

[COMPLETE] RAG and Vector

What You Will Learn About Building a RAG System

In this chapter, you will learn the steps involved in building a Retrieval-Augmented Generation (RAG) system using a vector database.

- The chapter explains how to set up a **vector database** using Postgres and the pgvector-php library. The sources use the example of a database about the sport of pickleball.
- You will learn how to ingest data into the database, including how to **chunk** it into smaller pieces for better search results.
- The sources describe the process of **embedding** the data using an LLM API, such as Ollama or OpenAI. Embedding converts text into a numerical representation that can be processed by a machine learning model.
- The chapter provides a step-by-step guide on how to query the vector database using an embedded question. You will learn about the different distance types used in vector searches, such as cosine distance.
- The sources also cover how to incorporate images and audio into a RAG system. They explain how to extract text from images using OpenAI's API and how to transcribe audio using OpenAI's Whisper model.
- Finally, the chapter discusses the importance of structuring data in a RAG system to improve the accuracy and relevance of search results.

Retrieval-Augmented Generation (RAG) is a technique that enhances AI language models by combining their built-in knowledge with relevant information retrieved from an external database. This technique allows the AI to provide more accurate, up-to-date, and context-specific responses by drawing on both its general knowledge and specific, retrieved facts. —Claude

Why RAG?

This is a good question, and one I sometimes skip over because I am like, “Why not?” Most applications use search, like Algolia, and it does an amazing job finding content - searching “Laravel Facades” will give me a list of document pages with that information.



In this case, I want my user to be able to ask questions about the data (that might be private or public), however I want the LLM to confine its results to just that data.

For example:

“I want to make a facade from my LLMDriver class. How can I do that?” I want the response to include not only documents that have the word “Facade”, but also more about the context of the question, allowing the system to return a greater depth of results that I then pass to the LLM to answer the question.

Remember, this is not just about chatting with your data; you can also automate tasks around it. You could use your Vector database to find content related to other content that otherwise might be hard to do on keyword search alone. For example, your system gets a support question about how to refund a customer; your RAG system does a vector query, finds several answers, and then creates a very nicely formatted and worded email to help the customer step-by-step.

One thing I want to emphasize is that I think RAG is great, but I believe it is not an end in itself, but rather a way to store your data and use it for other tasks like automating tags, finding related articles, user shopping preferences, customer support, email automation, and more.

What'll we do with RAG?

I recently wrote a Medium article [Laravel RAG System in 4 Steps](https://medium.com/@alnutile/laravel-rag-in-4-steps-6264b4df173a)¹. This chapter will dig into the details of each step mentioned in the article, but also go further in a few ways. First, we will go into processing images, then how to make tables as HTML, and lastly, hybrid search so we get a bit more out of our search steps.

Here's a quick overview:

1. Set up the vector database so we can do searches. Unlike keyword searches, vector searches can understand context and meaning, finding related items even if they don't share exact words.
2. Pass those results to the LLM using a prompt that keeps it from drifting or hallucinating.
3. Allow the user to have a threaded discussion.

So how do we get there?

1. We use Postgres as our vector database.
2. We install `pgvector-php` to give us the ability to do vector searches.
3. We import the required data, chunking it so that when we get results, it is small sets of data and not large full pages. This way, we know which section of the page matches.
4. We take the user's question and embed it using an LLM API, and then use it to search.
5. We then take the results and ask the LLM to stay within them to answer the original question.

Later, we will add a step to maximize images and tables, so that we can show them differently in the UI and utilize them better in our responses.

Getting Your Application Database Set up

We are using Postgres as our main database here. (You can use it as a second database if you already use MySQL.) Laravel has an easy way to set up Postgres, and it will be our central database for now. Once Postgres is installed, the library will do much of the heavy lifting regarding vector. When the database is ready, we'll get data into the system. We'll ingest a bunch of markdown files from a folder. Later, we'll try more complex PDFs.

The easiest ways to run Postgres locally for this is either [Laravel HERD](https://herd.laravel.com/)², the paid version, or [DBngin](https://docs.tableplus.com/utilities/dbngin)³ for free.

To deploy, you can easily use the [Forge](https://forge.laravel.com/)⁴ Postgres option on your existing server or deploy it as a standalone server. Then, if you follow the installation instructions of <https://github.com/pgvector/pgvector-php> you will be all set. Keep in mind, the migration that it adds to your system needs permissions to run:

¹<https://medium.com/@alnutile/laravel-rag-in-4-steps-6264b4df173a>

²<https://herd.laravel.com/>

³<https://docs.tableplus.com/utilities/dbngin>

⁴<https://forge.laravel.com/>

Figure 149. Migration

```
1 DB::statement('CREATE EXTENSION IF NOT EXISTS vector');
```

So as long as you are using the default user and permissions from Herd, DBngin and Forge, you should be fine.

For example, my local `.env` looks like this:

Figure 150. Local `.env`

```
1 DB_CONNECTION=pgsql
2 DB_HOST=127.0.0.1
3 DB_PORT=5432
4 DB_DATABASE=laravel_11ms
5 DB_USERNAME=postgres
6 DB_PASSWORD=
```

NOTE: When using Postgres, Laravel does not seem to be able to create the schema for you, as with MySQL. Using a tool like [TablePlus](https://tableplus.com/)⁵, you can easily add the new schema.

Now that the database is ready, let's get some data and ingest it into the database.

Data and Ingestion

We're going to use the sport Pickleball as our example source data. We are going to start with a small collection of publicly-available files I have about Pickleball in markdown format; you can see them in `storage/data/pickleball`. As with most of this data, we save it in two tables. One will be called `documents` and the other `chunks`. Then, we embed the data using the Ollama API (the same as OpenAI) and store those embeddings. Once done, we will have one document per file and multiple chunks per file size.

