HNS analytical software



LLM-MVC

A simple design pattern for structuring standard LLM tasks with generic prompts



Dr. Robert Bauer Senior Data Scientist HMS Analytical Software GmbH

data2day 19.09.2024

Introduction

Problem Statement

- Unstructured text from LLM input/output is difficult to handle and maintain, especially for large systems
- In some cases, working with structured data (models) is the better approach and many productive LLM systems use this
- However: there are different ways how such model centric systems can be designed and controlled: Do we use chains? Agents? DAGs? And do we always need the "big guns"?

Topic of the Talk

We introduce a **simple design pattern for model centric LLM systems** build upon the concepts of the well known Model-View-Controller architecture

HMS Analytical Software GmbH

.

Agenda

About HMS Model Decomposition 01 05 Model-Driven Prompting Code Examples 02 06 Motivation for Design Lessons Learned and 03 07 Summary Patterns MVC pattern for LLM 04 workflows

HMS Analytical Software GmbH





HMS Analytical Software GmbH









Model-driven Prompting



Direct Prompting

Task: Summarize tech blog posts and perform relevance scoring for different stakeholders (CEO, CSO, ...)



Model-driven Prompting

Task: Summarize tech blog posts and perform relevance scoring for different stakeholders (CEO, CSO, ...)



Model-driven Prompting

```
from typing import Union
# Pydantic
class Joke(BaseModel):
    """Joke to tell user."""
    setup: str = Field(description="The setup of the joke")
    punchline: str = Field(description="The punchline to the joke")
    rating: Optional[int] = Field(
        default=None, description="How funny the joke is, from 1 to 10"
class ConversationalResponse(BaseModel):
    """Respond in a conversational manner. Be kind and helpful."""
    response: str = Field(description="A conversational response to the user's query")
class Response(BaseModel):
    output: Union[Joke, ConversationalResponse]
structured llm = llm.with structured output(Response)
structured llm.invoke("Tell me a joke about cats")
```

https://python.langchain.com/v0.2/docs/how_to/structured_output/





Design Patterns

Design Patterns **describe common problems** in one or two words and solve them in a **structured and proofen** way

And they help us avoiding situations like these 🙂



Original Source: http://www.bonkersworld.net/building-software/ Accessed via: https://blog.codeinside.eu/2011/11/27/pragmatische-softwareentwicklung/

Examples



Examples

8

Factory Metho



- Various patterns for RAG / Few-Shotting / Validation and other important LLM topics
- Panel of Experts







Example Workflow with Original Pattern



How does MVC help?



LLM-MVC



LLM-MVC















Decomposition and Indirect Prompting



Model-driven Prompting

Task: Summarize tech blog posts and perform relevance scoring for different stakeholders (CEO, CSO, ...)

Pecap



Indirect Prompting

Task: Summarize tech blog posts and perform relevance scoring for different stakeholders (CEO, CSO, ...)



Model Decomposition

Core Idea

- We do NOT resolve the complete model in one batch (this is simple)
- Instead we
 individually
 resolve the parts
 step by step

Why not in one batch?

- Schemas can get large quickly for complex (real world) use cases
- From our experience, even the latest LLM generations struggle with large schemata with lots of annotations
- Limited control over individual parts of the model

Dependencies between Models



Dependencies define the context that is used together with the schema of a composite element

Wrap-Up

Dependency Management

- There are various different ways how data modeling and dependency management can be done
- Far beyond this talk
- Keep in mind: you can also use the pattern without decomposition to keep things simple

Thinking in "Model Types"

- To get started with decomposition, it can be very helpful to explicitly distinguish models based on their conceptual role
- Next: The basic modeling approach we use for our projects





- > Describe static a-priori information
- Only used as additional LLM context
- Never resolved by the LLM (=LLM never creates an object with schema of this model)
- > May be filled from database or RAG system



- Models that are actively resolved by an LLM with modeldriven prompting
- Means the schema of the model with annotations is set as "output format"
- Are always resolved with their complete schema which is different from how composed models are treated



- > Models that are not directly resolved due to size/complexity
- Instead, sub-models of the composite are resolved recursively after some form of dependency resolution has taken place
- > Can consist of static, input and other composed models



- Models that are derived for view construction
- Never resolved by the LLM; if something has to be resolved, its either in the I or the C category!
- > Usually a reduced version of one of the other model types

LLM-MVC Overview with Model Types



























Some important observations

- Subsequent calls can use smaller context windows (i.e., not the complete article)
- Smaller context = more control / easier for the LLM to follow instructions/annotations
- Steps can be validated individually!
- Models can have dependencies: Scores_XYZ depend on summary, for example
- If parts are changed, only dependent parts need to be updated!
- Process can be fully automated for generic cases (no real need for chains/agents here)

Static (pre-	Unresolved	Resolved from	Resolved this	Schema =	
resovled)		previous Step	step	target output	









Lessons Learned

• We frequently use model-centric approaches Patterns and Robust and maintainable Templates do help • Easy to explain + Easy to get started with templates • Frameworks such as langchain can be overkill Simplicity is key • LLM-MVC pattern can be build with requests and pydantic alone and is still an powerful abstraction • If small parts of an application are changed, we do not Partial re-evaluation have to run everything from scratch is a big issue for us • Decomposition = re-run only changed + dependent

Partial re-evaluation

Everything New			New	New Metadata			New Scoring Algo			New Tag Algo		
Costs with Composition			Costs with	Costs with Composition			Costs with Composition			Costs with Composition		
Scoring1	435	98	Scoring1	435	98		Scoring1	435	98	Scoring1	0	0
Scoring2	605	60	Scoring2	605	60		Scoring2	605	60	Scoring2	0	0
Basics	2002	597	Basics	2002	597		Basics	0	0	Basics	0	0
Tags	413	33	Tags	0	0		Tags	0	0	Tags	413	33
Summe	3455	788	Summe	3042	755		Summe	1040	158	Summe	413	33
	+25%	+35%		+10%	+30%			-63%	-73%		-85%	-95%
Without Composition Without Composition		tion		Without Composition			Without Composition					
Summe	2752	585	Summe	2752	585		Summe	2752	585	Summe	2752	585

Further Notes and Conclusion

Streaming?	 So far we primarily applied the pattern to batch use cases / use cases without streaming requirements Streaming + Decomposition do not work well together out-of-the-box (probably still possible) 			
Use in existing projects	 Depends; basic ideas can be integrated easily More effort if you truly want to benefit from decomposition 			
Frameworks	 Approach is framework independent Can be used with or without major frameworks We do both in our projects 			



Vielen Dank für die Aufmerksamkeit!



Dr. Robert Bauer Senior Data Scientist HMS Analytical Software GmbH

data2day 19.09.2024

