# Machine Learning and Reasoning with Built Environment Data
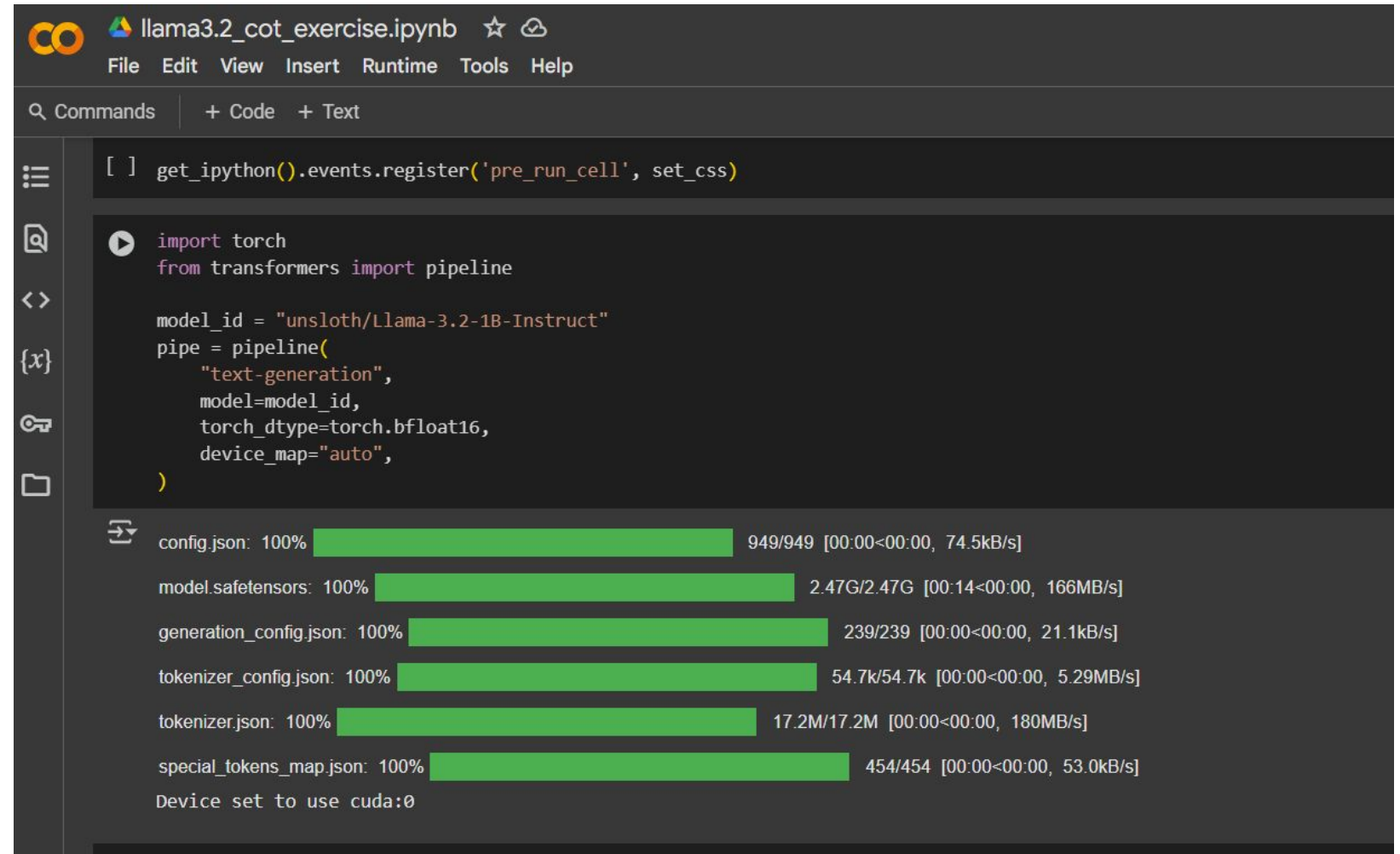
## Work with LLMs

Computing in Construction

Metropolia
University of Applied Sciences

# Update of Earlier Script

**- Instructions of Token**

**settings added**

**- model_id changed,**

**no separate**

**permission needed**

# Limitations

New information later than the cutoff dates off LLMs?

