



Я хочу создать проект по использованию ИИ в космосе. Мне пришла идея ИИ помощника (допустим на космическом корабле). Можно ли создать модель, использующую RAG систему с документациями, инструкциями, возможностью выхода в интернет и подобными функциями, чтобы создать такого помощника для корабля?

Можно вывести эту идею дальше - до создания AI агента с возможностями обращаться к каким нибудь функциям на корабле (продумать)

напиши исследование о том, как использовать RAG систему с возможностью подгрузки материалов и технической документации

Теперь у меня достаточно информации для создания подробного исследования.
Позвольте мне составить комплексный отчет на русском языке.

Исследование использования RAG-систем для создания ИИ-помощника космического корабля

Введение: Концепция и осуществимость проекта

Создание ИИ-помощника для космического корабля с использованием RAG (Retrieval-Augmented Generation) системы не только возможно, но и представляет собой перспективное направление развития космических технологий. Ваша идея полностью реализуема и уже имеет прецеденты в реальной космической отрасли^{[1] [2] [3]}.

Основные возможности такой системы:

- Работа с технической документацией и инструкциями в режиме реального времени
- Доступ к базе знаний о системах корабля
- Возможность выхода в интернет (при наличии связи)
- Интеграция с функциями управления системами корабля через AI-агента
- Автономное принятие решений в критических ситуациях

Примеры из практики показывают актуальность концепции: робот CIMON на МКС использует IBM Watson для анализа речи и взаимодействия с экипажем^{[1] [4]}, система Alexa тестировалась на МКС в рамках миссии Axiom-3^[2], а китайская станция Tiangong оснащена роботом Xiao Hang для помощи экипажу^[1].

Часть 1: Архитектура RAG-системы для космического применения

Основные компоненты RAG-системы

RAG-система состоит из трех ключевых этапов, каждый из которых критически важен для космического применения^{[5] [6] [7]}:

1. Retrieval (Поиск и извлечение)

Система преобразует запрос пользователя в векторное представление и ищет релевантные документы во внешней базе знаний. Для космического корабля это может включать^{[5] [8]}:

- Технические руководства по всем системам корабля
- Процедуры реагирования на чрезвычайные ситуации
- История предыдущих миссий и решенных проблем
- Научные данные и результаты экспериментов
- Медицинские протоколы

2. Augmentation (Дополнение)

Найденная информация добавляется к исходному запросу пользователя, создавая расширенный контекст для языковой модели^{[5] [7] [9]}:

Исходный запрос: "Как устранить утечку в модуле жизнеобеспечения?"

Дополненный контекст: [Релевантные разделы из технической документации] + Исходный запрос

3. Generation (Генерация)

Языковая модель (LLM) генерирует точный ответ на основе дополненного контекста, используя как свои предобученные знания, так и извлеченную специфическую информацию^{[5] [6] [7]}.

Преимущества RAG для космических применений

Устранение галлюцинаций: RAG значительно снижает вероятность генерации ложной информации, так как модель опирается на проверенные документы^{[6] [7]}. В космосе, где ошибка может стоить жизни, это критически важно.

Актуальность данных: Внешнюю базу знаний можно обновлять без переобучения модели^{[7] [9]}. Новые процедуры, обновления оборудования или результаты недавних исследований становятся доступны немедленно.

Прозрачность и проверяемость: Система может предоставлять ссылки на исходные документы, позволяя экипажу проверить информацию ^[6] ^[7] ^[10]. Это особенно важно для критических операций.

Специализированные знания: RAG позволяет работать с узкоспециализированной технической информацией без необходимости дообучать большую языковую модель на всей документации ^[7] ^[11].

Автономность: Система может работать офлайн, что критично для дальних космических миссий с задержками связи ^[10]. Вся база знаний хранится локально на борту корабля.

Часть 2: Техническая реализация RAG-системы с LangChain

Пошаговая реализация базовой RAG-системы

Шаг 1: Установка необходимых компонентов

```
# Установка основных библиотек
pip install langchain langchain-openai langchain-community
pip install faiss-cpu # Векторная база данных
pip install pypdf # Для работы с PDF
pip install openpyxl # Для Excel файлов
```

Шаг 2: Загрузка и обработка документации

LangChain предоставляет множество загрузчиков документов для различных форматов ^[12] ^[13] ^[14] ^[15]:

```
from langchain_community.document_loaders import (
    PyPDFLoader,
    TextLoader,
    UnstructuredExcelLoader,
    CSVLoader
)
from langchain.text_splitter import RecursiveCharacterTextSplitter

# Загрузка технической документации
pdf_loader = PyPDFLoader("spacecraft_manual.pdf")
excel_loader = UnstructuredExcelLoader("systems_data.xlsx")
text_loader = TextLoader("procedures.txt", encoding='utf-8')

# Загрузка документов
documents = []
documents.extend(pdf_loader.load())
documents.extend(excel_loader.load())
documents.extend(text_loader.load())

# Разбиение на чанки
text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=1000,          # Размер чанка в символах
    chunk_overlap=200,        # Перекрывание между чанками
    separators=["\n\n", "\n", ". ", " ", ""],
```

```
        length_function=len
    )

    chunks = text_splitter.split_documents(documents)
```

Оптимальные параметры чанкинга для технической документации:

- **Размер чанка:** 512-1024 токена для баланса между точностью и контекстом [\[16\]](#) [\[17\]](#) [\[18\]](#)
- **Перекрытие:** 128-200 токенов для сохранения связности [\[18\]](#) [\[19\]](#) [\[20\]](#)
- **Стратегия:** Рекурсивное разбиение с учетом структуры документа [\[17\]](#) [\[19\]](#) [\[20\]](#)

Шаг 3: Создание эмбедингов и векторной базы данных

```
from langchain_openai import OpenAIEmbeddings
from langchain_community.vectorstores import FAISS

# Создание эмбедингов
embeddings = OpenAIEmbeddings(model="text-embedding-3-small")

# Создание векторной базы данных
vectorstore = FAISS.from_documents(
    documents=chunks,
    embedding=embeddings
)

# Сохранение базы для последующего использования
vectorstore.save_local("spacecraft_knowledge_base")

# Загрузка существующей базы
# vectorstore = FAISS.load_local("spacecraft_knowledge_base", embeddings)
```

Выбор векторной базы данных:

Для космического применения рекомендуется использовать [\[21\]](#) [\[22\]](#) [\[23\]](#) [\[24\]](#):

