



LANGCHAIN ZOOMCAP SESSION 2

Harpreet sahota

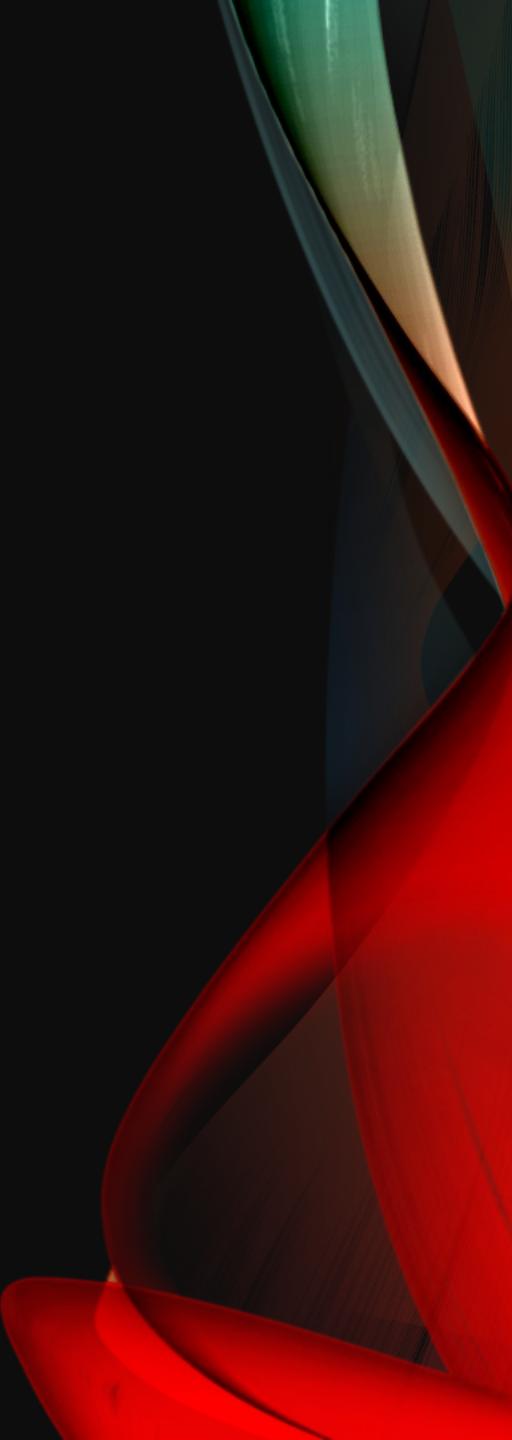


TODAY'S AGENDA

- The LLM landscape
- Prompt Engineering
- What is LangChain?
- Overview of core components in LangChain
 - Model I/O
 - Retrieval
 - Chains
 - Agents
 - Memory



THE LLM LANDSCAPE



OPEN SOURCE LLMS

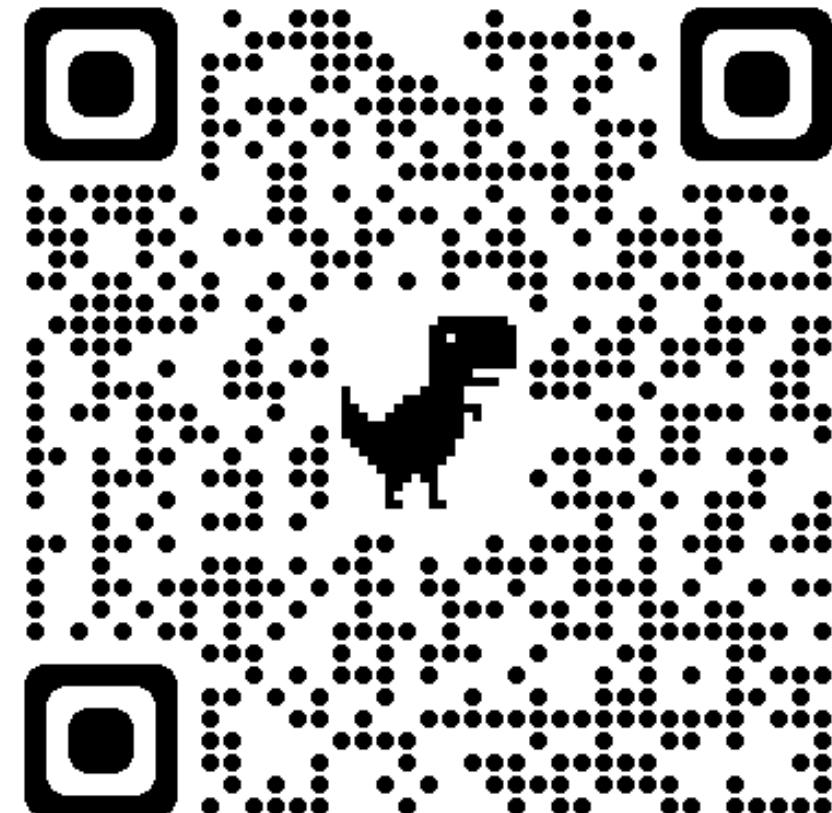
- Democratize access to advanced AI capabilities.
- Publicly accessible source codes for anyone to use, modify, and distribute.
- Fosters collaboration among developers, researchers, and enthusiasts.
- Encourages innovation, knowledge sharing, and collective development efforts.

NOTABLE OPEN SOURCE MODELS

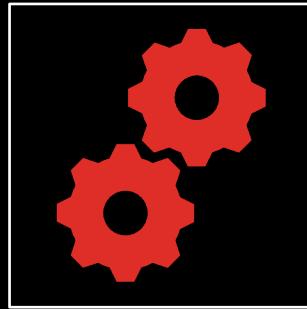
- LLaMA 2: Meta AI's model excels in dialogue situations, trained on extensive publicly available data.
- Alpaca: Known for its proficiency in instruction-following tasks.
- Falcon: Developed by the Technology Innovation Institute, optimized for text generation and chatbot functionality.
- Deci: Fast and efficient models

I've written an in-depth blog about the landscape of open-source LLMs (from Jan – July 2023) which you can read by scanning the following QR code, or visiting:

<https://deci.ai/blog/list-of-large-language-models-in-open-source/>



BENEFITS OF OPEN-SOURCE ECOSYSTEM



Flexibility and customization for tailored models.



Cost efficiency with substantial long-term benefits.



Data privacy through self-hosting, crucial in the digital landscape.



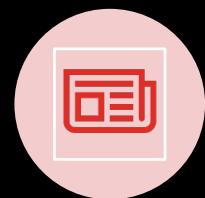
CLOSED-SOURCE LLMS

- Developed and maintained by organizations, the source code is not public.
- Architecture, training data, and algorithms are typically not accessible.
- Often commercial products may require licenses or subscriptions.

