Units SKU, a critical license for AI operations, and verify the entitlement within SAP BTP Global Account. Ensure that sufficient AI unit quotas are allocated to support high-demand scenarios for real-time processing and retrieval operations.

**Repository Access:** Configure Microsoft SharePoint with the necessary permissions and structure to serve as a document repository. Validate SharePoint's API connectivity with SAP systems, ensuring compatibility for document ingestion, classification, and retrieval workflows.

**SAP BTP Access:** Administrator privileges in SAP BTP are required to configure subaccounts, manage entitlements, and deploy services. Ensure Cloud Foundry is fully operational with defined spaces and quotas aligned to RAG and grounding processes.

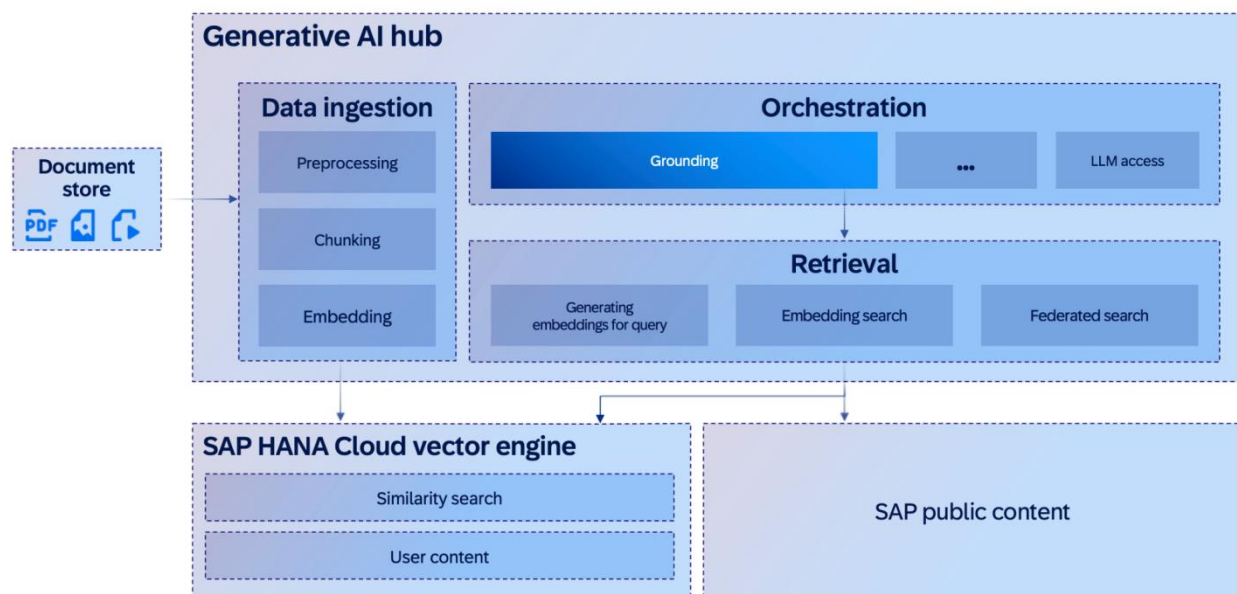
**Command-Line Tools:** Deploy GitBash CLI or a similar terminal interface for executing integration scripts and API commands. These tools are essential for performing advanced configurations such as service instance creation and secure authentication setup.

## IV. Technical Architecture of SAP Joule and RAG Integration

A robust architecture for integrating SAP Joule with RAG includes:

- **Data Layer:** Leverage SAP HANA for high-performance, in-memory data processing and Microsoft SharePoint for external document repository management. Establish federated access to unify these layers, enabling seamless data flow and retrieval for RAG workflows.
- **Integration Layer:** Employ APIs to facilitate seamless data retrieval.
- **Processing Layer:** Leverage AI models and RAG components to generate grounded responses.
- **Presentation Layer:** Deploy SAP Fiori applications and dashboards for user interaction.

SAP Joule is built on a modular, scalable architecture designed to seamlessly integrate retrieval-augmented AI capabilities. At its core, the Retrieval Index serves as a high-performance layer for accessing both structured and unstructured data from sources such as SAP HANA and Microsoft SharePoint. This index enables real-time querying and prioritizes document relevance using advanced ranking algorithms.



**Diagram 1: Architecture of the Grounding capability inside Generative AI hub. Courtesy: Sebastien np, SAP**

Diagram 1 is an illustration of the architecture of grounding capability. SAP's Document Grounding service utilizes Retrieval-Augmented Generation (RAG) to enrich Large Language Model (LLM) responses with critical context through a structured multi-step process comprising indexing, storage, and retrieval[5]:

- **Indexing Pipeline:** This stage ingests data from various sources, processes it into smaller, manageable chunks, and generates embeddings. These embeddings are stored in the SAP HANA Vector Database for efficient querying.
- **Vector Database:** Powered by the SAP HANA Vector Engine, the database enables high-performance querying of vector embeddings.
- **Retrieval:** When a user submits a query, the system generates vector embeddings using the LLM model. These embeddings are then matched against the HANA Vector Database to identify and retrieve the most relevant text chunks.

An orchestration layer ensures smooth integration of RAG models with SAP workflows, facilitating seamless interaction between data retrieval and generative outputs within complex enterprise environments, optimizing contextual data flow and output generation. To ensure enterprise applicability, fine-tuned language models are deployed, balancing accuracy and response performance for specific scenarios.

