



PHP AND LLMs

PRACTICAL INTEGRATIONS USING LARAVEL

**The LLM hype may be fading, but these
tools empower us to solve real
customer needs more efficiently.**

by Alfred Natile

.001

PHP AND LLMs

Practical Integrations using Laravel

Alfred Natile

This book is available at http://leanpub.com/php_and_llms

This version was published on 2024-08-19



The author generated this text in part with GPT-3, OpenAI's large-scale language-generation model. Upon generating draft language, the author reviewed, edited, and revised the language to their own liking and takes ultimate responsibility for the content of this publication.

© 2024 Alfred Natile

Tweet This Book!

Please help Alfred Natile by spreading the word about this book on [Twitter](#)!

The suggested tweet for this book is:

[Excited to learn practical applications of LLMs in my Laravel applications.](#)

The suggested hashtag for this book is [#laravelandllms](#).

Find out what other people are saying about the book by clicking on this link to search for this hashtag on Twitter:

[#laravelandllms](#)

Contents

[COMPLETED] Introduction - Embracing the Reality of LLMs in Development	1
Local LLMs Streamlining Your Development Workflow	3
Building the LLM Agnostic Driver	4
Prompting Overview before we Dig In	5
[COMPLETED] Parsing Opt-Out Text	6
Tools Overview Enhancing LLM Capabilities	10
[IN PROGRESS] - Web Scraping and Building an Events Database	11
Overview	11
Orchestration	16
Deep Dive	20
Wrapping It Up	35
Making Tools LLM-Agnostic	36
Chaining Tools	37
Web Searches with LLMs	38
Marketing Assistant	39
Hand Written Notes to Task List	40
RAG and Vector	41
SiteMap - Feed to RAG	42
Agents to Evaluate the Results	43
Standards Checking Tool	44
Final Thoughts	45

[COMPLETED] Introduction - Embracing the Reality of LLMs in Development

Yes, the hype around LLMs might be fading, but the reality is that these tools offer transformative possibilities for developers. My goal is to provide you with the knowledge and confidence to integrate LLMs into your projects, work effectively with them locally, and add a powerful new tool to your development arsenal. And we won't just skim the surface with "Hello World" examples. We're going to dig deep into the real-world applications that can genuinely enhance your day-to-day work as a developer.

PHP has a place in this seemingly Python + Ai era.

By the end of this book, I'll equip you with the skills to:

- **Work Locally with Ollama:** Learn how to save costs and iterate on your ideas effectively by working locally with Ollama. This approach will allow you to refine your concepts without the overhead of cloud services, giving you the flexibility to develop and test in your own environment.
- **Test and Mock with Confidence:** We'll explore how to set up robust testing and mocking frameworks, so you can deploy your code with the assurance that it will perform reliably in production. This chapter will help you build a foundation of confidence, knowing that your integrations are solid and ready for real-world use.
- **Integrate Multiple LLMs Seamlessly:** Discover how to write code that's flexible enough to work with any LLM service—whether it's Claude, Ollama, OpenAI, Groq, or another. We'll cover strategies for creating a unified interface, allowing you to switch between services or use multiple providers simultaneously, all within the same codebase.
- **Apply Real-World Solutions:** I'll share numerous examples of real-world solutions I've implemented (or plan to implement) that demonstrate the practical power of LLMs. These examples will show you how to save time, write efficient code, and tackle complex situations that might have seemed daunting before.
- **Master Tools for LLMs:** One of the most critical skills you'll learn is how to empower your prompts with tools. We'll dive into how tools can drive a wide range of tasks and automations, making your LLMs not just smart, but practically useful in day-to-day development.
- **Make Tools LLM-Agnostic:** Finally, we'll cover how to create tools that are independent of any specific LLM. This means you can write your prompts and logic once and use them across multiple LLMs, ensuring that your work remains flexible and future-proof.

By the end of this book, you'll not only have the skills and confidence to move forward with LLM technology but also a clear understanding of how to apply it practically in your work. This journey is about more than just learning a new tool; it's about expanding your capabilities as a developer and seeing the real potential that LLMs bring to solving complex problems efficiently and effectively.

Local LLMs Streamlining Your Development Workflow

Summary: We explore the benefits of working with local LLMs like Ollama, which allow you to develop and test your ideas without incurring cloud costs. This chapter includes practical tips for setting up local environments and integrating them into your workflow. This will include tips and tricks on how to test locally and in CI.

Building the LLM Agnostic Driver

Summary: This chapter covers how to integrate multiple LLMs into your projects, allowing you to switch between services like OpenAI, Claude, and Groq seamlessly. We discuss the advantages of this approach and provide examples of how to implement it effectively.

Prompting Overview before we Dig In

Summary: This chapter delves into the importance of crafting the right prompts and providing the correct context to LLMs. We explore different prompt iterations and demonstrate how to refine them to get the best results, focusing on a real-world example of parsing opt-out emails.

[COMPLETED] Parsing Opt-Out Text

So far, we've covered how to connect to LLMs, pass a request, and get a response.

In this chapter, I'll walk you through a problem I had to solve for a customer that perfectly showcases the power of these tools. The customer receives emails from various opt-out providers requesting to remove individuals from their mailing list. The catch? The data comes in all sorts of formats—some not even plain text, but just HTML!

In the past, I would have written a text parser specific to each domain to extract the Full Name, Phone Number, Address, and Reply-To Address (which, by the way, isn't always in the email header).

But one day, when a new domain started sending us requests, I had a lightbulb moment: why not use an LLM to handle the parsing?