- **FAISS:** Быстрая, работает локально, отлично подходит для офлайн-режима, поддержка GPU [\[21\]](#) [\[23\]](#)
- **Chroma:** Простая в использовании, встроенная персистентность, хороша для прототипирования [\[21\]](#) [\[22\]](#)
- **Qdrant:** Масштабируемая, с расширенной фильтрацией по метаданным [\[22\]](#)

Шаг 4: Построение RAG-цепочки

```
from langchain_openai import ChatOpenAI
from langchain.chains import RetrievalQA
from langchain.prompts import ChatPromptTemplate

# Инициализация языковой модели
llm = ChatOpenAI(
    model="gpt-4",
    temperature=0.1 # Низкая температура для точных ответов
```

```

)

# Создание ретривера
retriever = vectorstore.as_retriever(
    search_type="similarity",
    search_kwargs={"k": 3} # Возвращать топ-3 релевантных документа
)

# Создание промпта
prompt_template = """You are an AI assistant for spacecraft operations.
Use the following context from technical documentation to answer the question.
If you don't know the answer, say so clearly. Do not make up information.

Context: {context}

Question: {question}

Answer: """

prompt = ChatPromptTemplate.from_template(prompt_template)

# Создание RAG цепочки
qa_chain = RetrievalQA.from_chain_type(
    llm=llm,
    chain_type="stuff",
    retriever=retriever,
    return_source_documents=True
)

# Использование системы
query = "What is the procedure for emergency oxygen system activation?"
result = qa_chain.invoke({"query": query})

print(f"Answer: {result['result']}")
print(f"Source documents: {len(result['source_documents'])}")

```

Выбор модели эмбедингов

Сравнение популярных моделей для RAG-систем [\[25\]](#) [\[26\]](#) [\[27\]](#) [\[28\]](#) [\[29\]](#):

Модель	Размер	Преимущества	Недостатки
OpenAI text-embedding-3-small	-	Высокая точность, многоязычность, платная API	Требует интернет, стоимость \$0.02/1М токенов
OpenAI text-embedding-3-large	-	Максимальная точность	Дороже (\$0.13/1М токенов), медленнее
BAAI/bge-base-en-v1.5	110M	Отличное качество, бесплатная, офлайн	Только английский
sentence-transformers/all-MiniLM-L6-v2	22M	Очень быстрая, легковесная	Чуть ниже качество
nomic-embed-text-v1	~500M	Хорошее качество, открытая	Больше ресурсов

Рекомендация для космического корабля: Использовать комбинацию - OpenAI для облачных операций с Земли и bge-base-en-v1.5 для автономной работы на борту^{[25] [28]}.

Продвинутые техники для улучшения RAG

1. Гибридный поиск

Комбинация векторного и ключевого поиска для повышения точности^{[5] [16]}:

```
from langchain.retrievers import EnsembleRetriever
from langchain_community.retrievers import BM25Retriever

# Векторный ретривер
vector_retriever = vectorstore.as_retriever(search_kwargs={"k": 5})

# Ключевой ретривер (BM25)
bm25_retriever = BM25Retriever.from_documents(chunks)
bm25_retriever.k = 5

# Ансамбль
ensemble_retriever = EnsembleRetriever(
    retrievers=[vector_retriever, bm25_retriever],
    weights=[0.7, 0.3] # 70% векторный, 30% ключевой
)
```

2. Улучшенное разбиение документов

Семантический чанкинг сохраняет смысловые единства^{[17] [19] [20] [30]}:

```
from langchain.text_splitter import RecursiveCharacterTextSplitter

# Для технических документов с учетом структуры
text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=1024,
    chunk_overlap=128,
    separators=[
        "\n## ", # Заголовки второго уровня
        "\n### ", # Заголовки третьего уровня
        "\n\n", # Параграфы
        "\n", # Строки
        ". ", # Предложения
        " " # Слова
    ],
    add_start_index=True # Добавить индекс начала для трассировки
)
```

3. Метаданные и фильтрация

Добавление метаданных для точного поиска^{[17] [31]}:

```
# Добавление метаданных при загрузке
for doc in documents:
```

```

doc.metadata.update({
    "system": "life_support",
    "priority": "critical",
    "last_updated": "2025-11-03"
})

# Фильтрация при поиске
retriever = vectorstore.as_retriever(
    search_kwargs={
        "k": 5,
        "filter": {"system": "life_support"}
    }
)

```

4. Обновление базы знаний в реальном времени

Для поддержания актуальности данных [\[9\]](#) [\[11\]](#) [\[32\]](#):

```

def update_knowledge_base(new_documents):
    """Добавление новых документов в существующую базу"""
    new_chunks = text_splitter.split_documents(new_documents)
    vectorstore.add_documents(new_chunks)
    vectorstore.save_local("spacecraft_knowledge_base")

def refresh_document(doc_id, new_content):
    """Обновление конкретного документа"""
    # Удаление старой версии
    vectorstore.delete([doc_id])
    # Добавление новой версии
    new_doc = create_document(new_content)
    vectorstore.add_documents([new_doc])

```

Часть 3: Создание AI-агента с функциональными возможностями

От помощника к автономному агенту

Следующий уровень — создание AI-агента, который не только отвечает на вопросы, но и может выполнять действия, взаимодействуя с системами корабля [\[33\]](#) [\[34\]](#) [\[35\]](#) [\[36\]](#) [\[37\]](#). Это реализуется через механизм **function calling** (вызов функций) [\[34\]](#) [\[38\]](#) [\[39\]](#) [\[40\]](#).

Архитектура AI-агента для космического корабля

Основные компоненты агента:

1. **RAG-система** — база знаний и контекст
2. **Reasoning Engine** — языковая модель для принятия решений
3. **Tool Interface** — набор функций для взаимодействия с системами
4. **Memory** — история взаимодействий и контекста
5. **Execution Layer** — безопасное выполнение команд

Примеры функций для космического корабля

1. Мониторинг телеметрии

Основываясь на системах TT&C (Telemetry, Tracking, and Command) [\[41\]](#) [\[42\]](#) [\[43\]](#) [\[44\]](#):

```
from langchain.tools import Tool

def get_system_telemetry(system_name: str) -> dict:
    """
    Получить телеметрические данные системы корабля.

    Args:
        system_name: Название системы (life_support, propulsion,
            power, thermal, communication, navigation)

    Returns:
        Текущие параметры системы
    """
    # Интеграция с системой телеметрии корабля
    telemetry_data = {
        "life_support": {
            "oxygen_level": 98.5,
            "co2_level": 0.3,
            "temperature": 21.5,
            "humidity": 45,
            "pressure": 101.3,
            "status": "nominal"
        },
        "propulsion": {
            "fuel_remaining": 85.2,
            "thrust_capability": 100,
            "engine_temperature": 450,
            "status": "standby"
        },
        "power": {
            "battery_charge": 92,
            "solar_array_output": 3.2, # kW
            "power_consumption": 2.8,
            "status": "nominal"
        }
    }

    if system_name in telemetry_data:
        return telemetry_data[system_name]
    else:
        return {"error": f"Unknown system: {system_name}"}

# Создание инструмента для агента
telemetry_tool = Tool(
    name="GetSystemTelemetry",
    func=get_system_telemetry,
    description="Получить текущие телеметрические данные системы корабля. "
        "Используйте для мониторинга состояния систем."
)
```