Information not included in the training dataset of LLMs?

Data from local files which are not supposed to share?

Longer texts with more tokens than the context window?

Metropolia
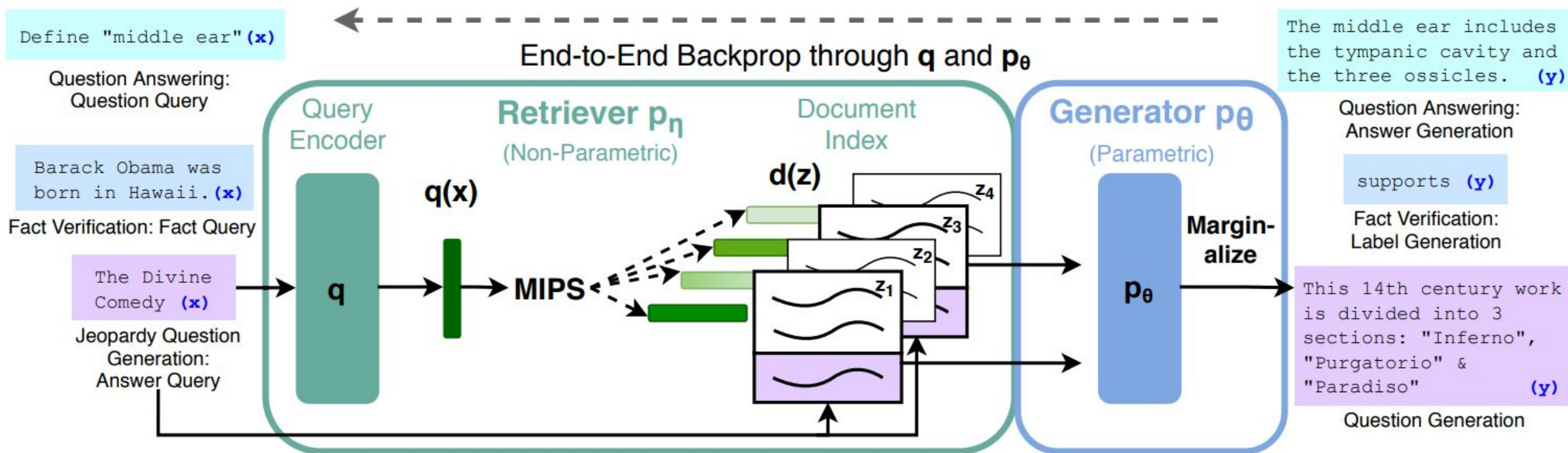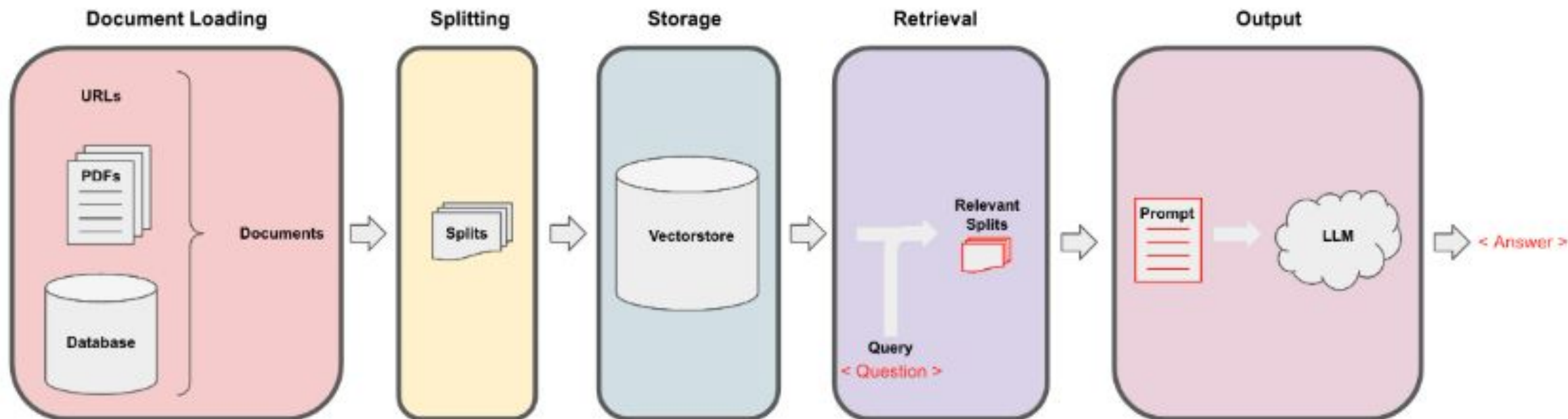University of Applied Sciences