First and foremost, you need to spend some time fine-tuning the prompt that you will pass to the llm with the email content. Ultimately, everything boils down to getting the right Prompt (text) and the right Context (text) to the LLM. This principle holds true no matter what problem you're trying to solve.

It's all about getting the right Prompt (text) and the right Context (text) to the LLM.

With this in mind, using the Driver system we built in a previous chapter, we can plug this into Tinkerwell and start testing our ideas.

```
1 <?php
2 $results = LlmDriverFacade::driver("ollama")
3     ->setFormat("json")
4     ->completion($prompt);
5
6 $results->content;
```

Note we are doing `setFormat` since many of these LLMs accept that. We then, depending on the LLM alter the request as needed.

The response we get is:

```
1 { "full_name": "John Doe", \n
2   "home_address": "123 Main Street, Anytown, USA", \n
3   "other_name": null, \n
4   "reply_to": "johndoe@example.com", \n
5   "valid_user_email": "johndoe@example.com", \n
6   "phone": "555-555-5555"\n
7 }
```

You will notice that in the sample email below, parsing this text is possible BUT to write this for every opt-out vendor, some which had harder to parse emails, it was more work than to do this once.

\$email

```
1 Subject: Request for Deletion of Personal Information
2
3 My name is John Doe, and I request that Your Company deletes all of the information \
4 it has collected about me, whether directly from me, through a third party, or throu
5 gh a service provider.
6
7 The following information is provided as required for validation of my request:
8
9 First Name: John
10 Last Name: Doe
11 Email: johndoe@example.com
12 Address: 123 Main Street, Anytown, USA
13 Phone Number: 555-555-5555
14
15 If you need any more information, please let me know as soon as possible. If you can\
16 not comply with my request—either in whole or in part—please state the reason why yo
17 u cannot comply. If part of my information is subject to an exception, please delete
18 all information that is not subject to an exception. If my request is incomplete, p
19 lease provide me with specific instructions on how to complete my request.
20
21 I have used a consumer identity protection service to send this email, but please re\
22 spond directly to my personal email address johndoe@example.com.
23
24 I am sending this request exercising my rights to my Personal Data under applicable \
25 privacy laws, including, but not limited to, the General Data Protection Regulation
26 (“GDPR”) and the California Consumer Privacy Act (“CCPA”)
```

\$prompt

```
1 $prompt = <<<PROMPT
2 <role>
3 Help parse opt out emails so we can remove the person from our system or ignore the \
4 email if it is
5 not related to opt out.
6
7 <task>
8 Review the email below and then if it is an opt out request pull the data out as def\
9 ined by the format section.
10 valid_user_email might be tricky since sometimes the email is generatic and going ba\
11 ck to the company sending these opt outs.
12 But if you are confident the email looks like a valid personal or business email add\
13 that to valid_user_email.
14
15 <format>
16 JSON FORMAT with key values below, using null if no value for the key. This is valid\
17 JSON not markdown.
18 If there is no data in this about a person just return the JSON below with all field\
19 s null. No surrounding text.
20
21 {
22   "full_name": null,
23   "home_address": null,
24   "other_name": null,
25   "reply_to": null,
26   "valid_user_email": null,
27   "phone": null
28 }
29
30 Only return valid JSON
31
32 {$email}
33 PROMPT;
```

NOTE this is like mu 4-5 th prompt attempt, “valid JSON” phrase came from the the www.deeplearning.ai trainings.

Then we pass that to the LlmDriver:

```
1  <?php
2  $results = LlmDriverFacade::driver("claude")->completion($prompt);
3  $results->content;
```

Note we will swap out drivers to show how it works.

And we still get:

```
1  {\n
2    "full_name": "John Doe", \n
3    "home_address": "123 Main Street, Anytown, USA", \n
4    "other_name": null, \n
5    "reply_to": null, \n
6    "valid_user_email": "johndoe@example.com", \n
7    "phone": "555-555-5555" \n
8  }
```

As you can see, we use the driver, which can be any LLM—whether it’s “ollama”, “claude,” “openai,” or “groq.” We then pass it the prompt along with the HTML/Text of the email.

We also enforce a JSON format for the output, which gives us some confidence in the results. From here, we can decide how to proceed:

- Ignore the output if all the fields are null.
- Process the output if the required fields are present.

Also, note that the reply-to address is embedded in the email. Many of these opt-outs require us to reply to the email that sent the request, but as you can see here, we had to adapt to this one.

This approach saved me from manually parsing the text out of these emails and significantly reduced the amount of code, as I previously had parsers for two other vendors.

Later I will talk about how to evaluate the accuracy of the results, as well as ideas on testing.

Tools Overview Enhancing LLM Capabilities

Summary: This chapter introduces the concept of tools that enhance the functionality of LLMs. We explore how to build tools that can work with any LLM and how to integrate them into your projects. A real-life example of an event scraping tool is used to demonstrate these concepts.

[IN PROGRESS] - Web Scraping and Building an Events Database

Overview

Here, we explore how use data scraped from a website and use LLMs to process the content. The chapter includes examples of tools that can extract relevant information from web pages and convert it into Events in our database.

We don't need to delve into scraping here since it's pretty common and there are many ways to do it. What I want to focus on is building a Tool that can handle a complex array of objects (Events) for the LLM to pass into it.

Once we have that tool in place, the following prompt can do a lot:

```
1 <role>
2 You are an assistant helping to parse data from HTML content and find the events on \
3 the page.
4 <task>
5 Using the context below, find all the events and pass them into the create event too\
6 l.
7 After that, save it to a document.
8 <context>
9 [INSERT OUR HTML HERE]
```

As you can see above, we are telling the prompt to use our tool. Keep in mind users would not be so specific; they might write "Find and create events from this page," but we can test that as well.