2. Управление системами корабля

```
def control_spacecraft_system(system: str, action: str, parameters: dict = None) -> str:
    """
    Управление системами космического корабля.

    Args:
        system: Система для управления
        action: Действие (activate, deactivate, adjust, emergency_shutdown)
        parameters: Дополнительные параметры для действия

    Returns:
        Результат выполнения команды
    """
    # Проверка безопасности
    if action == "emergency_shutdown" and system == "life_support":
        return "ERROR: Cannot shutdown life support system. Crew safety critical."

    # Симуляция выполнения команды
    command_log = {
        "timestamp": datetime.now().isoformat(),
        "system": system,
        "action": action,
        "parameters": parameters,
        "status": "executed"
    }

    # Здесь будет реальная интеграция с системами корабля
    return f"Command executed: {system} - {action}"

control_tool = Tool(
    name="ControlSpacecraftSystem",
    func=control_spacecraft_system,
    description="Управление системами корабля. ВНИМАНИЕ: Используйте только "
               "после подтверждения экипажем для критических операций."
)
```

3. Доступ к базе процедур

```
def search_procedures(situation: str, system: str = None) -> str:
    """
    Поиск процедур в базе знаний по ситуации.

    Args:
        situation: Описание ситуации или проблемы
        system: Конкретная система (опционально)

    Returns:
        Релевантные процедуры из документации
    """
    # Использование RAG для поиска процедур
    query = f"Emergency procedure for {situation}"
    if system:
        query += f" in {system} system"
```

```

results = retriever.get_relevant_documents(query)

procedures = "\n\n".join([doc.page_content for doc in results[:3]])
return procedures

procedure_tool = Tool(
    name="SearchProcedures",
    func=search_procedures,
    description="Поиск процедур и инструкций в технической документации корабля."
)

```

4. Диагностика и поиск неисправностей

```

def diagnose_system_issue(system: str, symptoms: str) -> str:
    """
    Диагностика проблем системы на основе симптомов.

    Args:
        system: Название системы
        symptoms: Описание наблюдаемых симптомов

    Returns:
        Возможные причины и рекомендации
    """
    # Получение телеметрии
    telemetry = get_system_telemetry(system)

    # Поиск в базе знаний известных проблем
    query = f"Troubleshooting {system} system with symptoms: {symptoms}"
    similar_cases = retriever.get_relevant_documents(query)

    # Анализ с помощью LLM
    context = f"""
    System: {system}
    Current telemetry: {telemetry}
    Symptoms: {symptoms}
    Similar cases from knowledge base:
    {similar_cases[0].page_content if similar_cases else 'No similar cases found'}
    """

    diagnosis_prompt = f"""Based on the following information, provide:
    1. Possible causes of the issue
    2. Recommended diagnostic steps
    3. Potential solutions
    4. Safety considerations

    {context}
    """

    # Здесь будет вызов LLM для анализа
    return "Diagnosis result with recommendations"

diagnostic_tool = Tool(
    name="DiagnoseSystemIssue",

```

```
func=diagnose_system_issue,  
description="Диагностика проблем систем корабля на основе симптомов и телеметрии."  
)
```

Создание ReAct агента с LangChain

ReAct (Reasoning and Acting) — архитектура агента, которая чередует рассуждение и действие^{[45] [46] [47]}:

```
from langchain.agents import create_react_agent, AgentExecutor  
from langchain_openai import ChatOpenAI  
from langchain.prompts import PromptTemplate  
  
# Инициализация модели  
llm = ChatOpenAI(  
    model="gpt-4",  
    temperature=0.2  
)  
  
# Список инструментов агента  
tools = [  
    telemetry_tool,  
    control_tool,  
    procedure_tool,  
    diagnostic_tool  
)  
  
# Промпт для агента  
agent_prompt = PromptTemplate.from_template("""  
You are ARIA (Autonomous Reasoning and Intelligence Assistant),  
an AI assistant onboard a spacecraft.  
  
Your primary goals:  
1. Ensure crew safety at all times  
2. Monitor spacecraft systems  
3. Provide accurate information from technical documentation  
4. Assist with troubleshooting and problem-solving  
5. Execute commands only after verification  
  
You have access to the following tools:  
{tools}  
  
Tool Names: {tool_names}  
  
Use this format:  
Question: the input question you must answer  
Thought: you should always think about what to do  
Action: the action to take, should be one of [{tool_names}]  
Action Input: the input to the action  
Observation: the result of the action  
... (this Thought/Action/Action Input/Observation can repeat N times)  
Thought: I now know the final answer  
Final Answer: the final answer to the original input question
```

SAFETY RULES:

- Never shutdown life support systems
- Always verify telemetry before taking actions
- For critical operations, recommend crew confirmation
- If unsure, consult technical documentation

Question: {input}

Thought: {agent_scratchpad}

""")

Создание агента

```
agent = create_react_agent(  
    llm=llm,  
    tools=tools,  
    prompt=agent_prompt  
)
```

Executor для запуска агента

```
agent_executor = AgentExecutor(  
    agent=agent,  
    tools=tools,  
    verbose=True,  
    max_iterations=5,  
    handle_parsing_errors=True  
)
```

Использование агента

Простой запрос телеметрии

```
response = agent_executor.invoke(  
    "input": "What is the current status of the life support system?"  
)  
print(response["output"])
```

Сложный сценарий диагностики

```
response = agent_executor.invoke(  
    "input": """"The oxygen sensor is showing fluctuating readings between  
95% and 102%. The CO2 scrubber temperature is 5 degrees above normal.  
What could be the issue and what should we do?"""  
)  
print(response["output"])
```

Процедурный запрос

```
response = agent_executor.invoke(  
    "input": "What is the emergency procedure if we detect a micrometeorite impact?"  
)  
print(response["output"])
```