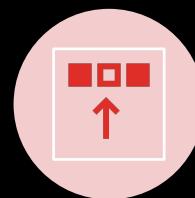
NOTABLE CLOSED SOURCE MODELS

- OpenAI's GPT Series: GPT-3.5 and GPT-4 excel in text completion, translation, and question-answering.
- Anthropic's Claude: Offers advanced NLP capabilities like editing, rewriting, and summarizing.
- Cohere's Command Series: Ideal for chatbots, excelling at interpreting instruction-like prompts.
- Other Closed-Source Models: Models like Jurassic-2 and PaLM making significant impact.

PROS AND CONS OF CLOSED SOURCE MODELS



Easy
Integration



Quick
Deployment



Regular
Updates



Limited
Customization



Data Privacy
Concerns



Scaling Costs

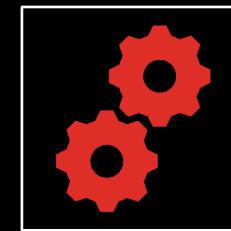
FACTORS FOR MODEL CHOICE



Expertise



Budget



Customization
needs



Data privacy

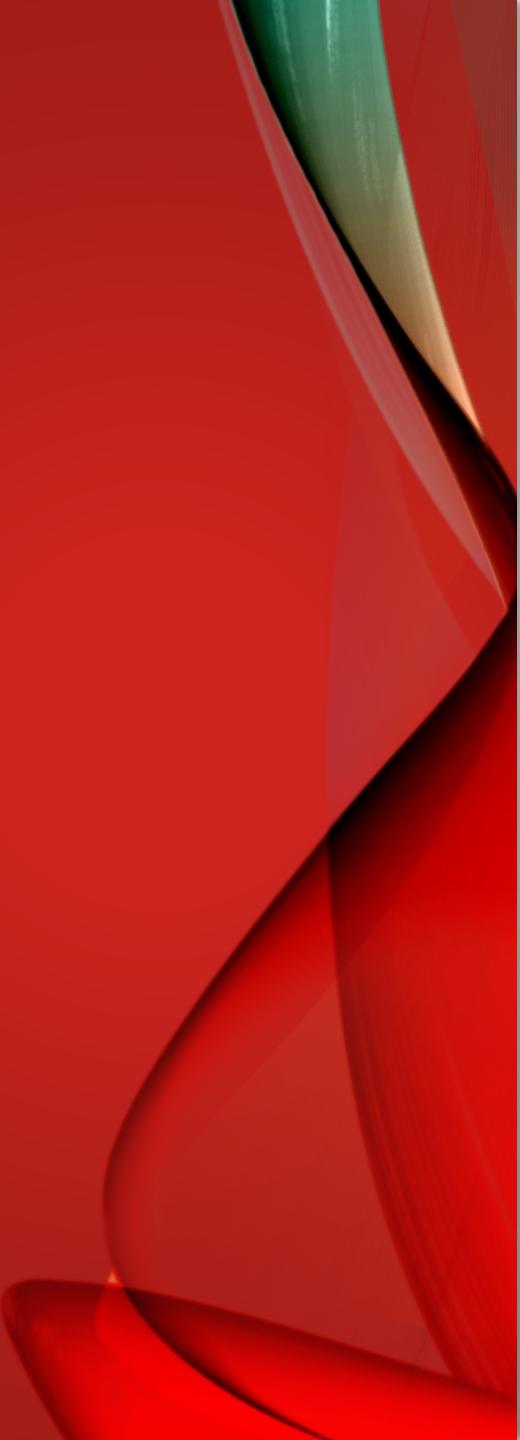
The background of the slide features a scenic landscape of rolling green hills under a dramatic sky. The top left corner shows a vibrant gradient from red to yellow, while the top right corner is a bright cyan color. A dark grey rectangular overlay covers the middle portion of the slide, containing the main title and bullet points.

THE ROAD AHEAD

- LLM Landscape is Dynamic
- Staying Informed is Crucial
- Long-Term Strategy Matters
- Open-Source or Proprietary
- LLMs are Here to Stay

PROMPT ENGINEERING



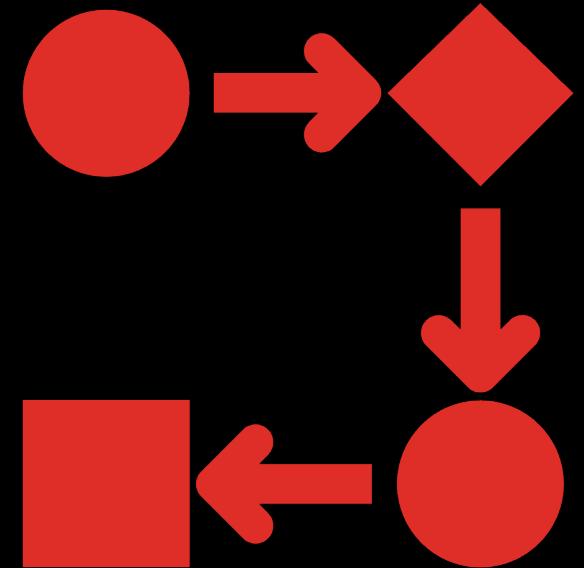


PROMPT ENGINEERING FOR LLMS

- Craft precise prompts for effective LLM guidance.
- Clear instructions improve final results.
- Think of it as guiding a skilled artisan.
- Ensure LLM understands and responds accurately.
- Precision enhances user intent.

THE EVOLUTION OF PROMPT ENGINEERING

- Prompt engineering has evolved significantly.
- Initially focused on formulating questions for desired answers.
- LLM advancements have improved prompting techniques.
- Models like GPT-3 and ChatGPT revolutionized prompt engineering.
- Utilizing zero-shot, one-shot, few-shot examples for effective guidance.
- These examples lead to accurate and diverse LLM responses.



THE IMPORTANCE OF PROMPT ENGINEERING

- Prompt engineering is the key to effective communication.
- It bridges human intent and machine understanding.
- Well-crafted prompts lead to accurate and context-aware LLM responses.
- Poorly designed prompts result in ambiguity or errors.

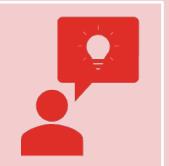
ADVANCED PROMPT ENGINEERING TECHNIQUES



Basic prompts may be insufficient in some cases.



Explore advanced techniques for improved results.



Chain of Thought (CoT) Prompting enhances reasoning.



Integration with external databases improves accuracy.