Document Table

Let us start with the document model. The model is a straightforward table - you can see it all here:

⁵<https://tableplus.com/>

Figure 151. database/migrations/2024_09_15_211042_create_documents_table.php

```
1  <?php
2
3  use Illuminate\Database\Migrations\Migration;
4  use Illuminate\Database\Schema\Blueprint;
5  use Illuminate\Support\Facades\Schema;
6
7  return new class extends Migration
8  {
9      /**
10       * Run the migrations.
11       */
12     public function up(): void
13     {
14         Schema::create('documents', function (Blueprint $table) {
15             $table->id();
16             $table->string('title')->nullable();
17             $table->longText('summary')->nullable();
18             $table->longText('content')->nullable();
19             $table->timestamps();
20         });
21     }
22
23     /**
24      * Reverse the migrations.
25      */
26     public function down(): void
27     {
28         Schema::dropIfExists('documents');
29     }
30 };
```

Once that is in place, we can now make the more complex model chunks.

Chunk Table

A lot of this will be new, so let us dig in one file at a time.

Chunks Table Migration

Figure 152. database/migrations/2024_09_15_211349_create_chunks_table.php

```

1 Schema::create('chunks', function (Blueprint $table) {
2     $table->id();
3     $table->string('guid');
4     $table->string('sort_order')->default(1);
5     $table->longText('content')->nullable();
6     $table->longText('summary')->nullable();
7     $table->vector('embedding_768', 768)->nullable();
8     $table->vector('embedding_1536', 1536)->nullable();
9     $table->vector('embedding_2048', 2048)->nullable();
10    $table->vector('embedding_3072', 3072)->nullable();
11    $table->vector('embedding_1024', 1024)->nullable();
12    $table->vector('embedding_4096', 4096)->nullable();
13    $table->integer('page_number')->nullable();
14    $table->foreignIdFor(\App\Models\Document::class, 'document_id');
15    $table->timestamps();
16 });

```

The unique field here is `vector` and the vector types. Here, I allow the system to have more than one vector type, depending on our model. In this initial system example, we use the one embedding type. But later on in other chapters, we will use other models for cost (local vs. OpenAI) or speed. When we save our data to the model, I will explain more about which LLM model we use and the embedding size.

The Factory Database/Factories/ChunkFactory.php

Figure 153. database/factories/ChunkFactory.php

```

1 public function definition(): array
2 {
3     $embeddings = get_fixture('embedding_response.json');
4
5     return [
6         'guid' => fake()->uuid(),
7         'content' => fake()->sentence(10),
8         'page_number' => fake()->numberBetween(1, 100),
9         'section_number' => fake()->numberBetween(1, 100),
10        'original_content' => fake()->sentence(10),
11        'summary' => fake()->sentence(5),
12        'embedding_3072' => data_get($embeddings, 'data.0.embedding'),
13        'embedding_1536' => null,
14        'embedding_2048' => null,

```

```

15         'embedding_4096' => null,
16         "document_id" => Document::factory(),
17     ];
18 }

```

Here, we can use the fixture file to generate a vector for the factory. A segment of that file looks like this:

Figure 154. tests/fixtures/embedding_response.json

```

1  {
2      "object": "list",
3      "data": [
4          {
5              "object": "embedding",
6              "index": 0,
7              "embedding": [
8                  -0.02932045,
9                  -0.019364825,
10                 // and on and on and on...

```

The Model App/Models/Chunk.php

And, of course, the model file:

Figure 155. app/Models/Chunk.php

```

1  <?php
2
3  namespace App\Models;
4
5  use Illuminate\Database\Eloquent\Factories\HasFactory;
6  use Illuminate\Database\Eloquent\Model;
7  use Pgvector\Laravel\HasNeighbors;
8  use Pgvector\Laravel\Vector;
9
10 class Chunk extends Model
11 {
12     use HasFactory;
13     use HasNeighbors;
14
15     protected $guarded = [];
16

```

```

17     protected $casts = [
18         'embedding_768' => Vector::class,
19         'embedding_3072' => Vector::class,
20         'embedding_1536' => Vector::class,
21         'embedding_2048' => Vector::class,
22         'embedding_1024' => Vector::class,
23         'embedding_4096' => Vector::class
24     ];
25
26     public function document(): \Illuminate\Database\Eloquent\Relations\BelongsTo
27     {
28         return $this->belongsTo(Document::class);
29     }
30 }

```

For the trait and the casts, we follow these docs: [here](#)⁶.

And we have many casts! But now we are ready to make a class to fill up this data.

Ingestion and Embedding

So, how you get the data in here is up to you: file uploads in the UI, web scraping, or a folder, etc. What I want to cover, though, is how to **chunk** the data and then how to **embed** it.

Chunking

We will keep this simple: you will see in [Larallama.io](#)⁷, and this will lead to this [ProcessFileJob](#)⁸.

Here is the goal:

Figure 156. app/Domains/Documents/ChunkContent.php

```

1  <?php
2
3  namespace App\Domains\Documents;
4
5  class ChunkContent
6  {
7      public function handle(string $content, string $title = null): string
8      {
9          // 1. Make a document model

```

⁶<https://github.com/pgvector/pgvector-php?tab=readme-ov-file#laravel>

⁷<https://github.com/LlmLaraHub/larallama/blob/main/app/Http/Controllers/CollectionController.php#L116>

⁸<https://github.com/LlmLaraHub/larallama/blob/main/app/Jobs/ProcessFileJob.php#L111>


```
10         // 2. Break the content into chunks
11         // 3. Save the chunks to the database
12         // 4. Embed the chunks
13         // 5. 🎉 Done
14     }
15 }
```

Break the Content into Chunks

We add this to the file:

Figure 157. app/Domains/Documents/ChunkContent.php

```
1 <?php
2
3 namespace App\Domains\Documents;
4
5 use App\Domains\Chunking\TextChunker;
6 use App\Models\Document;
7
8 class ChunkContent
9 {
10     public function handle(string $content, string $title = null): string
11     {
12         // Make a document model
13         $document = new Document();
14         $document->title = $title;
15         $document->content = $content;
16         $document->save();
17
18         // Break the content into chunks
19         $chunks = (new TextChunker)->handle($content);
20         // Save the chunks to the database
21         // Embed the chunks
22         // Done
23     }
24 }
```

Making the **Document** model is easy, and breaking the text into chunks is easy as well. Let's look into that class.