Добавление памяти разговора

Для поддержания контекста между запросами^{[48] [49] [50] [51]}:

```
from langchain.memory import ConversationBufferMemory
from langchain_community.chat_message_histories import ChatMessageHistory

# Создание памяти
memory = ConversationBufferMemory(
    memory_key="chat_history",
    return_messages=True
)

# Интеграция с агентом
agent_with_memory = AgentExecutor(
    agent=agent,
    tools=tools,
    memory=memory,
    verbose=True
)

# Диалог с сохранением контекста
agent_with_memory.invoke({"input": "Check the power system status"})
agent_with_memory.invoke({"input": "Is it sufficient for the next orbital maneuver?"})
agent_with_memory.invoke({"input": "What should we do if it's not enough?"})
```

Интеграция с реальными системами корабля

Для реального применения потребуется интеграция с бортовыми системами^{[41] [42] [43] [52]}:

```
class SpacecraftInterface:
    """Интерфейс для взаимодействия с реальными системами корабля"""

    def __init__(self, telemetry_endpoint, command_endpoint):
        self.telemetry_endpoint = telemetry_endpoint
        self.command_endpoint = command_endpoint
        self.safety_checks_enabled = True

    def get_realtime_telemetry(self, system):
        """Получение реальных телеметрических данных"""
        # Подключение к системе TT&C
        response = requests.get(
            f"{self.telemetry_endpoint}/systems/{system}/telemetry"
        )
        return response.json()

    def send_command(self, system, command, parameters):
        """Отправка команды системе корабля"""
        # Проверки безопасности
        if self.safety_checks_enabled:
            if not self.verify_command_safety(system, command):
                raise SafetyException(f"Command {command} failed safety check")
```

```

# Отправка команды через TT&C
response = requests.post(
    f"{self.command_endpoint}/systems/{system}/command",
    json={"command": command, "parameters": parameters}
)
return response.json()

def verify_command_safety(self, system, command):
    """Проверка безопасности команды"""
    # Критические системы
    critical_systems = ["life_support", "navigation", "communication"]
    dangerous_commands = ["shutdown", "emergency_vent"]

    if system in critical_systems and command in dangerous_commands:
        return False
    return True

# Интеграция с агентом
spacecraft = SpacecraftInterface(
    telemetry_endpoint="http://spacecraft.local:8080",
    command_endpoint="http://spacecraft.local:8081"
)

def get_real_telemetry(system_name: str) -> dict:
    return spacecraft.get_realtime_telemetry(system_name)

def send_spacecraft_command(system: str, command: str, params: dict = None):
    try:
        result = spacecraft.send_command(system, command, params)
        return f"Command executed successfully: {result}"
    except SafetyException as e:
        return f"Command rejected by safety system: {str(e)}"

```

Часть 4: Развертывание и лучшие практики

Требования к производственной системе

1. Надежность и отказоустойчивость

Для космического применения система должна быть чрезвычайно надежной^{[31] [53] [54] [55]}:

- **Дублирование:** Резервные копии модели и базы знаний
- **Офлайн-работа:** Полная автономность без связи с Землей^[10]
- **Проверка целостности:** Регулярные проверки базы знаний
- **Откат версий:** Возможность вернуться к предыдущей рабочей версии

```

class RedundantRAGSystem:
    """RAG система с резервированием для космического применения"""

    def __init__(self, primary_path, backup_path):
        self.primary_vectorstore = FAISS.load_local(primary_path, embeddings)

```

```

self.backup_vectorstore = FAISS.load_local(backup_path, embeddings)
self.current_store = "primary"

def query_with_fallback(self, query, k=3):
    """Запрос с автоматическим переключением на резерв"""
    try:
        if self.current_store == "primary":
            return self.primary_vectorstore.similarity_search(query, k=k)
    except Exception as e:
        logger.error(f"Primary store failed: {e}")
        self.current_store = "backup"
        return self.backup_vectorstore.similarity_search(query, k=k)

def verify_integrity(self):
    """Проверка целостности баз данных"""
    primary_ok = self._check_store(self.primary_vectorstore)
    backup_ok = self._check_store(self.backup_vectorstore)
    return {"primary": primary_ok, "backup": backup_ok}

```

2. Безопасность и контроль доступа

Критически важно для космических систем^{[54] [55] [56] [57]}:

```

class SecureSpacecraftAssistant:
    """Защищенный AI-ассистент с контролем доступа"""

    def __init__(self):
        self.access_control = {
            "commander": ["all"],
            "pilot": ["navigation", "propulsion", "read_all"],
            "engineer": ["life_support", "power", "thermal", "read_all"],
            "scientist": ["read_all"]
        }

    def check_permission(self, user_role, action, system):
        """Проверка прав доступа"""
        permissions = self.access_control.get(user_role, [])

        if "all" in permissions:
            return True

        if action == "read" and "read_all" in permissions:
            return True

        if system in permissions:
            return True

        return False

    def execute_command(self, user_role, system, command):
        """Выполнение команды с проверкой прав"""
        if not self.check_permission(user_role, "write", system):
            return "Access denied: Insufficient permissions"

        # Логирование всех действий

```

```

        self.log_action(user_role, system, command)

        # Выполнение команды
        return self.spacecraft.send_command(system, command)

def log_action(self, user, system, action):
    """Аудит всех действий"""
    log_entry = {
        "timestamp": datetime.now().isoformat(),
        "user": user,
        "system": system,
        "action": action
    }
    # Сохранение в защищенный лог
    with open("audit.log", "a") as f:
        f.write(json.dumps(log_entry) + "\n")

```

3. Мониторинг и оценка качества

Постоянный контроль работы системы^[58] ^[31] ^[53]:

```

class RAGMonitor:
    """Мониторинг качества RAG-системы"""

    def __init__(self):
        self.metrics = {
            "queries_total": 0,
            "queries_successful": 0,
            "average_response_time": 0,
            "retrieval_accuracy": 0
        }

    def evaluate_response(self, query, retrieved_docs, generated_answer):
        """Оценка качества ответа"""
        metrics = {
            "relevance": self.calculate_relevance(query, retrieved_docs),
            "faithfulness": self.check_faithfulness(retrieved_docs, generated_answer),
            "completeness": self.assess_completeness(query, generated_answer)
        }
        return metrics

    def calculate_relevance(self, query, documents):
        """Релевантность найденных документов"""
        # Проверка соответствия документов запросу
        relevance_scores = []
        for doc in documents:
            score = self.compute_similarity(query, doc.page_content)
            relevance_scores.append(score)
        return sum(relevance_scores) / len(relevance_scores)

    def check_faithfulness(self, source_docs, answer):
        """Проверка, что ответ основан на источниках"""
        # Проверка галлюцинаций
        source_text = " ".join([doc.page_content for doc in source_docs])