THE ART AND FUTURE OF PROMPT ENGINEERING



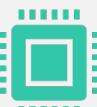
Prompt engineering is both a science and an art.



Deep understanding of model capabilities and desired outcomes is essential.



As LLMs advance, prompt engineering becomes increasingly critical.



LangChain leads with innovation and collaboration in prompt engineering.



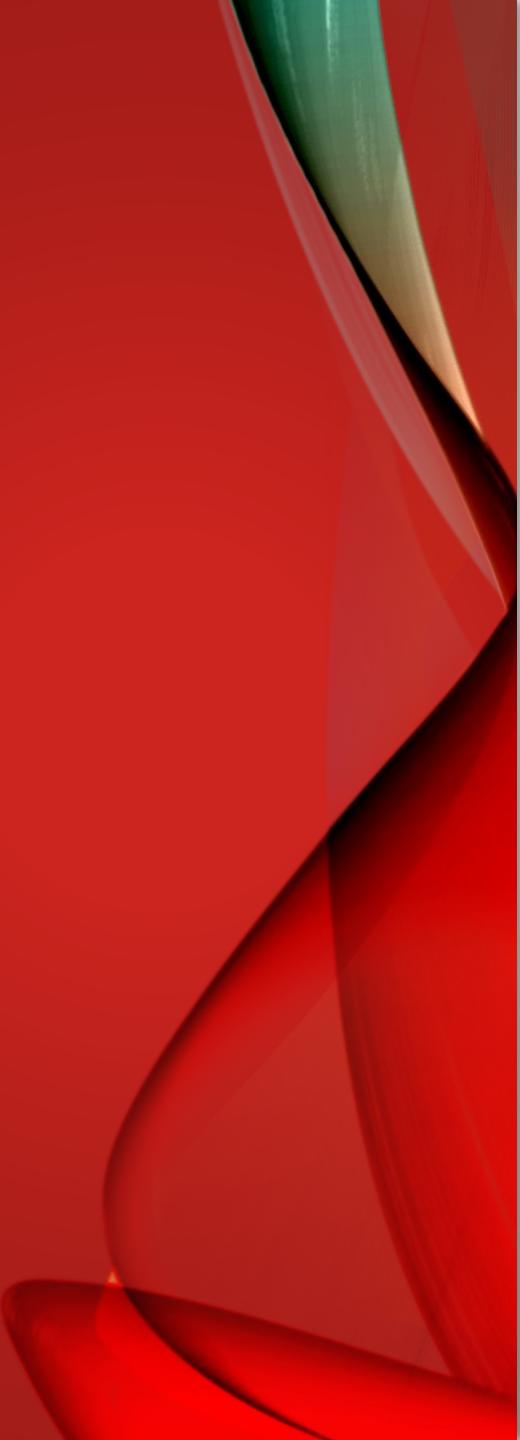
WHAT IS 🦜🔗 LANGCHAIN?

PRINCIPLES OF LANGCHAIN

Data-awareness

Agentic interaction

Modularity and
customization



DATA AWARENESS

- LangChain excels in connecting with diverse data sources.
- It goes beyond data retrieval, focusing on comprehension and context.
- Seamlessly integrates with databases, cloud storage, and online repositories.
- Ensures data accessibility and preparedness for processing.

AGENTIC



LangChain goes beyond data retrieval to enable "agentic" interactions.



This dynamic interaction with the environment creates responsiveness.



It's not just about answering queries; it's about context understanding and problem-solving.

MODULARITY AND CUSTOMIZATION

LangChain adopts a component-based architecture for modularity.

Developers can customize applications without unnecessary complexity.

Pick and choose components for chatbots, document summarizers, or code analyzers.

Tailor your solutions to fit your specific vision and requirements.



WHY LANGCHAIN MATTERS

- LLMs like ChatGPT and GPT-4 offer immense potential.
- They revolutionize human-machine interaction and understanding.
- Integrating these models into applications can be challenging.
- LangChain addresses this challenge and facilitates practical applications.

STREAMLINED LLM INTEGRATION



LangChain bridges the gap between LLMs and applications.



It offers essential tools and interfaces for seamless integration.



Enables efficiently incorporating models like GPT-4 or BLOOM into diverse applications.

TANGIBLE BENEFITS

LangChain translates LLM advantages into practical outcomes.

Legal firms can summarize documents rapidly, and customer support centers can offer real-time, context-aware chatbot responses.

The potential applications are diverse and extensive.

ADAPTING TO THE AI LANDSCAPE



In the evolving AI landscape,
staying updated is crucial.



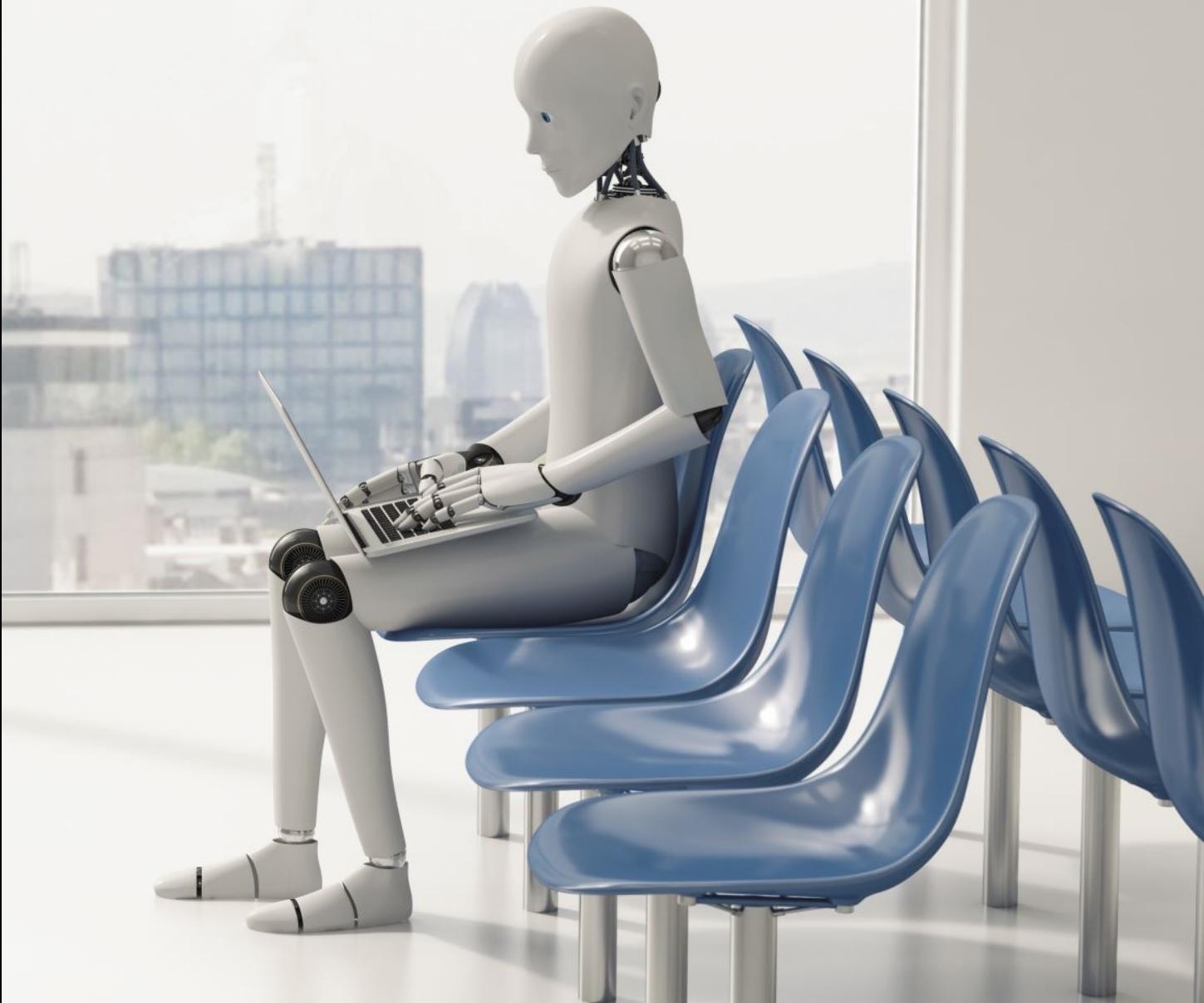
Mastery of LangChain provides an
advantage in adapting and
thriving.