So we need a simple Event model in the database for it to save the events. Then we need to make the tool and have it define the required arguments. When it is handed the function call from the LLM, it will create events.

As we saw in a previous chapter, Tools are ways to define a function and parameters that the LLM can use to complete a task. The task can be simple like `get_the_weather`, but we are going to do more complex tasks. To start with, this tool is going to have a complex parameter definition (this will be JSON for Claude; we'll show the abstracted version later in this chapter).

The `get_the_weather` is a Class we made to process the arguments. Below is the example structure for the `create_event_tool`

```

1  {
2      "name": "create_event_tool",
3      "description": "If the user needs to create one or more events this tool helps",
4      "input_schema": {
5          "type": "object",
6          "properties": {
7              "events": {
8                  "description": "Array of event objects",
9                  "type": "array",
10                 "items": {
11                     "type": "object",
12                     "properties": {
13                         "start_time": {
14                             "type": "string",
15                             "description": "Start time of the event"
16                         },
17                         "end_time": {
18                             "type": "string",
19                             "description": "End time of the event"
20                         },
21                         "title": {
22                             "type": "string",
23                             "description": "Title of the event"
24                         },
25                         "location": {
26                             "type": "string",
27                             "description": "Location of the event"
28                         },
29                         "summary": {
30                             "type": "string",
31                             "description": "Summary or description of the event"
32                         }
33                     }
34                 }
35             }
36         },
37         "required": [
38             "events"
39         ]
40     }
41 }

```

This took me a while to get right, but thanks to Claude's help, I was able to figure out how to do an

array of objects.

- name - this will be seen later registered in our application.
- description - here we do a good job of letting the LLM know how to use the tool.
- input_schema - this is the tricky part

You can see here the very simple tool in the docs:

```
1  {
2      "name": "get_weather",
3      "description": "Get the current weather in a given location",
4      "input_schema": {
5          "type": "object",
6          "properties": {
7              "location": {
8                  "type": "string",
9                  "description": "The city and state, e.g. San Francisco, CA"
10             }
11         },
12         "required": ["location"]
13     }
14 }
```

By now, we are going to take it up a notch. Our first property “events” has the type of “array”.

And in that object key, we start to define our properties. As you can see in the `curl` command below, we are just going to pass this to Claude.

```
1  curl --location 'https://api.anthropic.com/v1/messages' \
2      --header 'x-api-key: my_key' \
3      --header 'anthropic-version: 2023-06-01' \
4      --header 'content-type: application/json' \
5      --data '{
6          "model": "claude-3-5-sonnet-20240620",
7          "max_tokens": 4096,
8          "messages": [
9              {
10                 "role": "user",
11                 "content": "ALL THE HTML / TEXT WITH EVENTS"
12             }
13          ],
14          "tools": [
15              {
```

```

16         "name": "create_event_tool",
17         "description": "If the user needs to create one or more events this tool hel\
18 ps",
19         "input_schema": {
20             "type": "object",
21             "properties": {
22                 "events": {
23                     "description": "Array of event objects",
24                     "type": "array",
25                     "items": {
26                         "type": "object",
27                         "properties": {
28                             "start_time": {
29                                 "type": "string",
30                                 "description": "Start time of the event"
31                             },
32                             "end_time": {
33                                 "type": "string",
34                                 "description": "End time of the event"
35                             },
36                             "title": {
37                                 "type": "string",
38                                 "description": "Title of the event"
39                             },
40                             "location": {
41                                 "type": "string",
42                                 "description": "Location of the event"
43                             },
44                             "summary": {
45                                 "type": "string",
46                                 "description": "Summary or description of the event"
47                             }
48                         }
49                     }
50                 }
51             },
52             "required": [
53                 "events"
54             ]
55         }
56     }
57 ]
58 }'

```

And it will reply as follows:

```

1  {
2      "id": "msg_01GDtWoCBAnDpdLmLbJZGp2u",
3      "type": "message",
4      "role": "assistant",
5      "model": "claude-3-5-sonnet-20240620",
6      "content": [
7          {
8              "type": "text",
9              "text": "Certainly! I'll parse the events from the Dallas Cowboys schedu\
10 le and create them using the create_event_tool. I'll focus on the regular season gam
11 es, as there are quite a few events. Let's create these events:"
12          },
13          {
14              "type": "tool_use",
15              "id": "toolu_01B9W2pzR5RqDj11QciiseJq",
16              "name": "create_event_tool",
17              "input": {
18                  "events": [
19                      {
20                          "title": "Dallas Cowboys at Cleveland Browns",
21                          "start_time": "2024-09-08T15:25:00-05:00",
22                          "end_time": "2024-09-08T18:25:00-05:00",
23                          "location": "Cleveland Browns Stadium",
24                          "summary": "NFL Week 1: Dallas Cowboys at Cleveland Browns"
25                      },
26                      {
27                          "title": "Dallas Cowboys vs New Orleans Saints",
28                          "start_time": "2024-09-15T12:00:00-05:00",
29                          "end_time": "2024-09-15T15:00:00-05:00",
30                          "location": "AT&T Stadium",
31                          "summary": "NFL Week 2: Dallas Cowboys vs New Orleans Saints"
32                      },
33                      {
34                          "title": "Dallas Cowboys vs Baltimore Ravens",
35                          "start_time": "2024-09-22T15:25:00-05:00",
36                          "end_time": "2024-09-22T18:25:00-05:00",
37                          "location": "AT&T Stadium",
38                          "summary": "NFL Week 3: Dallas Cowboys vs Baltimore Ravens"
39                      },
40                      {
41                          "title": "Dallas Cowboys at New York Giants",