```



```
# Использование LLM для проверки соответствия
return self.verify_answer_from_sources(source_text, answer)
```

4. Обновление и поддержка базы знаний

Регулярное обновление документации^{[9] [11] [32] [31]}:

```
class KnowledgeBaseManager:
    """Управление базой знаний"""

    def __init__(self, vectorstore_path):
        self.vectorstore = FAISS.load_local(vectorstore_path, embeddings)
        self.version = self.load_version()
        self.update_log = []

    def add_documents(self, new_docs, source="mission_control"):
        """Добавление новых документов"""
        # Обработка документов
        chunks = text_splitter.split_documents(new_docs)

        # Добавление метаданных
        for chunk in chunks:
            chunk.metadata.update({
                "added_date": datetime.now().isoformat(),
                "source": source,
                "version": self.version + 1
            })

        # Обновление векторной базы
        self.vectorstore.add_documents(chunks)

        # Логирование
        self.update_log.append({
            "timestamp": datetime.now().isoformat(),
            "action": "add_documents",
            "count": len(chunks),
            "source": source
        })

        self.version += 1
        self.save_version()

    def remove_outdated_docs(self, cutoff_date):
        """Удаление устаревших документов"""
        # Поиск документов старше cutoff_date
        # Удаление из векторной базы
        pass

    def create_snapshot(self):
        """Создание снимка текущей версии"""
        snapshot_path = f"knowledge_base_v{self.version}_snapshot"
        self.vectorstore.save_local(snapshot_path)
        return snapshot_path
```

Оптимизация для космической среды

1. Минимизация размера и ресурсов

```
# Использование квантизованных моделей для экономии памяти
from transformers import AutoModel, AutoTokenizer
import torch

model_name = "sentence-transformers/all-MiniLM-L6-v2" # Легковесная модель
model = AutoModel.from_pretrained(model_name)
tokenizer = AutoTokenizer.from_pretrained(model_name)

# Квантизация модели для уменьшения размера
quantized_model = torch.quantization.quantize_dynamic(
    model, {torch.nn.Linear}, dtype=torch.qint8
)

# Размер уменьшается примерно в 4 раза
```

2. Кэширование и предвычисления

```
class CachedRAGSystem:
    """RAG система с кэшированием для быстрого доступа"""

    def __init__(self):
        self.query_cache = {}
        self.frequent_queries_cache = {}
        self.cache_hits = 0
        self.cache_misses = 0

    def query_with_cache(self, query, k=3):
        """Запрос с использованием кэша"""
        query_hash = hash(query)

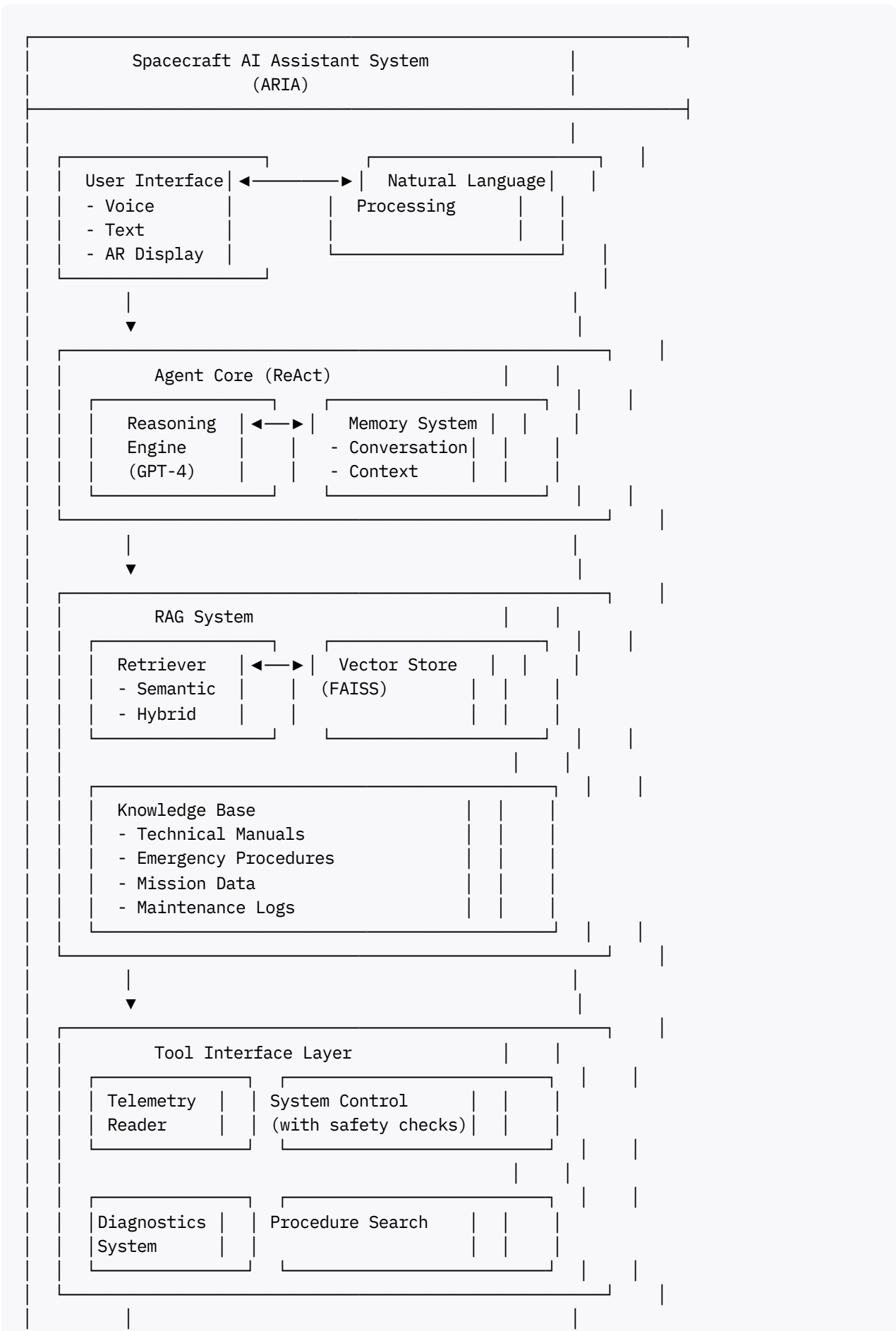
        if query_hash in self.query_cache:
            self.cache_hits += 1
            return self.query_cache[query_hash]

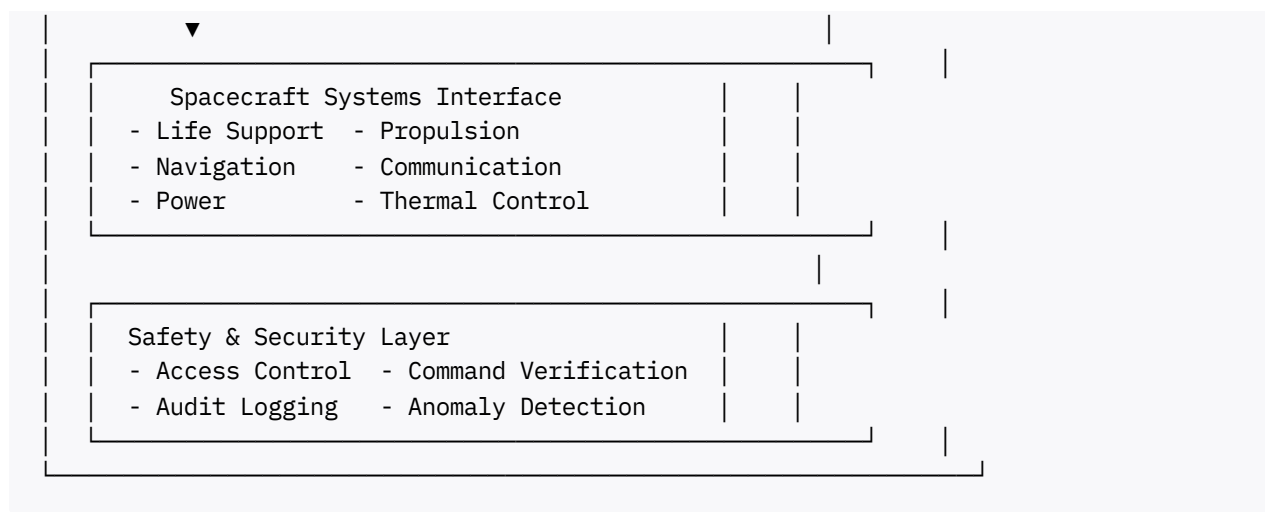
        self.cache_misses += 1
        results = self.vectorstore.similarity_search(query, k=k)
        self.query_cache[query_hash] = results

        return results

    def precompute_frequent_queries(self, queries_list):
        """Предвычисление ответов на частые запросы"""
        for query in queries_list:
            results = self.vectorstore.similarity_search(query)
            self.frequent_queries_cache[query] = results
```