Together, these components form a cohesive system, enabling dynamic and contextually aware AI outputs. **Core Components and Workflows** The architecture includes a sophisticated document pipeline that ingests, indexes, and classifies data for retrieval. This pipeline facilitates real-time querying by implementing robust APIs and algorithms, ensuring only the most relevant documents are accessed. The workflow dynamically generates responses by integrating retrieved data with generative AI outputs, tailoring results to enterprise needs. **Security and Performance Considerations** Security is paramount in SAP Joule's design. All data exchanges are encrypted using TLS, ensuring secure communications across endpoints. Role-based access control, coupled with OAuth tokens, provides stringent authentication mechanisms. Performance is optimized through caching techniques and parallel processing, which reduce latency and ensure smooth operation even under high query volumes.

### Workflow Customization in SAP Systems

- **Service Activation:** Activate the Document Grounding service in SAP BTP Subaccount and assign entitlements.
- **Service Instance and Key Creation:** Create a service instance and generate a service key for integration.
- **Authentication Configuration:** Establish secure trust between SAP BTP and Microsoft SharePoint.
- **Process Modeling:** Embed RAG workflows in existing SAP processes using Business Process Modeling tools.

## V. Challenges Implementing and Limitations

### Security and Compliance Limitations

Despite robust encryption and access control, maintaining compliance with industry regulations such as GDPR and CCPA across global operations remains a challenge. Organizations must ensure to get the clearance of usage of these tools to the dataset that is being processed. This could take a challenging path depending on the nature of the business and countries that it deals with.

### Limitations of Document Grounding

Document Grounding is limited to a maximum of 2,000 documents, restricted to PDF or Word formats, which must be in English and plain text. These constraints may exclude critical non-English or multimedia-based resources, impacting comprehensiveness. Additionally, the system relies on a daily content refresh schedule, which might not accommodate real-time updates required by dynamic operations. Microsoft SharePoint is configured as the primary repository, and while robust, its reliance on specific API configurations can lead to integration challenges when handling diverse enterprise setups.

### Dependency on Technical Expertise

Successful deployment and maintenance require advanced technical expertise in configuring enterprise-grade AI solutions. This includes setting up SAP AI Core and managing SAP BTP services to align with organizational needs. Teams must possess in-depth knowledge of machine learning frameworks, data pipeline optimization, and integration strategies to ensure seamless operation. For instance, tuning AI models to handle industry-specific requirements or customizing workflows for RAG retrieval tasks necessitates domain-specific expertise. Without such proficiency, organizations risk suboptimal implementation, reduced performance, and operational bottlenecks. technical skills in configuring SAP AI

Core, managing BTP services, and tuning AI models for specific use cases. This dependency on specialized expertise can be a bottleneck for organizations with limited technical resources.

## VI. Conclusion

The integration of SAP Joule with Retrieval-Augmented Generation (RAG) and document grounding provides an unparalleled framework for addressing the complexities in Document based conversational AI or workflow that needs to be built. This brings a new platform, that innovative thoughts can flow in solving an organization problem, rather than looking at solving in classical ways. Unlike traditional document processing methods, which often rely on static data and limited adaptability, RAG introduces dynamic context-specificity and scalability by leveraging external repositories and real-time data indexing. Benefits of RAG in document processing are,

1. **Dynamic Relevance:** Traditional document retrieval methods often fail to contextualize outputs dynamically. RAG, by contrast, uses advanced ranking algorithms and real-time indexing to ensure the most relevant and up-to-date information is retrieved.
2. **Enhanced Accuracy:** By grounding AI responses in validated and authoritative documents, RAG significantly reduces errors and enhances decision-making reliability. Ex- An employee can ask in conversational AI to provide the policy for international travel and approvals he has to make. It can read through documents and summarize it to a factual content.
3. **Scalability and Adaptability:** Unlike static systems, RAG seamlessly integrates structured and unstructured data from multiple sources, ensuring scalability for diverse and large-scale enterprise applications.
4. **Traceability:** The ability to trace AI outputs back to their originating data sources enhances audit readiness and compliance, which is often limited in traditional systems. Ex- A Conversational AI can be built to prepare audit trace documents for a 'Accounts payable invoice for a supplier'.
5. **Enterprise-Specific Optimization:** Traditional methods struggle with domain-specific adjustments. RAG integrates with SAP's ecosystem, offering tailored outputs aligned with enterprise workflows.

In summary, SAP Joule with RAG and document grounding not only enhances operational efficiency but also establishes a robust foundation for innovation. The methodologies described in this document enable enterprises to harness the full potential of AI, paving the way for smarter, faster, and more reliable decision-making processes.

## VII. References

1. Nagesh Caparthy. '*Joule - Unified Setup: Bridging Simplicity and Performance*'. Journal in SAP Technology community. Joule fundamentals in Nov 2024.
2. SAP SE. '*SAP Joule: AI-driven innovations in enterprise software*'. 2023.
3. Brown, T. '*AI and Enterprise Applications: A Modern Approach*'. TechPress. 2022.
4. Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., & Yih, W. T. (2020). '*Dense passage retrieval for open-domain question answering*'. 2020.
5. Sebastien nb. '*Harness retrieval-augmented generation in Joule and Generative AI Hub*'. Journal in SAP Technology community, Oct 2024.