Empowers learners to navigate
the dynamic AI environment
effectively.

LANGCHAIN AND HUMAN-MACHINE INTERACTION

- Simplifies LLM integration for real-world applications.
- Enabling natural, intuitive, and productive human-machine interactions.
- LangChain and prompt engineering are central to this transformation.
- A call to action for developers, businesses, and curious minds to unlock language model potential.



LANGCHAIN OVERVIEW

CORE MODULES IN LANGCHAIN

Model
I/O

Retrieval

Chains

Memory

Agents

Callbacks

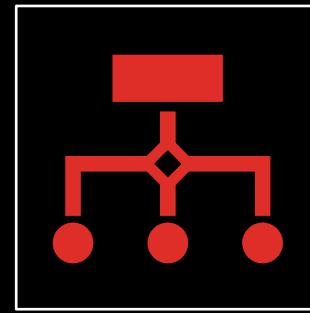
MODEL I/O



Prompts for managing prompt messages.



Interfaces for calling different types of models.



Output parsers for extracting structured information from free-form model outputs.



RETRIEVAL

- Many LLM applications require domain-specific knowledge.
- The retrieval module sources prepares, and retrieves external data.
- Features include basic retrieval and advanced techniques like RAG.
- RAG leverages LLMs and domain-specific data for precise answers.
- Components: document loaders, transformers, embeddings, vector stores, retrievers.

CHAINS

- LLMs are sufficient for simple apps, but complex ones benefit from chaining.
- LangChain's Chains module ensures modularity and maintainability.
- Chaining simplifies complex app development with multiple components.
- LLMChain combines message, template, and LLM for basic building blocks.
- The Chain interface allows users to create coherent applications by combining components.

AGENTS



Agents module creates AI agents for natural conversations and actions.



Manages conversation flow, tool integration, and memory.



Combines language model capabilities with executable tools.

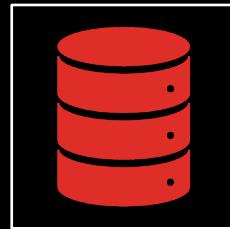


Core components: base agent, tool abstractions, agent executor, conversational agent.

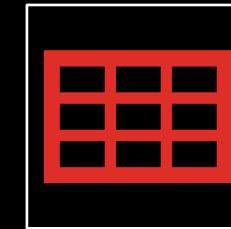
MEMORY



Memory module maintains context for LangChain agents.



Interfaces for storing agent states across multiple inferences.



Simplifies complex storage options through a unified API.

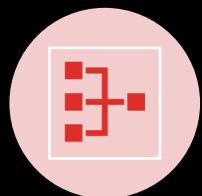


Essential for LangChain agents to maintain reasoning capabilities.

CALLBACKS

- Callbacks module facilitates event integration in LangChain.
- Enables functionalities like logging, monitoring, streaming, and custom integrations.
- We won't get into callbacks for this course

OVERVIEW OF LANGCHAIN MODULES



Model I/O Module streamlines model integration.



Retrieval Module ensures data enrichment,



Chains Module simplifies complex app development,



Agents Module creates AI agents for natural interactions,



Memory Module helps maintain context, and



Callbacks Module: Facilitates event integration for logging



MODEL I/O

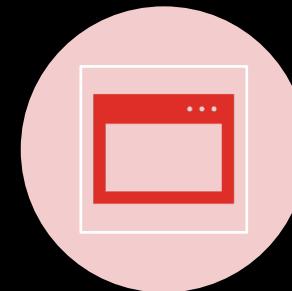
KEY COMPONENTS OF MODEL I/O

- Language models: Interfaces for LLMs and Chat models.
- Prompts: Templates for parametrizing and reusing prompts.
- Output Parsers: Enable text output extraction and structuring, useful for tasks like QA.

TYPICAL LANGCHAIN WORKFLOW



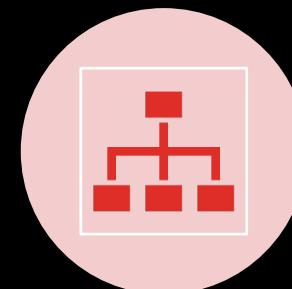
Select LLM or Chat model based on your specific use case.



Create a customized prompt as the input.



Transmit the input to the chosen LLM or Chat model.



Optionally, employ output parsers to process outputs as necessary.



WORKING LLMS

- The LLM class serves as a universal interface to LLM providers.
- It abstracts provider-specific APIs, offering predict and generate methods.
- Use "predict" for text completion by providing a prompt template; the output is plain text.
- "generate" accepts prompt lists and yields detailed LLMResult objects with completions and metadata.



CRAFTING EFFECTIVE PROMPTS

Instructions: Specify the model's intended response or behavior.

Context: Offer supplementary information, often with illustrative examples.

User Input: The user's actual query or input.

Output Indicator: Signifies the commencement of the model's response.