```

```

42         "start_time": "2024-09-26T19:15:00-05:00",
43         "end_time": "2024-09-26T22:15:00-05:00",
44         "location": "MetLife Stadium",
45         "summary": "NFL Week 4: Dallas Cowboys at New York Giants"
46     },
47     {
48         "title": "Dallas Cowboys at Pittsburgh Steelers",
49         "start_time": "2024-10-06T19:20:00-05:00",
50         "end_time": "2024-10-06T22:20:00-05:00",
51         "location": "Acrisure Stadium",
52         "summary": "NFL Week 5: Dallas Cowboys at Pittsburgh Steeler\
53 s"
54     }
55 ]
56 }
57 }
58 ],
59 "stop_reason": "tool_use",
60 "stop_sequence": null,
61 "usage": {
62     "input_tokens": 2652,
63     "output_tokens": 593
64 }
65 }

```

So our Tool will have to iterate over these results and make events. It really is that easy, BUT there is a catch.

Orchestration

It is really easy to call the LLM with PHP code like:

```

1  <?php
2
3  $result = $client->chat()->create([
4      'model' => 'gpt-4',
5      'messages' => [
6          ['role' => 'user', 'content' => 'Hello!'],
7      ],
8  ]);

```

But we want to do more than that; we want to be able to chain the results of that to either one or more tools and a final response or just the final response. That is where an Orchestration-like class comes in.

```

1  <?php
2  public function handle(Chat $chat, string $prompt): Message
3  {
4      $chat->addInput(
5          message: $prompt,
6          role: RoleEnum::User
7      );
8
9      $messageInDto = MessageInDto::from([
10         'content' => $prompt,
11         'role' => 'user',
12     ]);
13
14     $response = LlmDriverFacade::driver($chat->getDriver())
15         ->chat([
16             $messageInDto,
17         ]);
18
19     if (! empty($response->tool_calls)) {
20         Log::info('Orchestration Tools Found', [
21             'tool_calls' => collect($response->tool_calls)
22                 ->pluck('name')->toArray(),
23         ]);
24
25         $count = 1;
26         foreach ($response->tool_calls as $tool_call) {
27             Log::info('[LaraChain] - Tool Call '.$count, [
28                 'tool_call' => $tool_call->name,
29                 'tool_count' => count($response->tool_calls),
30             ]);
31
32             $message = $chat->addInput(
33                 message: $response->content ?? 'Calling Tools', //@NOTE ollama, open\
34                 ai blank but claude needs this : (
35                 role: RoleEnum::Assistant,
36                 tool: $tool_call->name,
37                 tool_id: $tool_call->id,
38                 args: $tool_call->arguments,
39             );

```

```

40
41     $tool = app()->make($tool_call->name);
42     $results = $tool->handle($message);
43     $message->updateQuietly([
44         'role' => RoleEnum::Tool,
45         'body' => $results->content,
46     ]);
47     $count++;
48 }
49
50 $messages = $chat->getChatResponse();
51
52 $response = LlmDriverFacade::driver($chat->getDriver())
53     ->chat($messages);
54
55 return $chat->addInput(
56     message: $response->content,
57     role: RoleEnum::Assistant,
58 );
59
60 } else {
61     Log::info('No Tools found just gonna chat');
62     $assistantMessage = $chat->addInput(
63         message: $response->content ?? 'Calling Tools',
64         role: RoleEnum::Assistant
65     );
66
67     return $assistantMessage;
68 }
69 }

```

As we read in earlier chapters, we have to maintain state, so I do that typically using a Chat and Message model. This allows me to chain all of this together based on that.

Let's look at this line by line:

```

1  <?php
2  $messages = $chat->getChatResponse();

```

Here we are just telling the Chat model to give us the thread of messages so we can continue our chat with the LLM.

```

1  <?php
2  $response = LlmDriverFacade::driver($message->getDriver())
3      ->chat($messages);

```

Here we pass our message to the LLMs. But in here is the LLM Driver loading the tools we have.

```

1  <?php
2  $this->getFunctions();

```

We will look at this one more shortly, but just consider it will load the tool we are going to make, register it with our Driver so it knows to offer this tool as an option in the chat interaction.

Let's zoom into chat with tools:

```

1  <?php
2  foreach ($response->tool_calls as $tool_call) {
3      $message = $chat->addInput(
4          message: $response->content ?? 'Calling Tools', //ollama, openai blank but c\
5  laude needs this :(
6          role: RoleEnum::Assistant,
7          show_in_thread: false,
8          meta_data: MetaDataDto::from([
9              'tool' => $tool_call->name,
10             'tool_id' => $tool_call->id,
11             'driver' => $chat->getDriver(),
12             'args' => $tool_call->arguments,
13         ]),
14      );
15      $tool = app()->make($tool_call->name);
16      $results = $tool->handle($message);
17      $message->updateQuietly([
18          'is_chat_ignored' => true,
19          'role' => RoleEnum::Tool,
20          'body' => $results->content,
21      ]);
22      $count++;
23  }

```

You can see it just iterates over tool(s) until it is done, then finishes off with:

```

1  <?php
2  $messages = $chat->getChatResponse();
3
4  /**
5   * @NOTE
6   * I have to continue to pass in tools once used above
7   * Since Claude needs them.
8   */
9  $response = LlmDriverFacade::driver($chat->getDriver())
10     ->setToolType(ToolTypes::Source)
11     ->chat($messages);
12
13  $assistantMessage = $chat->addInput(
14      message: $response->content,
15      role: RoleEnum::Assistant,
16      show_in_thread: true,
17      meta_data: null,
18      tools: null
19  );

```

Getting the message thread again, now with the updated results and finishing off the request to the LLMs now that all the tools are processed.