Figure 158. app/Domains/Documents/Chunking/TextChunker.php

```

1  <?php
2
3  namespace App\Domains\Chunking;
4
5  class TextChunker
6  {
7      public static function handle(string $text, int $chunkSize = 600, int $overlapSize = 100): array
8      {
9          {
10             $chunks = [];
11             $textLength = strlen($text);
12
13             // Calculate where the first chunk starts and the subsequent chunks.
14             for ($start = 0; $start < $textLength; $start += ($chunkSize - $overlapSize))
15             {
16                 if ($start + $chunkSize > $textLength) {
17                     // Get the remaining text if it's shorter than the chunk size.
18                     $chunks[] = substr($text, $start);
19                     break;
20                 }
21
22                 // Get the chunk from the text.
23                 $chunks[] = substr($text, $start, $chunkSize);
24             }
25
26             return $chunks;
27         }
28     }

```

We are taking that text and considering a chunk size of 600 and an overlap of 100, and then making an array of these chunks.

We want to break the data into small enough chunks that when we embed it and search for it, we get just enough data back and not too much. For example, if I make the entire page a chunk, but the user is searching for something that is really in the middle of the page, then we are going to return the entire page, giving the LLM too much context to really answer the users's prompt (question) well.

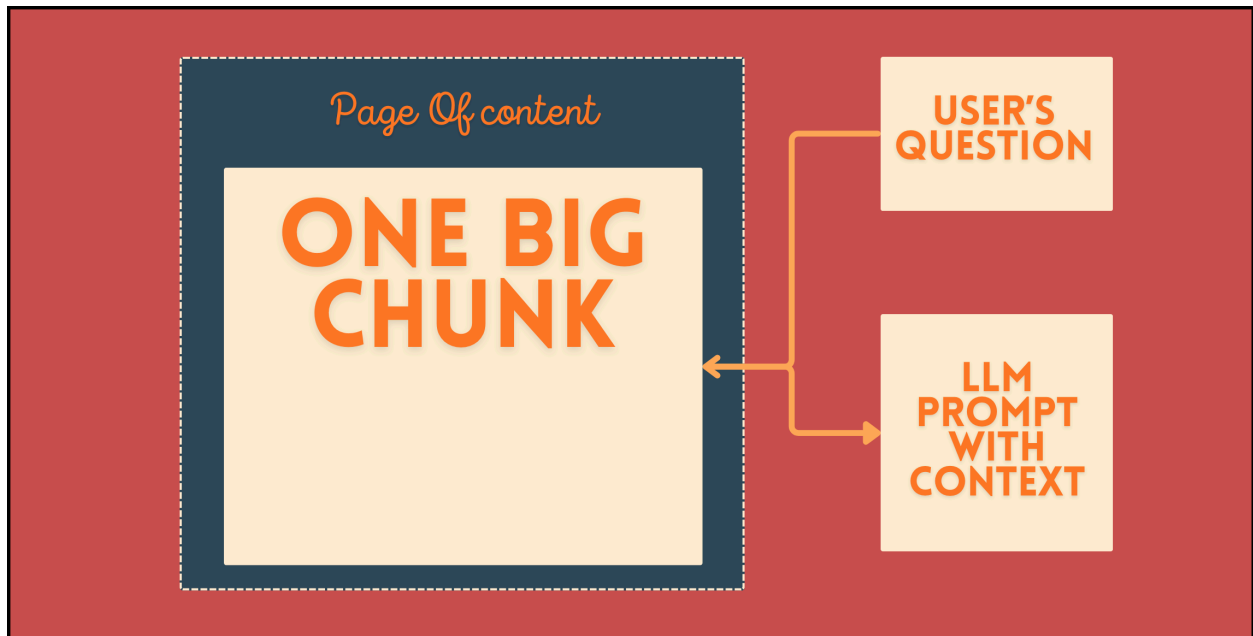


Figure 159. Too much Context

Versus smaller chunks with overlap:

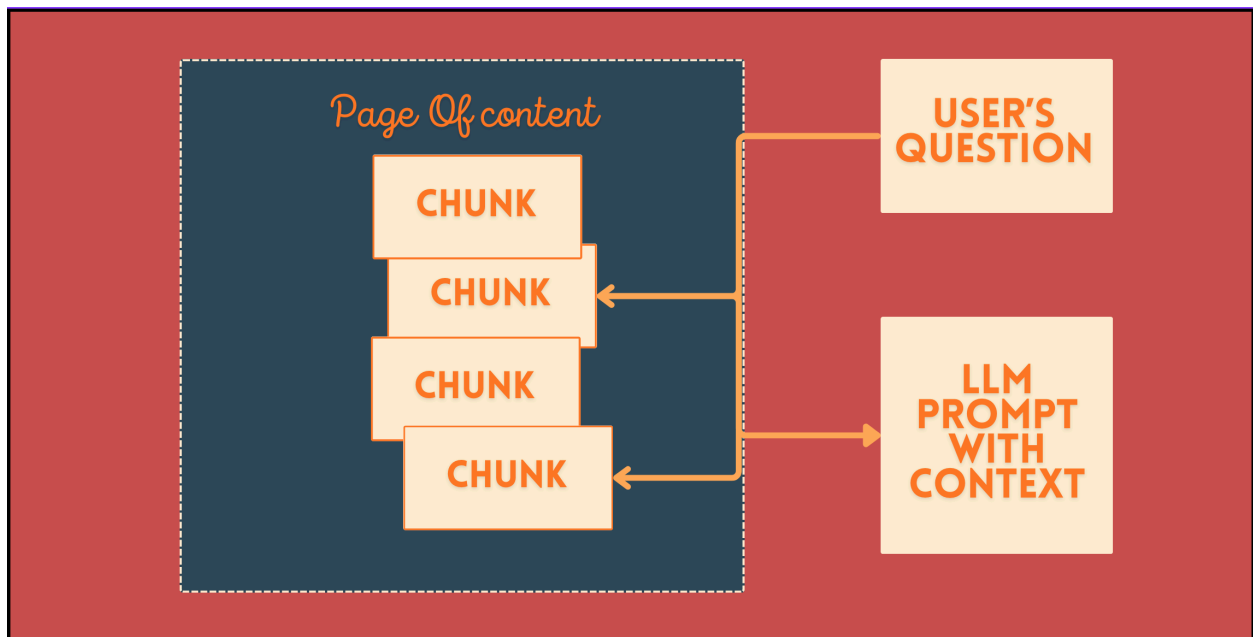


Figure 160. Just the right Chunk

We are starting to give the LLM better context so it can focus on the answer to the user's prompt. Now, let's save the chunks.