# Retrieval-Augmented Generation (RAG)



Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query $x$, we use Maximum Inner Product Search (MIPS) to find the top-K documents $z_i$. For final prediction $y$, we treat $z$ as a latent variable and marginalize over seq2seq predictions given different documents.

# How does it work?

## Retrieval augmented generation

In retrieval augmented generation (RAG), an LLM retrieves contextual documents from an external dataset as part of its execution.

This is useful if we want to ask question about specific documents (e.g., our PDFs, a set of videos, etc).

# Why Use Open Source LLMs

# Large or Small?

**Authors**

**Rina Diane Caballar**
Staff Writer

# What are small language models?

Small language models (SLMs) are artificial intelligence (AI) models capable of processing, understanding and generating natural language content. As their name implies, SLMs are smaller in scale and scope than large language models (LLMs).

In terms of size, SLM parameters range from a few million to a few billion, as opposed to LLMs with hundreds of billions or even trillions of parameters. Parameters are internal variables, such as weights and biases, that a model learns during training. These parameters influence how a machine learning model behaves and performs.

Small language models are more compact and efficient than their large model counterparts. As such, SLMs require less memory and computational power, making them ideal for resource-constrained environments such as edge devices and mobile apps, or even for scenarios where AI inferencing—when a model generates a response to a user's query—must be done offline without a data network.

Metropolia
University of Applied Sciences

AI FIGHTER

VS

Large Language Models

LLM quantization!?

quantization  performance

inference speed

Airtrain

# Explore Hugging Face

# Popular Python Libraries

**LangChain is an open-source framework designed to help developers build applications powered by large language models (LLMs). It focuses on streamlining the process of creating complex, language-driven workflows.**
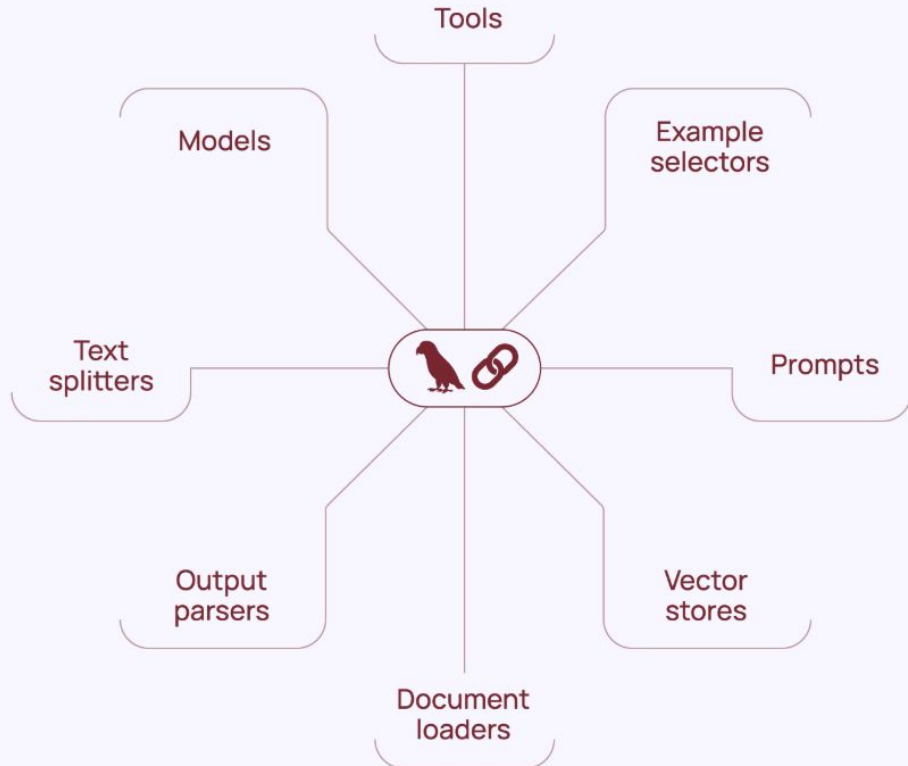
# Popular Python Libraries