Deep Dive

Ok, now that the overview is seen, let's put all of this together.

Start with the Request

The request comes in via a Job, Controller, or whatever. It then does the work, in this case to scrape a webpage. Then we take that content and pass it into the Orchestrate class like this:

<https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/Orchestration/Orchestrate.php#L14>

```

1  <?php
2  use Facades\App\Services\LlmServices\Orchestration\Orchestrate;
3  Orchestrate::handle($chat, $prompt);

```

NOTE: These classes have tests if you want to read how they work

So what is happening here is we are always dealing with Chat. I do this so we have a root to the thread. Sometimes we will use Message, but in this case, we are working at the Chat level. You will see in a prior chapter we made the Chat model.

And the \$prompt is the start of this work, which by now has gone through a few revisions. The prompt is basically the question the user is asking of the system. The question has a tool in it by happenstance in that the user might not even know the tool exists; they just asked for something that happened to have the needed tools.

The Orchestrate class then hands the question and the tools to the LLM to see if it can answer it with or without tools.

If the response comes back tool-based, e.g.:

```

1  {
2      "id": "msg_01GDtWoCBAnDpdLmLbJZGp2u",
3      "type": "message",
4      "role": "assistant",
5      "model": "claude-3-5-sonnet-20240620",
6      "content": [
7          {
8              "type": "text",
9              "text": "Certainly! I'll parse the events from the Dallas Cowboys schedu\
10 le and create them using the create_event_tool. I'll focus on the regular season gam
11 es, as there are quite a few events. Let's create these events:"
12          },
13          {
14              "type": "tool_use",
15              "id": "toolu_01B9W2pzR5RqDj11QciiseJq",
16              "name": "create_event_tool",
17              "input": {
18                  "events": [
19                      {
20                          "title": "Dallas Cowboys at Cleveland Browns",
21                          "start_time": "2024-09-08T15:25:00-05:00",
22                          "end_time": "2024-09-08T18:25:00-05:00",
23                          "location": "Cleveland Browns Stadium",
24                          "summary": "NFL Week 1: Dallas Cowboys at Cleveland Browns"
25                      },
26                      {
27                          "title": "Dallas Cowboys vs New Orleans Saints",
28                          "start_time": "2024-09-15T12:00:00-05:00",
29                          "end_time": "2024-09-15T15:00:00-05:00",
30                          "location": "AT&T Stadium",
31                          "summary": "NFL Week 2: Dallas Cowboys vs New Orleans Saints"

```

```

32         },
33         {
34             "title": "Dallas Cowboys vs Baltimore Ravens",
35             "start_time": "2024-09-22T15:25:00-05:00",
36             "end_time": "2024-09-22T18:25:00-05:00",
37             "location": "AT&T Stadium",
38             "summary": "NFL Week 3: Dallas Cowboys vs Baltimore Ravens"
39         },
40         {
41             "title": "Dallas Cowboys at New York Giants",
42             "start_time": "2024-09-26T19:15:00-05:00",
43             "end_time": "2024-09-26T22:15:00-05:00",
44             "location": "MetLife Stadium",
45             "summary": "NFL Week 4: Dallas Cowboys at New York Giants"
46         },
47         {
48             "title": "Dallas Cowboys at Pittsburgh Steelers",
49             "start_time": "2024-10-06T19:20:00-05:00",
50             "end_time": "2024-10-06T22:20:00-05:00",
51             "location": "Acrisure Stadium",
52             "summary": "NFL Week 5: Dallas Cowboys at Pittsburgh Steeler\
53 s"
54         }
55     ]
56 }
57 }
58 ],
59 "stop_reason": "tool_use",
60 "stop_sequence": null,
61 "usage": {
62     "input_tokens": 2652,
63     "output_tokens": 593
64 }
65 }

```

Then we process the tools requests.

The LLM Client

All of this happens in the Client, we will focus on the Claude client for now. <https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/ClaudeClient.php>

We cover this in an earlier chapter

```

1  <?php
2  $response = LlmDriverFacade::driver($chat->getDriver())
3      ->chat([
4          $messageInDto,
5      ]);

```

The driver is set to claude so we then talk to the CluadeClient chat method:

```

1  <?php
2  /**
3   * @param MessageInDto[] $messages
4   */
5  public function chat(array $messages): CompletionResponse
6  {
7      $model = $this->getConfig('claude')['models']['completion_model'];
8      $maxTokens = $this->getConfig('claude')['max_tokens'];
9
10     Log::info('LlmDriver::Claude::chat');
11
12     $messages = $this->remapMessages($messages);
13
14     $payload = [
15         'model' => $model,
16         'max_tokens' => $maxTokens,
17         'messages' => $messages,
18     ];
19
20     $payload = $this->modifyPayload($payload);
21
22     $results = $this->getClient()->post('/messages', $payload);
23
24     if (! $results->ok()) {
25         $error = $results->json()['error']['type'];
26         $message = $results->json()['error']['message'];
27         Log::error('Claude API Error Chat', [
28             'type' => $error,
29             'message' => $message,
30         ]);
31
32         throw new \Exception('Claude API Error Chat');
33     }
34

```

```

35     return ClaudeCompletionResponse::from($results->json());
36 }

```

We cover a lot of this in the other chapters but for now I want to focus on:

```

1  <?php
2  $payload = $this->modifyPayload($payload);

```

This is how we turned the “CreateEventTool”’s into an LLM agnostic tool payload.