Figure 161. app/Domains/Documents/ChunkContent.php

```
1 <?php
2
3 namespace App\Domains\Documents;
4
5 use App\Domains\Chunking\TextChunker;
6 use App\Models\Chunk;
7 use App\Models\Document;
8
9 class ChunkContent
10 {
11     public function handle(string $content, string $title = null): string
12     {
13         $document = new Document();
14         $document->title = $title;
15         $document->content = $content;
16         $document->save();
17
18         $chunks = (new TextChunker)->handle($content);
19
20         foreach ($chunks as $chunkSection => $chunk) {
21             $page_number = 1;
22             Chunk::create(
23                 [
24                     'section_number' => $chunkSection,
25                     'content' => $chunk,
26                     'document_id' => $document->id,
27                     'sort_order' => $page_number,
28                 ]
29             );
30         }
31     }
32 }
```

So, we simply iterate through the chunk, saving the chunk's order and content.

And we can see in our test that there are many chunks for one file!

Figure 162. tests/Feature/ChunkContentTest.php

```
1 <?php
2
3 test('can chunk content', function () {
4     $data = get_fixture("example_markdown.txt", false);
5
6     $results = (new \App\Domains\Documents\ChunkContent)->handle($data);
7
8     expect($results->content)->not->toBeNull();
9     expect($results->chunks->count())->toBe(23);
10
11 });
```

That is one file in the docs about Artisan!

But we still need to embed the chunk - let's move on to the most important part.

Embedding

Now, we will embed the chunk of content. This can take time if you are using your machine and Ollama; it depends on your machine's speed. If you are not multi-threading the work, as with a batch system, then it will be limited to one thread at a time. For this example, I will use Ollama to start, and then OpenAI.

Embedding with Ollama

First, we need to pull down a model that enables us to embed. If we visit <https://ollama.com/search?c=embedding>, we see the embedding models. Typically, I use **mxbai-embed-large**, but for this one, I want to try **nomic-embed-text**.

We need to first **pull** down the model from Ollama:

Figure 163. ollama

```
1 ollama pull nomic-embed-text
```

We are ready to use the model in our app once we have it.

And in our code, per the ollama docs <https://github.com/ollama/ollama/blob/main/docs/api.md#generate-embeddings>, we generate embeddings using the abstracted driver system we built in the previous chapter. But first, for a moment, let us just use it directly in this code so we can see a very simple use case:

Figure 164. app/Domains/Documents/ChunkContent.php

```
1 <?php
2
3 namespace App\Domains\Documents;
4
5 use App\Domains\Documents\Chunking\TextChunker;
6 use App\Models\Chunk;
7 use App\Models\Document;
8 use Illuminate\Support\Facades\Http;
9
10 class ChunkContent
11 {
12     public function handle(string $content, string $title = null): Document
13     {
14         $document = new Document();
15         $document->title = $title;
16         $document->content = $content;
17         $document->save();
18
19         $chunks = (new TextChunker)->handle($content);
20
21         foreach ($chunks as $chunkSection => $chunk) {
22             $embedding = Http::withHeaders([
23                 'Content-Type' => 'application/json',
24                 'Authorization' => 'Bearer ollama',
25             ])->post('http://localhost:11434/api/embed', [
26                 'model' => 'nomic-embed-text',
27                 'input' => $chunk,
28             ])->json();
29
30             Chunk::create(
31                 [
32                     'content' => $chunk,
33                     'document_id' => $document->id,
34                     'sort_order' => $chunkSection,
35                     'embedding_768' => $embedding['embeddings'][0],
36                 ]
37             );
38         }
39
40         return $document;
41     }
```

42 }

We see above that we are just an API call away from embedding our data! We just send it to the API we use (we will abstract this out shortly) and then save it to the correct column.

“How did I know the correct column?” - this is a good question. I could read the docs, or try it and get an error like this:

```
SQLSTATE[ 22000]: Data exception: 7 ERROR: expected 2048 dimensions, not 768
(Connection: pgsql, SQL: insert into "chunks" ("content", "document_id", "sort_order",
"embedding_2048", "updated_at", "created_at")
```

Then I add or use that correct column.

Before we refactor, let's ask our first question.

Our First RAG Question

We are going to build this “vector search” as a Tool shortly, but for now we will use [TinkerWell](https://tinkerwell.app/)⁹ to show how this works so we can keep it simple and all in one place.

Figure 165. Tinkerwell

```
1 use App\Models\Chunk;
2 use App\Services\LlmServices\LlmDriverFacade;
3 use App\Services\LlmServices\Requests\MessageInDto;
4 use Illuminate\Support\Facades\Http;
5 use Pgvector\Laravel\Distance;
6 use Pgvector\Vector;
7
8 /**
9  * 1) The user's question
10 */
11 $question = "What are the top two tips to being a good teammate";
12
13 /**
14  * 2) We embed the question so we can use it in the query
15 */
16 $embedding = Http::withHeaders([
17     "Content-Type" => "application/json",
18     "Authorization" => "Bearer ollama"
19 ])
```