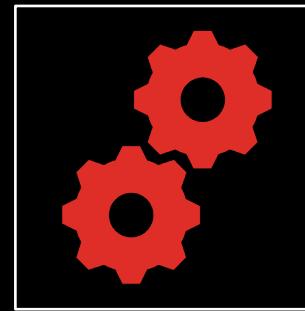
PROMPT TEMPLATES

- `PromptTemplate` employs placeholders (e.g., `{adjective}`, `{content}`) within template strings, facilitating final prompt string generation.
- Input variable validation against the template.
- Versatile input options, including dictionaries and dataclasses.
- Compatibility with various templating engines like Python's `str.format` or `Jinja2`.
- User-friendly extension and custom template creation.

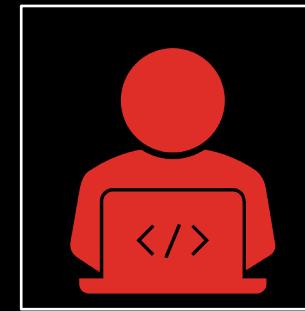
OUTPUT PARSERS



Raw text from language models
may not always suit downstream
requirements.



Output parsers structure model
text.



Convert to JSON, Python
dataclasses, and other useful
structures.

USE CASES OF OUTPUT PARSERS



Structuring unstructured text into structured data (e.g., JSON or Python objects).



Guiding language models in response formatting by injecting instructions into prompts via a "get_format_instructions()" method.

MODEL I/O SUMMARY

- Model I/O is the cornerstone of LangChain's functionality.
- Comprises three vital components:
 - Language models (LLMs) for interaction.
 - Prompt templates for customization.
 - Output parsers for text manipulation.
- Enables streamlined workflows for diverse use cases.

A small, white and tan dog with a blue and red textured collar is running towards the camera on a dirt path. The path is covered with fallen autumn leaves. In the background, there are tall trees with orange and yellow foliage. The overall atmosphere is bright and natural.

Retrieval



Retrieval is useful because it allows you to incorporate external data into your language model.

BENEFITS OF RETRIEVAL IN LANGCHAIN

- Incorporates user-specific data for personalized responses
- Supplements model knowledge with current information
- Aids in answering questions from large document sets
- Essential when accessing external data based on user queries
- Enhances language model responses and accuracy of answers.

USING DOCUMENT LOADERS



Multiple types available: PDF, Webpage, CSV, Directory, and more



Import the desired loader from `langchain.document_loaders`

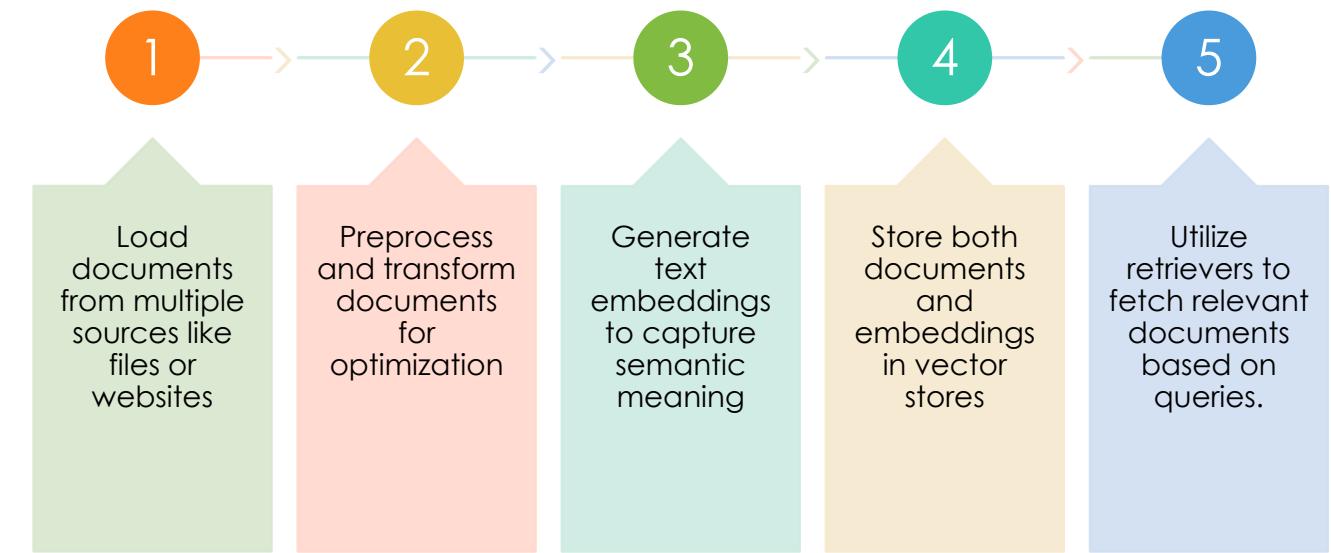


Instantiate the loader class with necessary arguments



The `load()` method converts files into LangChain's Document format.