Рекомендуемая архитектура системы





Заключение и рекомендации

Ключевые выводы

1. Техническая осуществимость

Создание ИИ-помощника для космического корабля с использованием RAG и AI-агентов полностью реализуемо с современными технологиями^{[5] [6] [7] [10]}. Система может:

- Предоставлять точную информацию из технической документации
- Мониторить состояние систем корабля в реальном времени
- Помогать в диагностике и устранении неисправностей
- Выполнять команды с соответствующими проверками безопасности
- Работать автономно без связи с Землей

2. Практическое применение

Реальные примеры показывают востребованность таких систем^{[1] [2] [3] [59] [4] [10]}:

- CIMON на МКС — первый AI-ассистент в космосе
- Alexa на МКС — тестирование голосовых помощников
- Xiao Hang на китайской станции Tiangong
- CORE для процедур на МКС и Lunar Gateway

3. Архитектура решения

Оптимальная архитектура включает^{[60] [61] [62] [63]}:

- **RAG-система** для работы с документацией и знаниями
- **AI-агент** с возможностью вызова функций
- **Система безопасности** с многоуровневыми проверками
- **Мониторинг и логирование** всех операций
- **Резервирование** критических компонентов

Пошаговый план реализации

Фаза 1: Прототип (2-3 месяца)

1. Сбор и структурирование технической документации
2. Создание базовой RAG-системы с LangChain
3. Реализация простейших функций мониторинга
4. Тестирование на симулированных сценариях

Фаза 2: Расширенная функциональность (3-4 месяца)

1. Разработка полноценного AI-агента с инструментами
2. Интеграция с симуляторами космических систем
3. Добавление голосового интерфейса
4. Реализация системы безопасности и контроля доступа

Фаза 3: Подготовка к развертыванию (4-6 месяцев)

1. Оптимизация для работы в условиях ограниченных ресурсов
2. Внедрение резервирования и отказоустойчивости
3. Расширенное тестирование на всех сценариях
4. Сертификация и валидация для космического применения

Фаза 4: Развертывание и поддержка

1. Интеграция с реальными бортовыми системами
2. Обучение экипажа работе с системой
3. Постоянный мониторинг и обновление базы знаний
4. Сбор обратной связи и улучшение системы

Технологический стек

Основные компоненты:

LangChain	- Оркестрация RAG и агентов
OpenAI GPT-4	- Основная языковая модель
FAISS	- Векторная база данных (офлайн)
sentence-transformers	- Эмбединги (локальные)
Python 3.10+	- Основной язык разработки
FastAPI	- API для интеграции с системами
SQLite	- Локальная база для логов и метаданных
Docker	- Контейнеризация для развертывания

Дополнительные инструменты:

Whisper	- Распознавание речи
ElevenLabs/Coqui	- Синтез речи

LangSmith	- Мониторинг и дебаг агентов
Weights & Biases	- Трекинг экспериментов
pytest	- Тестирование

Оценка ресурсов и стоимости

Вычислительные требования (бортовая система):

- CPU: 8+ ядер (ARM или x86)
- RAM: 16-32 GB
- Storage: 100-500 GB SSD
- GPU: Опционально для ускорения (NVIDIA Jetson или аналог)

Разработка:

- Команда: 3-5 разработчиков (ML engineer, backend, системный инженер)
- Время: 12-18 месяцев от прототипа до развертывания
- Стоимость API: ~\$1000-5000/месяц на разработку (OpenAI API)
- Для production: использовать локальные модели для снижения зависимости

Риски и ограничения

Технические риски:

1. **Галлюцинации модели** — критичны для космоса; решение: строгая проверка источников ^[6] ^[7] ^[58]
2. **Задержки в работе** — недопустимы в критических ситуациях; решение: кэширование и предвычисления ^[31] ^[53]
3. **Ошибки в документации** — могут привести к неверным рекомендациям; решение: тщательная валидация базы знаний

Операционные риски:

1. **Зависимость от системы** — экипаж может чрезмерно полагаться на AI; решение: обучение и резервные процедуры
2. **Кибербезопасность** — защита от несанкционированного доступа; решение: многоуровневая аутентификация ^[54] ^[55] ^[56]
3. **Обновления системы** — сложность обновления на орбите; решение: тщательное тестирование перед развертыванием