⁹<https://tinkerwell.app/>

```
20   ->post("http://localhost:11434/api/embed", [  
21     "model" => "nomic-embed-text",  
22     "input" => $question  
23   ])  
24   ->json();  
25  
26 /**  
27  * 3) Using the Vector library we instantiate the class  
28  *    with the question vectorized 🗨️  
29  */  
30 $embedding = new Vector($embedding["embeddings"][0]);  
31  
32 /**  
33  * 4) We use the Vector library to then search  
34  *    More on this shortly  
35  */  
36 $searchresults = Chunk::query()  
37   ->select("chunks.content", "chunks.sort_order", "documents.title")  
38   ->join("documents", "chunks.document_id", "=", "documents.id")  
39   ->orderBy("chunks.sort_order")  
40   ->nearestNeighbors("embedding_768", $embedding, Distance::Cosine)  
41   ->limit(3)  
42   ->get()  
43   ->map(function ($chunk) {  
44     return "File Name: " .  
45       $chunk->title .  
46       "\n" .  
47       "Sort Order: " .  
48       $chunk->sort_order .  
49       "\n" .  
50       $chunk->content;  
51   })  
52   ->implode("\n\n");  
53  
54  
55 /**  
56  * 5) Then we take this context and  
57  *    give it to the LLM with the original text version of the question  
58  *    to give a full answer to the question  
59  */  
60 $messages = [  
61   MessageInDto::from([  
62     "role" => "system",
```



```

63     "content" =>
64     "You are an assistant with our collection of data about Pickleball. Do not ans\
65 wer questions outside the content given. If you do not know the answer due to a lack
66 of information, ask the user to rephrase their question. The results will include F
67 ile Name and Sort Order. Please reference the File Name in your results, make a list
68 of Files, and sort them in order as citations in Markdown at the bottom of the resu
69 lts."
70   ]),
71   MessageInDto::from([
72     "role" => "user",
73     "content" => $question
74   ])
75 ];
76
77 $payload = [
78   "model" => "llama3",
79   "messages" => $messages,
80   "stream" => false,
81   "options" => [
82     "temperature" => 0
83   ]
84 ];
85
86 //HTTP HERE WITH OPTIONS
87 $results = Http::withHeaders([
88   "content-type" => "application/json"
89 ])
90   ->timeout(120)
91   ->baseUrl("http://localhost:11434/api/")
92   ->post("/chat", $payload);
93
94 $results->json();

```

NOTE: We set the temperature to 0 - this is key! We are saying, “Do not be creative, just stick to the facts”. You could set it higher to experiment, say 0.2, and so on. You can read more about this in the [OpenAi Docs](https://platform.openai.com/docs/api-reference/chat/create#chat-create-temperature)¹⁰

Results:

¹⁰<https://platform.openai.com/docs/api-reference/chat/create#chat-create-temperature>

Figure 166. Results from question

```

1  [
2      "model" => "llama3",
3      "created_at" => "2024-09-21T10:56:53.398165Z",
4      "message" => [
5          "role" => "assistant",
6          "content" => ""
7      ]
8      Based on our collection of data about Pickleball, here are the top two tips to bein\
9      g a good teammate:\n
10     \n
11     1. **Communicate Effectively**: Good communication is key to success in Pickleball.\
12     As a teammate, you should be able to convey your thoughts and intentions to your pa
13     rtner. This includes calling out shots, signaling where you want the ball to go, and
14     providing encouragement or constructive feedback.\n
15     2. **Support Each Other's Strengths**: A good Pickleball team complements each othe\
16     r's strengths and weaknesses. As a teammate, you should be aware of your partner's s
17     kills and play styles and work together to create opportunities for them to shine. T
18     his might mean setting up your partner for an easy shot or covering their weaknesses
19     .\n
20     \n
21     Files:\n
22     \n
23     * [Pickleball_Teamwork.pdf](#) (Sort Order: 1)\n
24     * [Effective_Communication_in_Pickleball.md](#) (Sort Order: 2)\n
25     \n
26     Citations:\n
27     [Markdown]\n
28     [Pickleball_Teamwork.pdf](#) - "The Importance of Teamwork in Pickleball"\n
29     [Effective_Communication_in_Pickleball.md](#) - "Why Communication is Key to Succes\
30     s in Pickleball"
31     "",
32 ],
33 "done_reason" => "stop",
34 "done" => true,
35 "total_duration" => 13753143833,
36 "load_duration" => 20333792,
37 "prompt_eval_count" => 109,
38 "prompt_eval_duration" => 2630338000,
39 "eval_count" => 250,
40 "eval_duration" => 11100467000,
41 ]

```

Let's review each one of these steps in detail.

Step 1: The User's Question

Not much more info is needed here, so let us look at the next step.

Step 2: Embedding the Question

This is where things get interesting. We take the user's question and embed it using the API. It's just a normal HTTP request so we can see how this all works.

Step 3: Vectorize the Embeddings

Then we take that JSON result, which becomes an array of integers in PHP, and pass it into the Vector class to use in our query step.

Step 4: Query the Database

Then we query the database with this. Let's dig in here:

Figure 167. Vector Query

```
1 $searchresults = Chunk::query()
2   ->select("chunks.content", "chunks.sort_order", "documents.title")
3   ->join("documents", "chunks.document_id", "=", "documents.id")
4   ->orderBy("chunks.sort_order")
5   ->nearestNeighbors("embedding_768", $embedding, Distance::Cosine)
6   ->limit(3)
7   ->get()
8   ->map(function ($chunk) {
9       return "File Name: " .
10          $chunk->title .
11          "\n" .
12          "Sort Order: " .
13          $chunk->sort_order .
14          "\n" .
15          $chunk->content;
16   })
17   ->implode("\n\n");
```

Distance Types

Since I am not an AI or data expert, I am going to ask Claude for help explaining Cosine and the other options:

Figure 168. Vector Distances

1	L2 (Euclidean) Distance:
2	
3	Description: Measures the straight-line distance between two points in Euclidean space.
4	
5	Use case: Best for when the magnitude of the vectors is important, and you want to find the closest vectors in absolute terms.
6	
7	When to use: Ideal for applications where the scale of features matters, such as in geographical coordinates or physical measurements.
8	
9	
10	
11	Cosine Distance:
12	
13	Description: Measures the cosine of the angle between two vectors, ignoring magnitude and focusing on direction.
14	
15	Use case: Best when you want to find vectors pointing in a similar direction, regardless of their magnitude.
16	
17	When to use: Ideal for text similarity, document classification, and recommendation systems where the relative importance of features matters more than their absolute values.
18	
19	
20	
21	
22	Inner Product Distance:
23	
24	Description: Measures the dot product of two vectors, considering both direction and magnitude.
25	
26	Use case: Useful when both the direction and magnitude of vectors are important, but you want to emphasize similarity in larger magnitudes.
27	
28	When to use: Often used in machine learning models, especially in neural networks, where the inner product represents a similarity score.
29	
30	
31	
32	
33	In your RAG (Retrieval-Augmented Generation) system lookup, you used Cosine distance.
34	This is a good choice for text-based applications because it focuses on the similarity of direction between vectors, which often represents semantic similarity in text embeddings, regardless of the length of the text or the specific magnitudes of the embedding values.
35	
36	
37	

So we will stick to Cosine for now but I was curious when the heck I would use the other ones!