**Ollama is a platform designed to simplify running and deploying large language models locally. It focuses on providing a streamlined experience for developers who want to leverage LLMs without relying solely on cloud-based APIs.**

## Get up and running with large language models.

Run Llama 3.3, DeepSeek-R1, Phi-4, Mistral, Gemma 3, and other models, locally.

**Download ↓**

Available for macOS, Linux, and Windows

# Popular Python Libraries

LlamaIndex (formerly known as GPT Index) is a data framework that helps connect unstructured data sources (like documents, databases, and APIs) with large language models (LLMs). It enables efficient retrieval and reasoning over private or structured data, making it a powerful tool for Retrieval-Augmented Generation (RAG) applications.

# Popular Python Libraries

**Haystack is an open-source framework focused on building end-to-end search systems and question-answering (QA) applications using natural language processing techniques and language models.**

# Popular Python Libraries

**Llama.cpp is a lightweight, efficient C++ implementation of Meta's LLaMA (Large Language Model Meta AI) models, designed for running large language models (LLMs) on local devices with minimal dependencies. It enables fast inference on CPUs and GPUs without requiring cloud services or heavy frameworks like PyTorch.**

## llama.cpp



license MIT  Server passing

Roadmap / Project status / Manifesto / ggml

Inference of Meta's LLaMA model (and others) in pure C/C++

> 📣 **Important**

New `llama.cpp` package location: ggml-org/llama.cpp

Update your container URLs to: `ghcr.io/ggml-org/llama.cpp`

## Python Bindings for `llama.cpp`

docs passing  Tests passing  pypi v0.3.8  python 3.8 | 3.9 | 3.10 | 3.11 | 3.12  license MIT  downloads 341k/month
Github Downloads 467k

Simple Python bindings for **@ggerganov's** `llama.cpp` library. This package provides:

- Low-level access to C API via `ctypes` interface.
- High-level Python API for text completion
  - OpenAI-like API
  - LangChain compatibility
  - LlamaIndex compatibility
- OpenAI compatible web server
  - Local Copilot replacement
  - Function Calling support
  - Vision API support
  - Multiple Models

# GGUF

**Running LLMs with CPUs**

1. **Efficiency**: GGUF makes LLMs more compact and faster to load. This is crucial for local deployment, where storage space and RAM might be limited compared to cloud environments. The format uses advanced compression techniques to reduce model size without sacrificing performance.

2. **Compatibility**: GGUF improves how LLMs work across different platforms and devices. It provides a standardized way to package model weights, architecture information, and metadata, making it easier for various software to interpret and use the model consistently.

3. **Local Deployment**: It's particularly useful for running LLMs on personal computers or local servers. GGUF's optimizations allow even large models to run on consumer-grade hardware, democratizing access to powerful AI capabilities.

4. **Customization**: GGUF allows for easy fine-tuning and modification of models. Users can adjust parameters, add custom tokens, or modify model behavior without needing to re-train the entire model from scratch.