That modify payload calls the parent class to get the functions and then it’s own remapFunctions all the clients have their own remapping though most now just do what ChatGPT does.

The code below does not matter as much as realizing we are taking the formatting of the tools parameters (will show after this) and turning them into the JSON format Claude API requires.

```

1  <?php
2  public function remapFunctions(array $functions): array
3  {
4      return collect($functions)->map(function ($function) {
5          $properties = [];
6          $required = [];
7
8          $type = data_get($function, 'parameters.type', 'object');
9
10         foreach (data_get($function, 'parameters.properties', []) as $property) {
11             $name = data_get($property, 'name');
12
13             if (data_get($property, 'required', false)) {
14                 $required[] = $name;
15             }
16
17             $subType = data_get($property, 'type', 'string');
18             $properties[$name] = [
19                 'description' => data_get($property, 'description', null),
20                 'type' => $subType,
21             ];
22
23             if ($subType === 'array') {
24                 $subItems = $property['properties'][0]->properties; //stop at th\
25 is for now
26                 $subItemsMapped = [];
27                 foreach ($subItems as $subItemKey => $subItemValue) {
28                     $subItemsMapped[$subItemValue->name] = [

```

```

29         'type' => $subItemValue->type,
30         'description' => $subItemValue->description,
31     ];
32     }
33
34     $properties[$name]['items'] = [
35         'type' => 'object',
36         'properties' => $subItemsMapped,
37     ];
38     }
39 }
40
41 $itemsOrProperties = $properties;
42
43 return [
44     'name' => data_get($function, 'name'),
45     'description' => data_get($function, 'description'),
46     'input_schema' => [
47         'type' => 'object',
48         'properties' => $itemsOrProperties,
49         'required' => $required,
50     ],
51 ];
52 }->values()->toArray();
53 }

```

At this point we just make an HTTP request to the API and return it.

NOTE: I talk about `MessageInDto` and `ClaudeCompletionResponse` in another section but the goal here is ALL llms take the same Data Object and return the same data object. In this case `ClaudeCompletionResponse` extends `CompletionResponse`. This is really powerful in that we can now talk to and return the same object from any LLM we plug in.

The CrateEventTool

Ok now that we see how Orchestrate works and how the `ClaudeClient` reformats the functions lets look at the foundation to tools and in this case the `CreateEventTool`

<https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/Functions/CreateEventTool.php>

You will see this extends `FunctionContract` so let's start with that below:

<https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/Functions/FunctionContract.php>

```
1  <?php
2
3  namespace App\Services\LlmServices\Functions;
4
5  use App\Models\Message;
6
7  abstract class FunctionContract
8  {
9      protected string $name;
10
11      protected string $description;
12
13      protected string $type = 'object';
14
15      abstract public function handle(
16          Message $message,
17      ): FunctionResponse;
18
19      public function getFunction(): FunctionDto
20      {
21          return FunctionDto::from(
22              [
23                  'name' => $this->getName(),
24                  'description' => $this->getDescription(),
25                  'parameters' => [
26                      'type' => $this->type,
27                      'properties' => $this->getProperties(),
28                  ],
29              ]
30          );
31      }
32
33      public function getName(): string
34      {
35          return $this->name;
36      }
37
38      public function getKey(): string
39      {
40          return $this->name;
```

```

41     }
42
43     public function getDescription(): string
44     {
45         return $this->description;
46     }
47
48     public function getParameters(): array
49     {
50         return $this->getProperties();
51     }
52
53     /**
54      * @return PropertyDto[]
55      */
56     abstract protected function getProperties(): array;
57 }

```

Our tools have to have a name and good description

See some docs at <https://docs.anthropic.com/en/docs/build-with-claude/tool-use>

Then we have the handle that all tools have that will return the object FunctionDto

```

1  <?php
2
3  namespace App\Services\LlmServices\Functions;
4
5  class FunctionDto extends \Spatie\LaravelData\Data
6  {
7      public function __construct(
8          public string $name,
9          public string $description,
10         public ParametersDto $parameters,
11     ) {}
12 }

```

Ok now lets see how the CreateEventTool extends this:

```
1  <?php
2
3  namespace App\Services\LlmServices\Functions;
4
5  use App\Models\Event;
6  use App\Models\Message;
7  use Carbon\Carbon;
8  use Illuminate\Support\Facades\Log;
9
10 class CreateEventTool extends FunctionContract
11 {
12     protected string $name = 'create_event_tool';
13
14     protected string $description = 'If the user needs to create one or more events \
15 this tool help';
16
17     public function handle(
18         Message $message): FunctionResponse
19     {
20         Log::info('CreateEventTool called');
21
22         $args = $message->args;
23
24         $eventArray = data_get($args, 'events', []);
25
26         foreach ($eventArray as $event) {
27             $start_date = null;
28             $end_date = null;
29             $description = data_get($event, 'description', null);
30             $start_time = data_get($event, 'start_time', null);
31             $end_time = data_get($event, 'end_time', null);
32             $location = data_get($event, 'location', null);
33             $title = data_get($event, 'title', 'No Title Found');
34             $all_day = data_get($event, 'all_day', false);
35
36             if ($start_time != '' || $start_time != null) {
37                 $start_date = Carbon::parse($start_time)->format('Y-m-d');
38                 $start_time = Carbon::parse($start_time)->format('H:i:s');
39             }
40
41             if ($end_time != '' || $end_time != null) {
42                 $end_date = Carbon::parse($end_time)->format('Y-m-d');
43                 $end_time = Carbon::parse($end_time)->format('H:i:s');
```



```

44         }
45
46         Event::updateOrCreate([
47             'title' => $title,
48             'start_date' => $start_date,
49             'start_time' => $start_time,
50             'location' => $location,
51         ],
52         [
53             'description' => $description,
54             'end_date' => $end_date,
55             'end_time' => $end_time,
56             'all_day' => $all_day,
57         ]);
58     }
59
60     return FunctionResponse::from([
61         'content' => json_encode($eventArray),
62         'prompt' => $message->body,
63     ]);
64 }
65
66 /**
67  * @return PropertyDto[]
68  */
69 protected function getProperties(): array
70 {
71     return [
72         new PropertyDto(
73             name: 'events',
74             description: 'Array of event objects',
75             type: 'array',
76             required: true,
77             properties: [
78                 new PropertyDto(
79                     name: 'items',
80                     description: 'Event object',
81                     type: 'object',
82                     required: true,
83                     properties: [
84                         new PropertyDto(
85                             name: 'start_time',
86                             description: 'Start time of the event',