IMPLEMENTING RETRIEVAL IN LANGCHAIN



SUMMARY OF RETRIEVAL



Enables personalized and context-aware model responses



Facilitates access to up-to-date information and facts



Streamlines question answering from vast document sets



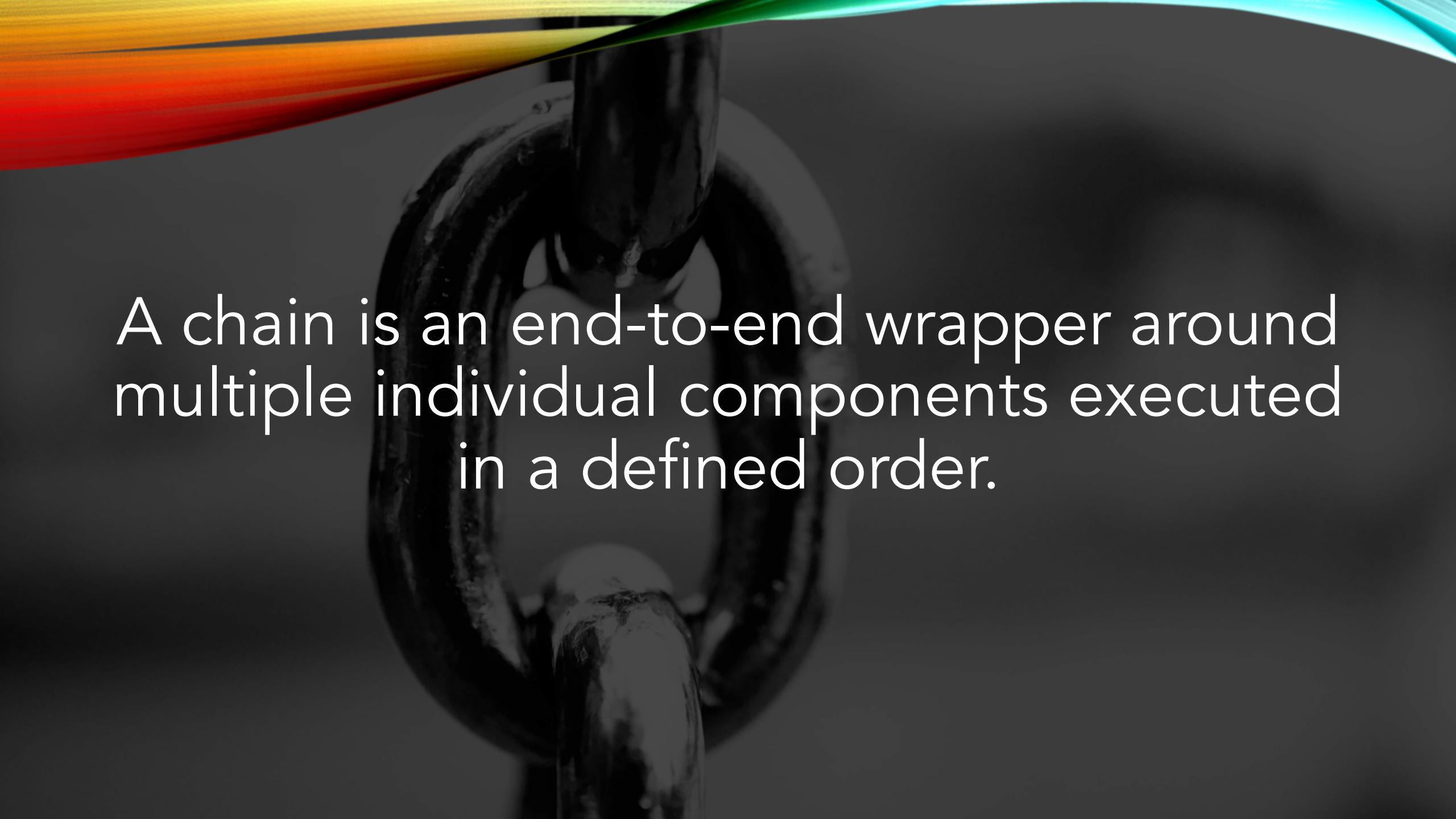
Incorporates a variety of document loaders for diverse sources



Enhances LLM applications for improved performance and understanding.



CHAINS



A chain is an end-to-end wrapper around multiple individual components executed in a defined order.

BENEFITS OF USING CHAINS



Breaks complex tasks into manageable, sequential steps.



Leverages strengths of different systems.



Provides context and state through sequential calls.



Enables additional processing and validation between calls.



Simplifies debugging and monitoring of call sequences.

CORE CHAIN TYPES IN LANGCHAIN

-  LLMChain: Chains language models for complex prompt handling.
-  RouterChain: Provides conditional routing for branching logic.
-  SimpleSequentialChain: Enables linear chaining of multiple chains.
-  TransformChain: Facilitates data transformations between chains.
-  Foundation for building more complex chain structures.

UNDERSTANDING LLMCHAIN IN LANGCHAIN

- LLMChains enhance language model functionality.
- Integrated across LangChain, including in other chains and agents.
- Comprises a PromptTemplate, a language model, and an optional output parser.
- Facilitates formatting user input for LLM responses.
- Enables the creation of intricate chains through combination.



ROUTER CHAINS IN CHATBOTS

- Route inputs to appropriate destination chains based on text.
- Examination and decision-making done by the router chains.
- Destination chains manage the specific task execution.
- Central to creating versatile multi-purpose chatbots and assistants.
- Enables chatbots to handle varied user requests effectively.

SEQUENTIAL CHAINS

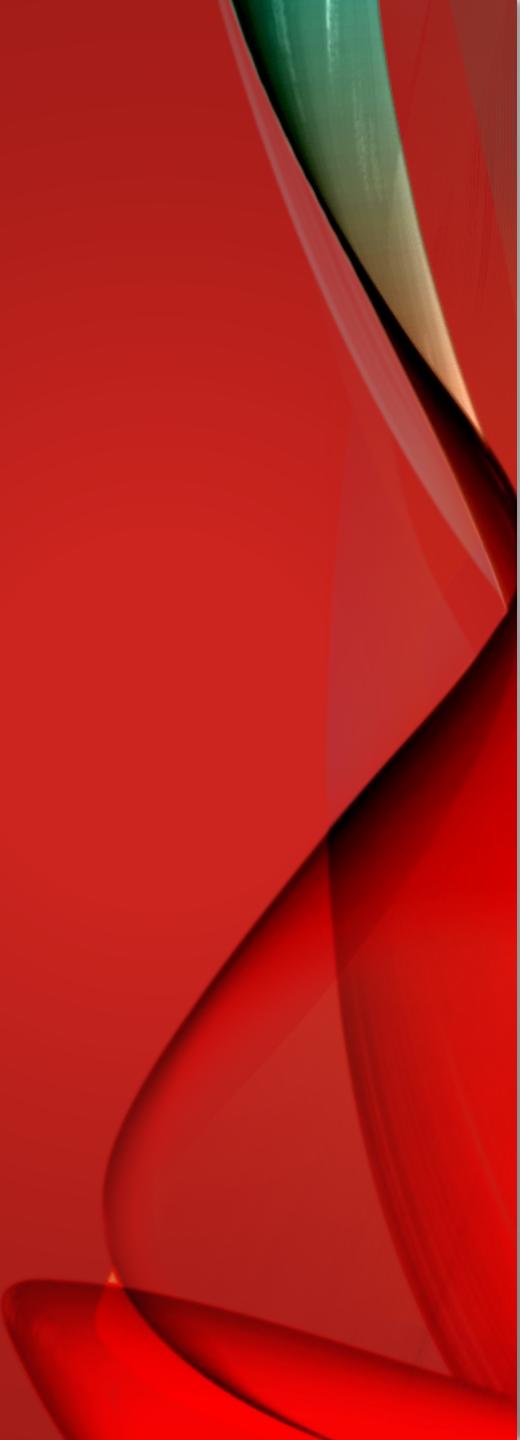
Enable series of calls, using output from one as input to another.

Compose multiple chains into specific execution pipelines.

SimpleSequentialChain: Singular input/output, linear flow.

SequentialChain: Allows for multiple inputs/outputs.

Useful for multi-step processing in language models.