# Hands-on Practice

## Transformers and Colab AI Features

# Python with LLMs on Local Computer

## Nvidia GPU

# Python with LLMs on Local Computer

## Install CUDA Toolkit

# Python with LLMs on Local Computer

## Check CUDA version

# Python with LLMs on Local Computer

## Install pytorch according to your system and CUDA version

# Python with LLMs on Local Computer

## OpenVINO for CPU, other GPUs or NPU

# Python with LLMs on Local Computer

## Install Transformers

# RAG with LangChain



## RAG from webpage

```python
# Load documents from webpage
urls = ["https://www.metropolia.fi/fi/tutkimus-kehitys-ja-innovaatiot/yhteistyoalustat/smart-lab"]
loader_html = AsyncHtmlLoader(urls)
docs_html = loader_html.load()

bs_transformer = BeautifulSoupTransformer()
docs_transformed = bs_transformer.transform_documents(docs_html, tags_to_extract=["div"])

# Split it into chunks
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=100)
docs_transformed_split = text_splitter.split_documents(docs_transformed)


print(len(docs_transformed_split))
print(len(docs_transformed[0].page_content))
```

```
Fetching pages: 100%|##########| 1/1 [00:00<00:00,  4.03it/s]
```
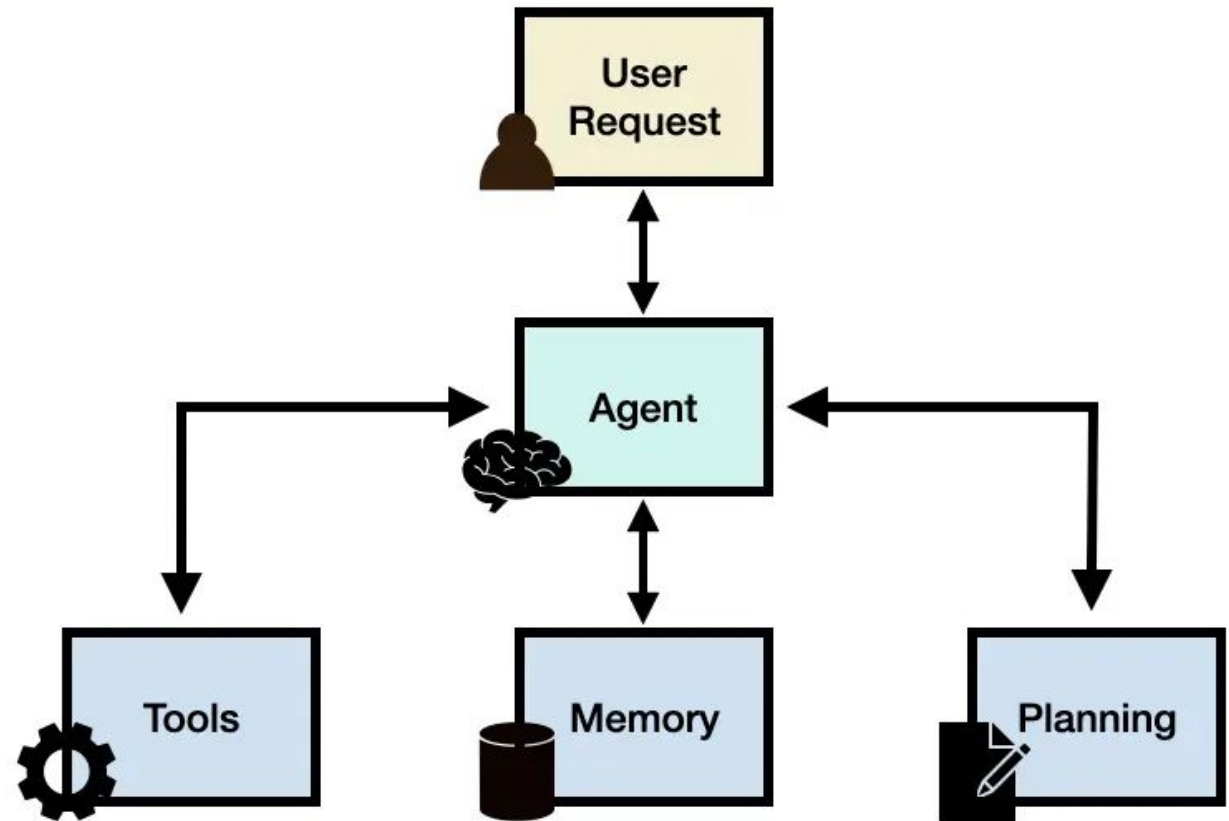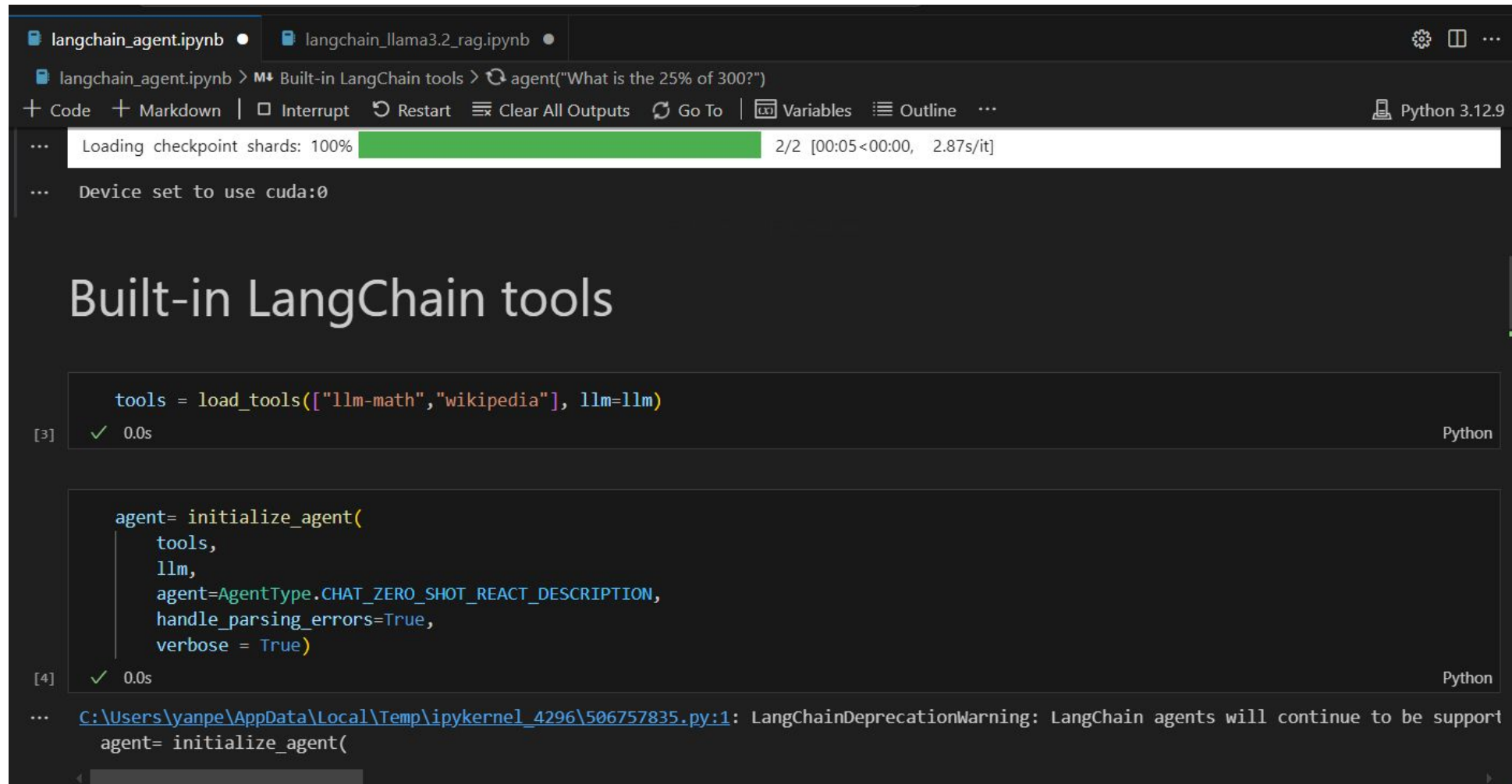
AI Ascent