Figure 169. When to use other types

```

1 There might be specific cases where you'd consider alternatives:
2
3 If the absolute magnitude of features is crucial (e.g., if the price range is a crit\
4 ical factor in product similarity), you might consider L2 distance.
5 If you want to emphasize both direction and magnitude (e.g., in a very specialized p\
6 roduct recommendation system where both similarity and popularity matter), you might
7 use Inner Product distance.
8
9 However, for general-purpose related product recommendations and document similarity\
10 in blog or documentation systems, Cosine distance is usually the most reliable and
11 widely-used choice.

```

Step 5: Use the Results for the Final Prompt

Context is key!

Now that we have some results, we need to use those results as our **context**. In this example, I take the user's question and append the context to it. And then we add a system level prompt to make sure the LLM understands the goal and limits of its role:

Figure 170. Message Array as a DataObject

```

1 $prompt = [
2     MessageInDto::from([
3         "role" => "system",
4         "content" => "You are an assistant with our collection of data about Pickleb\
5 all. Do not answer questions outside the content given. If you do not know the answe
6 r due to a lack of information ask the user to rephrase their question. Please inclu
7 de the document name and page in your response"
8     ]),
9     MessageInDto::from([
10         "role" => "assistant",
11         "content" => "Here is relevant information from our database:\n\n" . $search\
12 results
13     ]),
14     MessageInDto::from([
15         "role" => "user",
16         "content" => $question
17     ])
18 ];

```

Adding assistant array objects to the thread is fine. You might need to sort them based on the LLM you are using. (The Abstraction chapter covers this topic).

Pull in all the Files and Computer Considerations

Getting data in can be done in many ways: file upload, webhooks, API, etc. Below, I just wrote a quick command to get the files into the database.

Figure 171. app/Console/Commands/ChunkFolder.php

```

1  <?php
2
3  namespace App\Console\Commands;
4
5  use App\Domains\Documents\ChunkContent;
6  use Illuminate\Console\Command;
7  use Illuminate\Support\Facades\File;
8
9  class ChunkFolder extends Command
10 {
11     /**
12      * The name and signature of the console command.
13      *
14      * @var string
15      */
16     protected $signature = 'app:chunk-folder {absolute-path}';
17
18     /**
19      * The console command description.
20      *
21      * @var string
22      */
23     protected $description = 'Chunk all files in the folder';
24
25     /**
26      * Execute the console command.
27      */
28     public function handle()
29     {
30         $this->withProgressBar(File::allFiles($this->argument('absolute-path')), fun\
31 ction ($file) {
32             try {
33                 $fileName = $file->getFilenameWithoutExtension();
34                 $content = File::get($file);
35                 $this->info('Chunking '.$file);
36                 app(ChunkContent::class)->handle($content, $fileName);

```

```

37         } catch (\Throwable $e) {
38             $this->error('Error chunking '.$file);
39         }
40     });
41 }
42 }

```

This will take some time on my M3, so keep that in mind. I could use the OpenAI API, but this is a great example of how we can play with ideas without spending extra money. OK, that laptop is extra money, but I would have bought it anyway. These are just some thoughts on that while I wait.

Figure 172. TOPs

```

1  Tensor Processing Units (TOPs, or Tera Operations Per Second) and how this impacts y\
2  our computers ability to run a local LLM.
3
4  Here's a comparison of the TOPs (Tensor Operations Per Second) for the M1, M2, and M\
5  3 series Mac chips:
6
7  1. M1 Series:
8      - M1: Up to 11 TOPs
9      - M1 Pro: Up to 22 TOPs
10     - M1 Max: Up to 44 TOPs
11     - M1 Ultra: Up to 88 TOPs
12
13  2. M2 Series:
14     - M2: Up to 15.8 TOPs
15     - M2 Pro: Up to 31.6 TOPs
16     - M2 Max: Up to 63.2 TOPs
17     - M2 Ultra: Up to 126.4 TOPs
18
19  3. M3 Series:
20     - M3: Up to 18 TOPs
21     - M3 Pro: Up to 36 TOPs
22     - M3 Max: Up to 72 TOPs
23     - M3 Ultra: predicting 144 TOPs

```

Adding Images

Now for images! If you look in the folder `storage/data/pickleball/images`, you will see many images that are charts and data related to the documents. This can also be data we get out of

PowerPoints, PDFs, and more. All those seem to have decent tools we can use to not only strip the text out of them but also get the images and tables.

First let us refactor the `app/Console/Commands/ChunkFolder.php` class so when we come to an image we take a moment to call the OpenAI API (was not impressed with Lava 🙄).

Figure 173. `app/Console/Commands/ChunkFolder.php`

```

1  <?php
2
3  namespace App\Console\Commands;
4
5  use App\Domains\Documents\ChunkContent;
6  use App\Services\LlmServices\LlmDriverFacade;
7  use Illuminate\Console\Command;
8  use Illuminate\Support\Facades\File;
9  use Illuminate\Support\Facades\Http;
10 use Illuminate\Support\Facades\Log;
11
12 class ChunkFolder extends Command
13 {
14     /**
15      * The name and signature of the console command.
16      *
17      * @var string
18      */
19     protected $signature = 'app:chunk-folder {absolute-path}';
20
21     /**
22      * The console command description.
23      *
24      * @var string
25      */
26     protected $description = 'Chunk all files in the folder';
27
28     /**
29      * Execute the console command.
30      */
31     public function handle()
32     {
33         $this->withProgressBar(File::allFiles($this->argument('absolute-path')), fun\
34 ction ($file) {
35             try {
36                 $fileName = $file->getFilenameWithoutExtension();
37