TRANSFORM CHAINS

- Apply text transformations such as splitting, filtering, and more.
- Define custom data transformation logic within the pipeline.
- Facilitate sequential application of multiple transformations.
- Useful for data preprocessing before advancing to next steps.
- Harnesses the `TransformChain` class for structured transformation.

SUMMARY OF CHAINS

- Serve as foundational building blocks for complex operations.
- LLMChain focuses on enhancing language model interactions.
- RouterChain offers dynamic routing based on input, ideal for versatile chatbots.
- Sequential Chains streamline multi-step processes, from simple to complex workflows.
- Transform Chains specialize in sequential text data transformations.



AGENTS

AGENTS VS CHAINS

- Agents use LLMs to decide on a sequence of actions.
- Chains have hardcoded sequences defined in code.
- Agents employ language models as reasoning engines.
- Agents interact with other tools or environments through language models.
- Agents decide, act, observe results, and iterate for outcomes.



COMPARING AGENTS AND CHAINS

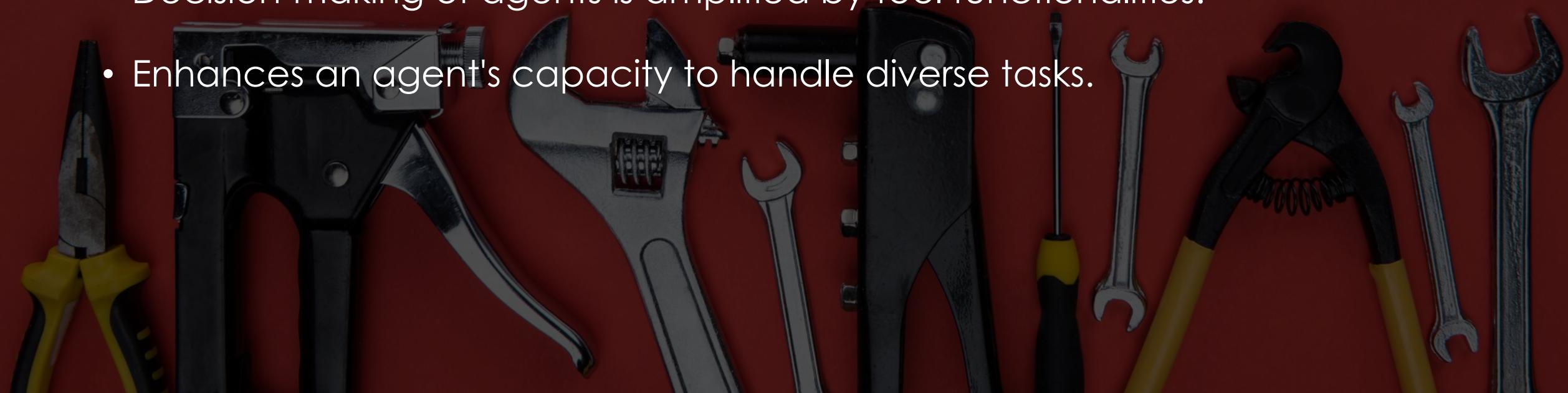
- Chains emphasize a sequence of calls and data flow.
- Agents center on decision-making and environment interaction.
- Agents serve diverse applications from chatbots to data querying.
- An LLM acts as the reasoning engine in agents.
- Agents uniquely integrate with tools and memory components.

TOOLS AND TOOLKITS

- Tools are components dedicated to specific tasks.
- They can fetch data or process it for various needs.
- Toolkits group multiple tools together.
- Designed to ensure tools within work harmoniously.
- Offer a broader set of combined functionalities.

WHY DO AGENTS NEED TOOLS?

- Tools empower agents to execute and implement solutions.
- Agents and tools together create versatile systems.
- Decision-making of agents is amplified by tool functionalities.
- Enhances an agent's capacity to handle diverse tasks.



WHAT AGENTS CAN DO WITH TOOLS

- Tools grant access to data from external sources.
- They aid in data cleaning, transformation, and analysis.
- Facilitate integration with systems like chatbots or CRMs.
- Allow customization for unique agent needs and tasks.
- Enable integration with specialized or proprietary systems.

THE IMPORTANCE OF MEMORY FOR AGENTS

- Supports storage and retrieval of information.
- Maintains context for effective decision-making.
- Remembers past interactions for continuity.
- Ensures personalized responses to users.
- Crucial for coherent and consistent communication.

HOW MEMORY ENHANCES AGENTS PERFORMANCE

- Facilitates understanding of context from past inputs.
- Accumulates long-term knowledge for precise responses.
- Personalizes interactions based on user history.
- Ensures consistent and natural user dialogue.
- Amplifies user experience through engaging conversations.

SUMMARY: AGENTS IN LANGCHAIN

Agents harness decision-making for effective task execution.

Tools and memory are crucial for an agent's enhanced performance.