Andrew
Ng

# LLM Agent Framework

**Generally speaking, an LLM agent framework can consist of the following core components:**

- **User Request** - a user question or request
- **Agent/Brain** - the agent core acting as coordinator
- **Planning** - assists the agent in planning future actions
- **Memory** - manages the agent's past behaviors

# LangChain Agents

# Model Context Protocol (MCP)

## Model Context Protocol

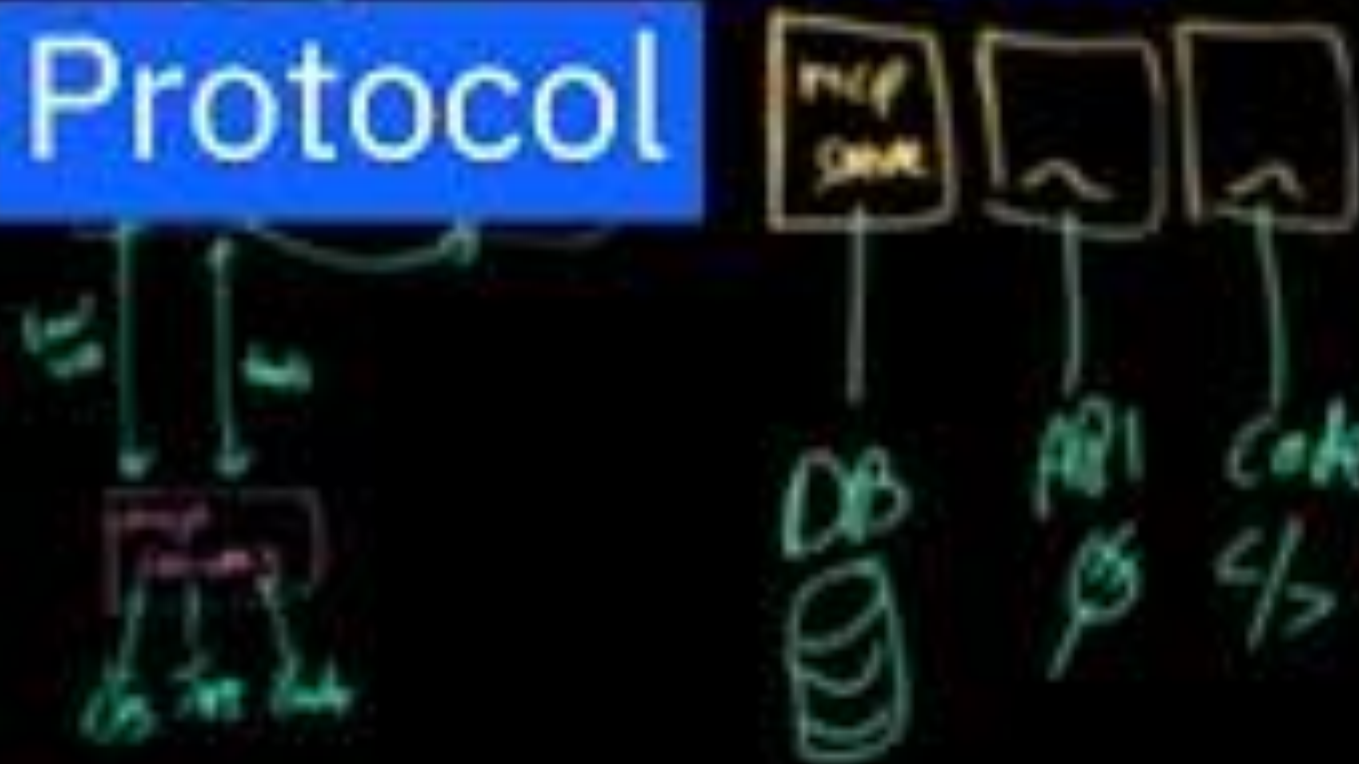The Model Context Protocol is an open standard that enables developers to build secure, two-way connections between their data sources and AI-powered tools. The architecture is straightforward: developers can either expose their data through MCP servers or build AI applications (MCP clients) that connect to these servers.

Today, we're introducing three major components of the Model Context Protocol for developers:

- The Model Context Protocol specification and SDKs
- Local MCP server support in the Claude Desktop apps
- An open-source repository of MCP servers

Metropolia
University of Applied Sciences

How to build a
Model Context
Protocol

# BlenderMCP

README    MIT license

**Contributors**

**Languages**

● Python 100.0%

# BlenderMCP - Blender Model Context Protocol Integration

BlenderMCP connects Blender to Claude AI through the Model Context Protocol (MCP), allowing Claude to directly interact with and control Blender. This integration enables prompt assisted 3D modeling, scene creation, and manipulation.

Full tutorial

## Join the Community

Give feedback, get inspired, and build on top of the MCP: Discord

## Supporters

**Top supporters:**

CodeRabbit

**All supporters:**

Support this project

## Release notes (1.1.0)

- Added support for Poly Haven assets through their API
- Added support to prompt 3D models using Hyper3D Rodin

## Metropolia
University of Applied Sciences

# Thank you!

Metropolia
University of Applied Sciences