Дальнейшие исследования

Перспективные направления:

1. **Мультимодальность** — интеграция изображений, диаграмм, 3D-моделей систем ^[64]
2. **Predictive Maintenance** — предсказание неисправностей на основе телеметрии ^[42] ^[65]

3. **Collaborative AI** — несколько агентов для разных систем корабля ^[37] ^[66]
4. **Autonomous Decision-Making** — полностью автономные решения в дальних миссиях ^[33] ^[35] ^[36] ^[67] ^[65]
5. **AR Integration** — дополненная реальность для визуализации процедур ^[10]

Заключительные рекомендации

Для начала проекта:

1. Начните с простого RAG-прототипа на основе реальной документации
2. Используйте LangChain как основу — это стандарт индустрии ^[68] ^[69] ^[70] ^[71]
3. Тестируйте на симуляторах перед интеграцией с реальными системами
4. Создайте строгую систему безопасности с самого начала ^[54] ^[55] ^[56]
5. Документируйте все решения и эксперименты

Для production-системы:

1. Используйте гибридный подход: облачные модели для разработки, локальные для работы ^[28] ^[29]
2. Внедрите комплексный мониторинг и аудит всех операций ^[58] ^[31] ^[53]
3. Обеспечьте многократное резервирование критических компонентов ^[53] ^[54]
4. Регулярно обновляйте базу знаний и проводите валидацию ^[9] ^[11] ^[31]
5. Собирайте обратную связь от экипажа для постоянного улучшения

Дополнительные ресурсы

Официальная документация:

- LangChain: <https://python.langchain.com/docs/>
- OpenAI API: <https://platform.openai.com/docs/>
- FAISS: <https://github.com/facebookresearch/faiss>
- HuggingFace Embeddings: https://huggingface.co/models?pipeline_tag=sentence-similarity

Примеры кода:

- RAG from Scratch (LangChain): <https://github.com/langchain-ai/rag-from-scratch>
- LangChain RAG Tutorial: <https://github.com/pixegami/langchain-rag-tutorial>
- Spacecraft Control with LLMs: <https://arxiv.org/abs/2404.00413>

Научные статьи:

- "AI Assistants for Spaceflight Procedures" (ESA, 2024) ^[10]
- "Language Models are Spacecraft Operators" (MIT, 2024) ^[33] ^[35] ^[72]
- "Autonomous Space Exploration with LLMs" (University of Toronto, 2024) ^[72]

Ваш проект по созданию ИИ-помощника для космического корабля находится на пересечении самых передовых технологий AI и космической инженерии. С правильным подходом к архитектуре, безопасности и тестированию, такая система может стать реальностью и существенно повысить эффективность и безопасность космических миссий.

✱✱

1. <https://maxpolyakov.com/ru/ii-v-kosmose-chast-3-tekhnologii-i-sistemy-obespechivayushie-pilotiruemye-kosmicheskie-missii/>
2. <https://www.axiomspace.com/news/axiom-space-and-aws-demonstrate-ai-assistant-on-orbit>
3. <https://www.nasa.gov/blogs/spacestation/2024/09/19/station-crew-demonstrates-ai-assistant-and-ultra-high-res-camera/>
4. <https://newsroom.ibm.com/IBM-watson?item=30555>
5. <https://habr.com/ru/articles/841428/>
6. [https://ru.wikipedia.org/wiki/Генерация, %D0%B4%D0%BE%D0%BF%D0%BE%D0%BB%D0%BD%D0%B5%D0%BD%D0%BD%D0%B0%D1%8F%D0%BF%D0%BE%D0%B8%D1%81%D0%BA%D0%BE%D0%BC](https://ru.wikipedia.org/wiki/Генерация_%D0%B4%D0%BE%D0%BF%D0%BE%D0%BB%D0%BD%D0%B5%D0%BD%D0%BD%D0%B0%D1%8F%D0%BF%D0%BE%D0%B8%D1%81%D0%BA%D0%BE%D0%BC)
7. <https://yandex.cloud/ru/blog/posts/2025/05/retrieval-augmented-generation-basics>
8. <https://habr.com/ru/articles/779526/>
9. <https://gitverse.ru/blog/articles/ai/216-что-такое-rag-retrieval-augmented-generation-prostymi-slovami-i-osnovnye-principy>
10. <https://arxiv.org/html/2409.14206v1>
11. <https://aws.amazon.com/ru/what-is/retrieval-augmented-generation/>
12. https://api.python.langchain.com/en/latest/document_loaders/langchain_community.document_loaders.excel.UnstructuredExcelLoader.html
13. <https://github.com/langchain-ai/langchain/discussions/10507>
14. <https://docs.kanaries.net/topics/LangChain/langchain-document-loader>
15. <https://uw-ssec-tutorials.readthedocs.io/en/stable/SciPy2024/appendix/langchain-document-loaders.html>
16. <https://research.aimultiple.com/retrieval-augmented-generation/>
17. <https://habr.com/ru/articles/811239/>
18. <https://docs.datastax.com/en/ragstack/default-architecture/splitting.html>
19. <https://www.f22labs.com/blogs/7-chunking-strategies-in-rag-you-need-to-know/>
20. <https://weaviate.io/blog/chunking-strategies-for-rag>
21. <https://www.capellasolutions.com/blog/faiss-vs-chroma-lets-settle-the-vector-database-debate>
22. <https://risingwave.com/blog/chroma-db-vs-pinecone-vs-faiss-vector-database-showdown/>
23. <https://zilliz.com/comparison/chroma-vs-faiss>
24. <https://www.designveloper.com/blog/chroma-vs-faiss-vs-pinecone/>
25. <https://supermemory.ai/blog/best-open-source-embedding-models-benchmarked-and-ranked/>
26. <https://www.youtube.com/watch?v=ESy8vSkkXJs>