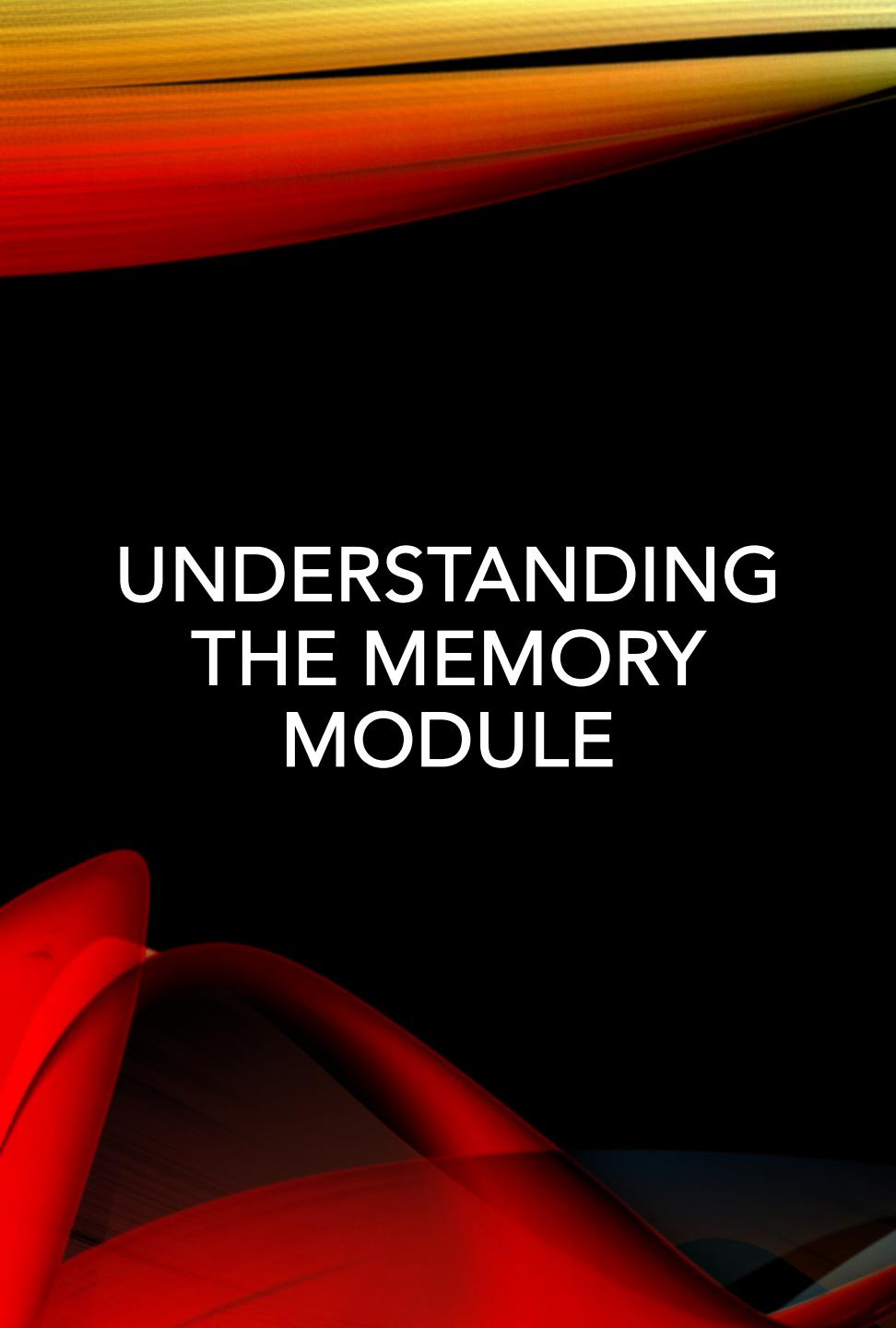
Agents ensure contextual understanding and personalized user interactions.

Integration capabilities with tools offer versatility in applications.

Memory fosters continuity, long-term knowledge, and tailored experiences.

A scene from the Pixar movie Finding Nemo. Dory, the blue tang fish with a yellow fin, is swimming towards the camera with a wide, joyful smile. She has large, expressive eyes. The background is a vibrant underwater environment filled with various coral formations, sea fans, and other small tropical fish. The lighting is bright and colorful, typical of the Pixar animation style.

MEMORY



UNDERSTANDING THE MEMORY MODULE

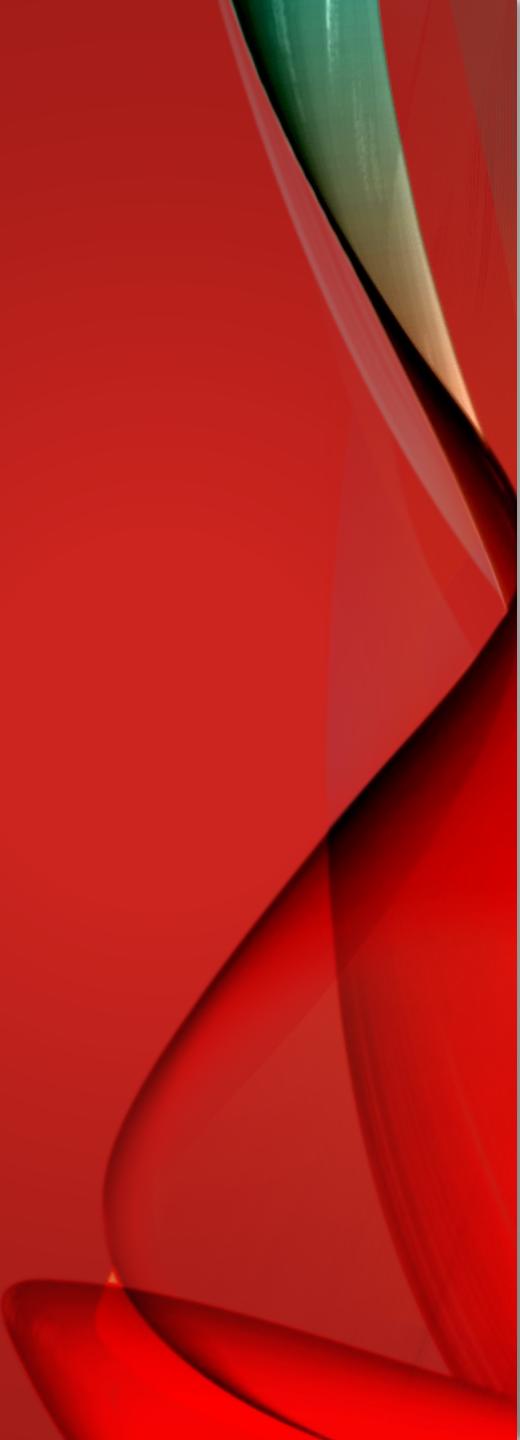
Memory module persists state between chain or agent calls.

Aids in recalling previous interactions for better decisions.

Offers a standard interface for state persistence.

Enables the language model to maintain memory and context.

Essential for personal assistants, autonomous agents, and simulations.



FUNCTIONS OF THE MEMORY MODULE

- Provides memory and context to the language model.
- Empowers LLM to make informed decisions.
- Remembers user inputs, system responses, and related data.
- Stores critical interaction details.
- Facilitates data access during future interactions.

WHEN AND HOW TO USE THE MEMORY MODULE

- Use for applications needing context and persistence.
- Ideal for personal assistants to recall user preferences and past queries.
- Memory system's primary functions: reading and writing.
- Every chain requires specific inputs, some user-derived and some from memory.
- Chains read memory before core logic execution and write post-processing.



KEY DECISIONS IN MEMORY SYSTEMS

- Two core decisions: how to store state and how to query it.
- Memory captures records of all chat interactions.
- LangChain provides multiple storage options, from temporary lists to databases.
- Designing algorithms for chat interpretation is complex.
- Systems vary: some display recent chats, others summarize the last 'K' messages.

IMPLEMENTING MEMORY IN LANGCHAIN

- Refined systems identify and focus on chat entities in current sessions.
- Different apps need distinct memory querying methods.
- Implementation Steps:
 1. Setup prompt and memory.
 2. Initialize LLMChain.
 3. Call LLMChain.

SUMMARY: MEMORY IN LANGCHAIN

- Central to maintaining context across sessions.
- Ranges from basic to advanced storage and querying methods.
- Versatile: From personal assistants to tailored applications.
- Key aspects: storage method and state querying.
- Simple initiation with potential for customization.