```

```

87         type: 'string',
88         required: true
89     ),
90     new PropertyDto(
91         name: 'end_time',
92         description: 'End time of the event',
93         type: 'string',
94         required: false
95     ),
96     new PropertyDto(
97         name: 'title',
98         description: 'Title of the event',
99         type: 'string',
100        required: true
101    ),
102    new PropertyDto(
103        name: 'location',
104        description: 'Location of the event',
105        type: 'string',
106        required: false
107    ),
108    new PropertyDto(
109        name: 'description',
110        description: 'Description of the event',
111        type: 'string',
112        required: true
113    ),
114    ],
115    ),
116    ],
117    ),
118    ];
119 }
120 }

```

Let's look first at the function/tool name `create_event_tool` this is key since we register this in <https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/LlmServiceProvider.php> which we registered in an earlier chapter.

```
1  <?php
2      public function boot(): void
3      {
4          $this->app->bind('llm_driver', function () {
5              return new LlmDriverClient;
6          });
7
8          $this->app->bind('create_event_tool', function () {
9              return new CreateEventTool;
10         });
11
12     }
```

Then we register the tool in the BaseClient: <https://github.com/alnutile/php-llms/blob/main/app/Services/LlmServices/BaseClient.php>

```
1  <?php
2      public function getFunctions(): array
3      {
4          $functions = collect(
5              [
6                  new CreateEventTool,
7              ]
8          );
9
10         return $functions->transform(
11             /** @phpstan-ignore-next-line */
12             function (FunctionContract $function) {
13                 return $function->getFunction();
14             }
15         )->toArray();
16     }
```

So now when any LLM Client runs and does getFunctions it will get the needed functions from here. And when Orchestrate iterates over the tools

```

1  <?php
2  //app/Services/LlmServices/Orchestration/Orchestrate.php
3  foreach ($response->tool_calls as $tool_call) {
4      Log::info('[LaraChain] - Tool Call ' . $count, [
5          'tool_call' => $tool_call->name,
6          'tool_count' => count($response->tool_calls),
7      ]);
8
9      $message = $chat->addInput(
10         message: $response->content ?? 'Calling Tools', //@NOTE ollama, openai blank\
11         but claude needs this :(
12         role: RoleEnum::Assistant,
13         tool: $tool_call->name,
14         tool_id: $tool_call->id,
15         args: $tool_call->arguments,
16     );
17
18     $tool = app()->make($tool_call->name);
19     $results = $tool->handle($message);
20     $message->updateQuietly([
21         'role' => RoleEnum::Tool,
22         'body' => $results->content,
23     ]);
24     $count++;
25 }

```

It will call that tool and use handle

Note too it sets the \$message to have args and other tool info that is key.

At this point we see that the tool gets the args and iterates over them, knowing the format will match it's parameters.

```

1  <?php
2  $eventArray = data_get($args, 'events', []);
3
4  foreach ($eventArray as $event) {
5      $start_date = null;
6      $end_date = null;
7      $description = data_get($event, 'description', null);
8      $start_time = data_get($event, 'start_time', null);
9      $end_time = data_get($event, 'end_time', null);

```

```

10     $location = data_get($event, 'location', null);
11     $title = data_get($event, 'title', 'No Title Found');
12     $all_day = data_get($event, 'all_day', false);
13
14     if ($start_time != '' || $start_time != null) {
15         $start_date = Carbon::parse($start_time)->format('Y-m-d');
16         $start_time = Carbon::parse($start_time)->format('H:i:s');
17     }
18
19     if ($end_time != '' || $end_time != null) {
20         $end_date = Carbon::parse($end_time)->format('Y-m-d');
21         $end_time = Carbon::parse($end_time)->format('H:i:s');
22     }
23
24     Event::updateOrCreate([
25         'title' => $title,
26         'start_date' => $start_date,
27         'start_time' => $start_time,
28         'location' => $location,
29     ],
30     [
31         'description' => $description,
32         'end_date' => $end_date,
33         'end_time' => $end_time,
34         'all_day' => $all_day,
35     ]);
36
37     return FunctionResponse::from([
38         'content' => json_encode($eventArray),
39         'prompt' => $message->body,
40     ]);
41 }

```

Finally all Tools return FunctionResponse so Orchestrate can use the response if needed.