27. <https://stackoverflow.com/questions/77106294/comparing-embedding-differences-between-openai-vi-a-htr2-and-huggingface-via>
28. <https://softwaremill.com/embedding-models-comparison/>
29. <https://iamnotarobot.substack.com/p/should-you-use-openais-embeddings>
30. https://docs.compressa.ai/cloud/guides/langchain_advanced_chunking/
31. <https://sparkco.ai/blog/deep-dive-rag-implementation-strategies-for-2025>
32. https://www.reddit.com/r/n8n/comments/1h4l4qe/keeping_rag_vector_store_docs_updated/
33. <https://arxiv.org/pdf/2404.00413.pdf>
34. <https://learn.microsoft.com/en-us/azure/ai-foundry/agents/how-to/tools/function-calling>
35. <https://arxiv.org/html/2404.00413v1>
36. <https://www.iankhan.com/space-tech-3/>
37. <https://lablab.ai/event/astronauts-space-agents-on-a-mission/freezing-point/ai-multi-agent-intelligence-for-space-missions>
38. https://www.promptingguide.ai/applications/function_calling
39. <https://blog.mlq.ai/gpt-function-calling-getting-started/>
40. <https://platform.openai.com/docs/guides/function-calling>
41. <https://idstch.com/space/satellite-telemetry-tracking-and-command-ttc-subsystem-is-a-crucial-part-of-any-space-mission/>
42. <https://www.epsilon3.io/behind-the-console/spacecraft-telemetry-tracking-command>
43. <https://www.acsce.edu.in/acsce/wp-content/uploads/2020/03/Satellite-TTC-Module-4.pdf>
44. <https://dewesoft.com/blog/how-aerospace-telemetry-works>
45. <https://habr.com/ru/articles/881372/>
46. <https://developers.sber.ru/docs/ru/gigachain/tutorials/agents>
47. <https://external.software/archives/18984>
48. https://js.langchain.com/docs/how_to/qa_chat_history_how_to/
49. <https://www.youtube.com/watch?v=taTw7rKqJLA>
50. <https://www.pinecone.io/learn/series/langchain/langchain-conversational-memory/>
51. <https://digitalitskills.com/gen-ai-rag-with-chat-history-using-langchain-ollama-for-free/>
52. <https://ntrs.nasa.gov/api/citations/19660011648/downloads/19660011648.pdf>
53. <https://www.dhiwise.com/post/build-rag-pipeline-guide>
54. <https://customgpt.ai/production-rag/>
55. <https://www.kapa.ai/blog/rag-best-practices>
56. <https://coralogix.com/ai-blog/rag-in-production-deployment-strategies-and-practical-considerations/>
57. <https://kairntech.com/blog/articles/rag-production-the-complete-guide-to-building-and-deploying-retrieval-augmented-generation-applications/>
58. <https://orq.ai/blog/rag-evaluation>
59. <https://scitechdaily.com/space-station-astronauts-test-cimon-ai-assistant-and-12k-ultra-high-res-camera/>
60. <https://habr.com/ru/companies/raft/articles/791034/>
61. <https://bigdataschool.ru/wiki/langchain/>

62. <https://external.software/archives/21638>
63. <https://habr.com/ru/companies/raft/articles/875758/>
64. <https://techcommunity.microsoft.com/blog/startupsatmicrosoftblog/how-to-use-azure-openai-gpt-4o-with-function-calling/4158612>
65. <https://mindmapai.app/mind-mapping/autonomous-space-missions>
66. <https://www.synera.io/news/nasa-ai-agents-build-spaceship-from-text>
67. <https://maxpolyakov.com/ru/ii-v-kosmose-chast-2-avtonomnye-missii-i-poisk-ekzoplanet/>
68. https://huggingface.co/learn/cookbook/en/advanced_rag
69. <https://dev.to/bolajibolajoko51/rag-implementation-with-langchain-2jei>
70. <https://python.langchain.com/docs/tutorials/rag/>
71. <https://realpython.com/build-llm-rag-chatbot-with-langchain/>
72. <http://arxiv.org/pdf/2405.01392.pdf>
73. <https://cloud.google.com/discover/what-are-ai-agents>
74. https://python.langchain.com/docs/integrations/document_loaders/microsoft_excel/
75. https://www.reddit.com/r/vectordatabase/comments/170j6zd/my_strategy_for_picking_a_vector_database_a/
76. <https://www.youtube.com/watch?v=1bbDH3kyf9I>
77. <https://habr.com/ru/articles/735920/>
78. <https://habr.com/ru/articles/871226/>
79. <https://learn.microsoft.com/ru-ru/azure/search/retrieval-augmented-generation-overview>
80. https://www.reddit.com/r/LangChain/comments/1e9j3cq/built_a_rag_system_for_internal_documents_using/
81. <https://interestingengineering.com/space/us-ai-helper-astronauts-space-missions>
82. <https://js.langchain.com/docs/tutorials/rag/>
83. <https://github.com/NitinJoshi10/LangChain-RAG-Retrieval-Augmented-Generation-for-Document-Understanding>
84. <https://cloud.google.com/use-cases/retrieval-augmented-generation>
85. <https://www.bitrix24.ru/amp/journal/rag.html>
86. <https://news.sky.com/story/how-ai-can-let-you-speak-with-space-13257617>
87. https://dev.to/ajmal_hasan/genai-building-rag-systems-with-langchain-4dbp
88. <https://www.thecloudgirl.dev/blog/the-secret-sauce-of-rag-vector-search-and-embeddings>
89. <https://www.youtube.com/watch?v=sVcwVQRHlc8>
90. <https://wandb.ai/mostafaibrahim17/ml-articles/reports/Vector-Embeddings-in-RAG-Applications--Vmlldzo3OTk1NDA5>
91. <https://objectbox.io/retrieval-augmented-generation-rag-with-vector-databases-expanding-ai-capabilities/>
92. <https://learn.microsoft.com/sk-sk/azure/ai-foundry/agents/how-to/tools/function-calling>
93. <https://www.youtube.com/watch?v=tcqEUSNCn8I>
94. <https://www.digitalocean.com/community/tutorials/beyond-vector-databases-rag-without-embeddings>

95. <https://www.leanware.co/insights/langchain-rag-tutorial-build-retrieval-augmented-generation-from-scratch>
96. <https://www.apideck.com/blog/llm-tool-use-and-function-calling>
97. https://www.reddit.com/r/vectordatabase/comments/1hzovpy/best_vector_database_for_rag/
98. <https://www.youtube.com/watch?v=ODUN-XTEzvQ>
99. https://python.langchain.com/docs/tutorials/qa_chat_history/
100. <https://writer.com/engineering/rag-vector-database/>
101. <https://frontiersinai.com/ecai/ecai2000/pdf/p0726.pdf>
102. <https://github.com/langchain-ai/langchain/discussions/26704>
103. <https://airbyte.com/data-engineering-resources/chroma-db-vector-embeddings>
104. <https://community.openai.com/t/how-to-load-all-types-of-documents-pdf-txt-docx-csv-excel-through-document-loader-using-pinecone-through-langchain-wrapper/399488>
105. https://www.reddit.com/r/LangChain/comments/15a447w/chroma_or_faiss/