```



```

38         $content = File::get($file);
39
40         if(str(File::mimeType($file))->contains('image')){
41             $content = $this->getImageContent($file);
42         }
43
44         $this->info('Chunking '.$file);
45         app(ChunkContent::class)->handle($content, $fileName);
46     } catch (\Throwable $e) {
47         $this->error('Error chunking '.$file);
48         Log::error($e);
49     }
50 });
51 }
52
53 protected function getImageContent($file) : string
54 {
55     $prompt = <<<PROMPT
56
57 This is an image with chart data about Pickleball. Please describe the data shown in\
58 the image and render it as a markdown table.
59 Do not add any flavor text.
60
61 PROMPT;
62
63     $results = LlmDriverFacade::driver("openai")->image(
64         prompt: $prompt,
65         base64Image: base64_encode(File::get($file))
66     );
67
68     return $results->content;
69 }
70 }

```

And now we end with more data like:

```

1 | Year | Total Players (millions) |
2 |-----|-----|
3 | 2017 | 3.1 |
4 | 2018 | 3.3 |
5 | 2019 | 3.5 |
6 | 2020 | 4.2 |
7 | 2021 | 4.8 |
8 | 2022 | 8.9 |

```

Then when we ask a question like “Total Players in 2021 compared to 2017”, our system can easily give us an answer.

```

1 According to our database, the total number of Pickleball players in 2017 was 3.1 mi\
2 llion.\n
3 \n
4 To find the total number of players in 2021, we can look at the "player_growth" file\
5 :\n
6 \n
7 | Year | Total Players (millions) |\n
8 |-----|-----|\n
9 | 2017 | 3.1 |\n
10 | 2018 | 3.3 |\n
11 | 2019 | 3.5 |\n
12 | 2020 | 4.2 |\n
13 | 2021 | 4.8 |\n
14 \n
15 According to this information, the total number of Pickleball players in 2021 was 4.\
16 8 million.\n
17 \n
18 To find the difference between the number of players in 2021 and 2017, we can subtra\
19 ct:\n
20 \n
21 4.8 (2021) - 3.1 (2017) = 1.7\n
22 \n
23 So, there were approximately 1.7 million more Pickleball players in 2021 compared to\
24 2017.

```

Pretty nice!

Adding Voice

Let’s do an audio example to show how easy this really is and how it is not far off from what we did above. In this case, I downloaded a short audio clip and then updated the `ChunkFolder` command as

below.

Figure 174. app/Console/Commands/ChunkFolder.php

```

1  <?php
2
3  namespace App\Console\Commands;
4
5  use App\Domains\Documents\ChunkContent;
6  use App\Services\LlmServices\LlmDriverFacade;
7  use Illuminate\Console\Command;
8  use Illuminate\Support\Facades\File;
9  use Illuminate\Support\Facades\Log;
10
11 class ChunkFolder extends Command
12 {
13     /**
14      * The name and signature of the console command.
15      *
16      * @var string
17      */
18     protected $signature = 'app:chunk-folder {absolute-path}';
19
20     /**
21      * The console command description.
22      *
23      * @var string
24      */
25     protected $description = 'Chunk all files in the folder';
26
27     /**
28      * Execute the console command.
29      */
30     public function handle()
31     {
32         $this->withProgressBar(File::allFiles($this->argument('absolute-path')), fun\
33 ction ($file) {
34             try {
35                 $fileName = $file->getFilenameWithoutExtension();
36
37                 if (str(File::mimeType($file))->contains('image')) {
38                     $content = $this->getImageContent($file);
39                 } elseif (str(File::mimeType($file))->contains('audio')) {
40                     $content = $this->getAudioContent($file);

```

```

41         } else {
42             $content = File::get($file);
43         }
44
45         $this->info('Chunking '.$file);
46         app(ChunkContent::class)->handle($content, $fileName);
47     } catch (\Throwable $e) {
48         $this->error('Error chunking '.$file);
49         Log::error($e);
50     }
51 });
52 }
53
54 protected function getAudioContent($file): string {
55
56     $results = LlmDriverFacade::driver("openai")->audio(filePath: $file);
57
58     return $results->content;
59 }
60
61
62 protected function getImageContent($file): string
63 {
64     $prompt = <<<'PROMPT'
65
66 This is an image with chart data about Pickleball. Please describe the data shown im\
67 age and render as a markdown table.
68 Do not add any flavor text.
69
70 PROMPT;
71
72     $results = LlmDriverFacade::driver('openai')->image(
73         prompt: $prompt,
74         base64Image: base64_encode(File::get($file))
75     );
76
77     return $results->content;
78 }
79 }

```

And added to the driver (see the chapter on this):

Figure 175. app/Services/LlmServices/OpenAiClient.php

```

1     public function audio(
2         string $filePath,
3         string $model = "whisper-1"): CompletionResponse
4     {
5         $token = $this->getConfig('openai')['api_key'];
6
7         if (is_null($token)) {
8             throw new \Exception('Missing open ai api key');
9         }
10
11        $file = new UploadedFile($filePath, basename($filePath), null, null, true\
12e);
13
14        $response = Http::withToken($token)
15            ->baseUrl('https://api.openai.com/v1')
16            ->timeout(240)
17            ->attach('file', $file->getContent(), $file->getClientOriginalName())
18            ->post('/audio/transcriptions', ['model' => $model]);
19
20        if ($response->failed()) {
21            Log::error('OpenAi API Error ', [
22                'error' => $response->body(),
23            ]);
24
25            throw new \Exception('OpenAi API Error Chat');
26        }
27
28        return CompletionResponse::from([
29            'content' => data_get($response->json(), 'text', ''),
30            'stop_reason' => 'stop',
31        ]);
32    }

```

I used OpenAI for this, making it so easy to send the file and get results.

NOTE: I use HTTP here to show how easy it is to talk to these services. The [PHP library for OpenAI](#)¹¹ is great, but I want to enable you to talk to any API and see how it works here.

Now, when we pull in the files, it will take care of this one, adding a chunk that has nothing to do with pickleball!

¹¹<https://github.com/openai-php/laravel>

Figure 176. Test Audio File Output

```
1 This is a test audio file to prove that this is working.
```

And that's it. Our RAG system now has audio and data from images.

Unstructured vs Structured

...almost like the use cases between NoSQL and SQL, I have to consider what can the LLM do for me and what do I have to do for it.