Let's look now at the `getProperties` of the `CreateEventTool`. As noted in a previous chapter we abstract out the parameters so that it starts as an array of `PropertyDto` objects. That later any LLM can transform.

```
1  <?php
2  /**
3   * @return PropertyDto[]
4   */
5  protected function getProperties(): array
6  {
7      return [
8          new PropertyDto(
9              name: 'events',
10             description: 'Array of event objects',
11             type: 'array',
12             required: true,
13             properties: [
14                 new PropertyDto(
15                     name: 'items',
16                     description: 'Event object',
17                     type: 'object',
18                     required: true,
19                     properties: [
20                         new PropertyDto(
21                             name: 'start_time',
22                             description: 'Start time of the event',
23                             type: 'string',
24                             required: true
25                         ),
26                         new PropertyDto(
27                             name: 'end_time',
28                             description: 'End time of the event',
29                             type: 'string',
30                             required: false
31                         ),
32                         new PropertyDto(
33                             name: 'title',
34                             description: 'Title of the event',
35                             type: 'string',
36                             required: true
37                         ),
38                         new PropertyDto(
39                             name: 'location',
40                             description: 'Location of the event',
41                             type: 'string',
42                             required: false
43                         ),
```

```
44         new PropertyDto(  
45             name: 'description',  
46             description: 'Description of the event',  
47             type: 'string',  
48             required: true  
49         ),  
50     ],  
51 ),  
52 ],  
53 ),  
54 ];  
55 }
```

Wrapping It Up

So we see here a few patterns that lay the foundation to all our interactions with LLMs. Some of this covered in the previous chapters.

- Our LLM Clients have a DataObject in and DataObject out that they all share.
- Our Tools have parameters that are objects that then any LLM Client can transform to the needed JSON for that API.
- The messages we get back have to be transformed as well and like above the different LLM Clients need a different format to the messages.
- Tools are powerful and with the Orchestrate class we can chain them together naturally or forced (more on that later)

Making Tools LLM-Agnostic

Summary: Building on the previous chapter, we focus on creating tools that are not tied to any specific LLM. This allows you to switch between LLMs like OpenAI, Ollama, and Claude without rewriting your code. The chapter includes examples of how to structure your tools to be flexible and adaptable.

Chaining Tools

Summary: This chapter explores how to iterate over tools in complex workflows. We discuss real-world scenarios where multiple tools need to be chained together to achieve a final result. Examples include processing web content and generating images, highlighting the importance of tool chaining.

Web Searches with LLMs

Summary: A case study of automating web searches using LLMs. The chapter walks through setting up a system that regularly searches the web for relevant content, processes the results, and takes further actions based on predefined criteria.^a

Marketing Assistant

Summary:

This tool will help us make events via a prompt about helping us with marketing. This will allow the system to then track the status of those events, let us know what is coming up, and ideally automate some of the actions of those events.

Hand Written Notes to Task List

Summary: In this chapter, we explore how to leverage LLMs to transform handwritten notes and voice files into actionable task lists. We'll discuss the process of uploading these inputs, using LLMs to interpret and organize the content, and generating structured tasks that can be integrated into your workflow. Whether you're dealing with meeting notes, brainstorming sessions, or voice memos, this chapter will show you how to automate the process, saving time and ensuring nothing falls through the cracks.

This chapter can demonstrate the practical utility of LLMs in day-to-day productivity, adding another layer of automation to your development toolkit.

RAG and Vector

Summary: In this chapter, we delve into the concept of Retrieval-Augmented Generation (RAG) systems, an architectural approach designed to enhance the effectiveness of LLM applications. We'll cover the foundational aspects of RAG, including how to set it up easily within your projects. However, it's important to understand that while RAG can significantly improve LLM performance, it's not an end in itself. We'll discuss the limitations of RAG and how it should be seen as one tool among many in your development toolkit, rather than a complete solution.

This chapter will provide valuable insights into the benefits and limitations of RAG systems, helping readers make informed decisions about when and how to use this approach in their projects.

SiteMap - Feed to RAG

Summary: In this chapter, we explore how to take a website's sitemap and integrate it into a Retrieval-Augmented Generation (RAG) system. We'll discuss the process of chunking the sitemap content to make it searchable, allowing you to interact with it through search, chat, or even by creating a chatbot. By the end of this chapter, you'll have the knowledge to feed customer content into a RAG system, enabling powerful new ways to engage with the material.

This chapter will guide readers through a practical application of RAG, demonstrating how to use it to enhance customer interactions and content management.

Agents to Evaluate the Results

Summary: This chapter focuses on evaluating the accuracy of LLM outputs. We discuss methods for assessing the quality of the results and provide strategies for improving reliability over time.

Standards Checking Tool

Summary: This chapter introduces a standards checking tool that evaluates documents against predefined standards. We cover how to set up the tool to work with any LLM, how to interpret the results, and how to integrate this tool into your content management workflow

Final Thoughts

By now, I've shared with you my experiences and insights on integrating LLMs into your day-to-day solutions. We've covered not only the practical aspects of working with these powerful tools but also the strategies for making them a reliable part of your development workflow.

From setting up local environments to save costs, to testing and mocking with confidence, to abstracting LLMs with data objects—these are skills that will empower you to harness the full potential of LLMs. By adopting these techniques, you'll be able to approach complex problems with a new perspective, equipped with tools that simplify and enhance your work.

As we conclude, it's clear that while the initial hype around LLMs might fade, the practical applications and possibilities they offer are here to stay. Whether you're solving age-old problems or tackling new challenges, the integration of LLMs into your toolkit opens up a world of opportunities to innovate and deliver more effective solutions to your clients.

I hope this book has not only provided you with the technical know-how but also inspired you to explore and push the boundaries of what's possible with LLMs. The skills you've gained will serve you well as you continue to evolve in this ever-changing landscape of technology.

Thank you for joining me on this journey, and I look forward to seeing how you apply these insights to create powerful, intelligent solutions in your work.

This expanded conclusion aims to reinforce the key takeaways from the book while encouraging your readers to apply what they've learned in their own projects.