As this data comes in, we are treating it as unstructured. You can read about a [very mature library](#)¹² that helps to add some structure to this data. Sure, we did pull out some data from images in the example above, but we did not consider the importance of simple Markdown or HTML elements like H1, H2, H3, Bold and so on. We also did not mark the chunks as type of content for example “H1”, “Image”, “Table” and so forth. If you are working on a project where this level of detail is important, then putting together the PHP code to do it is of course possible.

Let us look at some existing code on this to see how it could be done. If you visit <https://github.com/LlmLaraHub/larallama/blob/main/app/Jobs/ProcessFileJob.php#L17> you start to see that this job considers the different types of files. What this will lead to for a PPTX file is this: <https://github.com/LlmLaraHub/larallama/blob/main/app/Domains/Documents/Transformers/ProcessPpt.php>; let's take a moment to see the data type:

Figure 177. app/Domains/Documents/Transformers/ProcessPpt.php

```
1 if ($shape instanceof RichText) {
2     $pageContent = $shape->getPlainText();
3     $guid = $shape->getHashCode();
4
5     $content = $this->output(
6         type: StructuredTypeEnum::Narrative,
7         content: $pageContent,
8         page_number: $page_number,
9         guid: $guid,
10        element_depth: $shapeCount,
11        is_continuation: $shapeCount > 0,
12    );
13
14    yield $content;
15 } elseif ($shape instanceof Table) {
16     $table = $shape->getRows();
```

¹²<https://unstructured.io/>

```

17     $this->title = 'Table';
18     $this->subject = 'Table';
19
20     $content = $this->output(
21         type: StructuredTypeEnum::Table,
22         content: 'table data',
23         page_number: $page_number,
24         guid: $shape->getHashCode(),
25         element_depth: 0,
26         is_continuation: false,
27     );
28     yield $content;

```

We could take time to do a few different integrations based on the type of data. If a table, we could use the LLM to turn it into a Markdown table. Or if it is an image we could do as we did before to pull the information out of it. Then at the same time we could add a field to the chunks table called type to make it clear what type of data this is. As we see earlier in the search method, we then can surface this data to help the LLM consider it in its results:

Figure 178. Example Search

```

1  $searchresults = Chunk::query()
2      ->select("chunks.content", "chunks.sort_order", "chunks.type", "documents.title")
3      ->join("documents", "chunks.document_id", "=", "documents.id")
4      ->orderBy("chunks.sort_order")
5      ->nearestNeighbors("embedding_768", $embedding, Distance::Cosine)
6      ->limit(3)
7      ->get()
8      ->map(function ($chunk) {
9          return "File Name: " .
10             $chunk->title .
11             "\n" .
12             "Sort Order: " .
13             $chunk->sort_order .
14             "\n" .
15             "Type: " .
16             $chunk->type .
17             "\n" .
18             $chunk->content;
19      })
20      ->implode("\n\n");

```

By tagging different parts of your content (like headers, tables, or images), you're essentially giving your AI a roadmap to navigate the information.

Let the LLM do the work for us as we then fix our prompt to ask it to consider types and what to do with that information. Maybe something like this:

Figure 179. Update Prompt to consider Types

```
1 You are an assistant with our collection of data about Pickleball. Do not answer que\
2 stions outside the content given. If you do not know the answer due to a lack of inf
3 ormation ask the user to rephrase their question. The results will include File Name
4 and Sort Order. Please reference the File Name in your results and make a list of F
5 iles and sort order as citations in Markdown at the bottom of the results. There is
6 a Type attribute as well to help you give weight to the type of data for example, Ta
7 ble, H1, H2, H3, and so on.
```

While chunking files and throwing them at your LLM can work, taking the time to add a bit of structure to your data can seriously level up your game. It's not just about making your data look pretty – it's about giving your LLM the context it needs to return more concise results. By tagging different parts of your content (like headers, tables, or images), you're essentially giving your AI a roadmap to navigate the information. This means more accurate responses, better handling of complex documents, and the ability to squeeze more value out of those precious tokens. It is a funny balance, though, since as a developer, I really really want to structure my data, but almost like the use cases between NoSQL and SQL, I have to consider what the LLM can do for me and what I have to do for it.

Beyond RAG

I want to end this chapter by discussing going beyond RAG. First, there are so many use cases for vector searching that it would be a waste to associate it only with RAG. For example, you can share with a user “similar searches,” or you can use it to find “similar content,” “recommendation system,” “anomaly detection in time series data,” and more.

Lastly, people and businesses will want to do more than just chat with their data. Though RAG makes a great foundation for so many possibilities for users and their data, any RAG system can be a hub for numerous other use cases, or what I call a Content Retrieval System or CRS. Creating a system that can vectorize incoming data is excellent and makes a foundation for many other possible use cases. On that note, I was using [NotbookLM](https://notebooklm.google/)¹³ by Google, which is essentially a RAG system. But then they took it so much further, “generate a podcast”, “create a study guide”, and more. You can see how easy it is to set up this foundation for any data store you or your client might have.

What You Learned

The Basics of RAG Systems: Retrieval-Augmented Generation combines the power of large language models (LLMs) with the precision of information retrieval from external databases, helping

¹³<https://notebooklm.google/>

to overcome the limitations of LLMs by grounding their responses in accurate and up-to-date information.

The Importance of Vector Databases: Vector databases are crucial for RAG systems because they enable semantic search, allowing users to find related information even when the search terms don't exactly match the document content. Traditional keyword-based search engines often fall short in this area.

Chunking and Embedding: These are essential steps for preparing your data for a RAG system. Chunking divides large documents into smaller, more manageable units, ensuring that the LLM receives the most relevant context for a given query. Embedding transforms text into numerical vectors that capture the meaning and relationships between words, making it possible to measure semantic similarity.

Using LLMs for Embedding and Response Generation: LLMs are used in two key ways within a RAG system. First, they generate the embeddings used to represent text in the vector database. Second, they process retrieved information and user queries to generate coherent and informative responses.

Going Beyond Simple Text: The source shows how to incorporate images and audio data into a RAG system. Techniques such as extracting text from images using optical character recognition (OCR) and transcribing audio using speech-to-text models extend the capabilities of the RAG system beyond traditional text-based information retrieval.

The Power of Structured Data: Structuring your data improves the RAG system's accuracy and efficiency by adding semantic tags or metadata, guiding the LLM to better understand the content and provide more relevant results.