

# The LLM Mesh

An Architecture for Building Agentic Applications in the Enterprise



Kurt Muehmel



## **Efficient and Governed** Generative AI With Dataiku



#### THE LLM MESH

A common backbone for GenAI applications, enabling choice and flexibility among the growing number of models and providers.



#### On DATAIKU ANSWERS

A packaged, scalable web application to democratize enterprise-ready LLM chat and retrieval-augmented generation (RAG).



#### PROMPT STUDIOS

Iteratively design and evaluate LLM prompts, compare performance and cost across models, and operationalize GenAl in your data projects.



#### AI-POWERED ASSISTANTS

Go faster and farther with AI Prepare, AI Code Assistant, and AI Explain, all of which improve efficiency and the overall product.



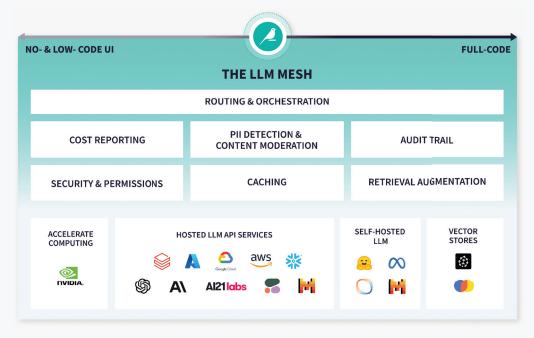
#### LLM-POWERED DATA

No-code text recipes enhanced with pre-trained Hugging Face models and LLMs for text summarization, classification, and other common language tasks.



#### **GENAL SOLUTIONS**

Pre-built Generative AI use cases and applications for even faster time to value.



### The LLM Mesh

# An Architecture for Building Agentic Applications in the Enterprise

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

Kurt Muehmel



#### The LLM Mesh

by Kurt Muehmel

Copyright © 2025 O'Reilly Media, Inc. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<a href="http://oreilly.com">http://oreilly.com</a>). For more information, contact our corporate/institutional sales department: 800-998-9938 or <a href="mailto:corporate@oreilly.com">corporate@oreilly.com</a>.

Editors: Jeff Bleiel and Aaron Black
Production Editor: Kristen Brown
Interior Designer: David Futato

Cover Designer: Susan Brown
Illustrator: Kate Dullea

August 2025: First Edition

#### Revision History for the Early Release

2024-08-02: First Release 2024-10-01: Second Release 2024-10-22: Third Release 2025-02-28: Fourth Release

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *The LLM Mesh*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the author and do not represent the publisher's views. While the publisher and the author have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the author disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

This work is part of a collaboration between O'Reilly and Dataiku. See our statement of editorial independence.

## **Table of Contents**

Bri	ef Table of Contents ( <i>Not Yet Final</i> )	Vii
1.	Using LLMs in the Enterprise	. 9
	What Is an LLM Mesh?	11
	The Right Model for the Right Application	14
	Bottom Line: Why the LLM Mesh?	27
2.	Objects for Building Agentic Applications	29
	The Potential of New Agentic Applications	30
	LLM Mesh-Related Objects: An Overview	35
	The Objects of an LLM Mesh in Detail	38
	Cataloging LLM-Related Objects	53
	Conclusion	54
3.	Quantifying and Optimizing the Cost of LLMs in the Enterprise	57
	Quantifying the Costs of Agentic Applications	58
	Techniques for Limiting Costs	71
	Cost-Efficient AI Operations in the Enterprise	77
	Conclusion	80
4.	Measuring and Monitoring the Performance of Agentic Application	s.
		83
	How an LLM Mesh Helps Measure Performance	85
	Measuring and Monitoring the Quality of Generated Text	86
	Intrinsic Quality Evaluation	88

Extrinsic Quality Evaluation	90
Implementing a Performance Architecture in an LLM Mesh	100
Measuring and Monitoring the Speed of Agentic	
Applications	104
Conclusion	105

vi Table of Contents

# Brief Table of Contents (*Not Yet Final*)

Chapter 1: Using LLMs in the Enterprise (available)

Chapter 2: Objects for Building LLM-Powered Applications (available)

Chapter 3: Quantifying and Optimizing the Cost of LLMs in the Enterprise (available)

Chapter 4: Measuring and Monitoring the Performance of Agentic Applications (available)

Chapter 5: Audit Trail, Security, and Permissions (unavailable)

Chapter 6: Retrieval Augmentation (unavailable)

Chapter 7: Conclusion (unavailable)

### **Using LLMs in the Enterprise**

#### A Note for Early Release Readers

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 1st chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at <code>jbleiel@oreilly.com</code>.

"May you live in times of rapid technological progress." This is the blessing and the curse of our current moment. Recent advances in AI and a growing interest in technology, thanks to the release of wildly popular consumer products, have led to a frenzy of interest in, and use of, AI, and Large Language Models (LLMs) in particular, in the enterprise.

However, AI and LLMs remain nascent in the enterprise, meaning that best practices for their use are being defined. At the same time, the core technologies — the models themselves, technologies to host and serve the models, etc. — are evolving rapidly.

Table 1-1 provides a brief timeline of the release of various models and technologies that could be relevant for enterprise use. The

diversity and speed of release create both opportunities and challenges when you are looking to use these technologies in production use cases.

Table 1-1. A (Non-Exhaustive) Timeline of Enterprise-Relevant Model and Product Releases

Developer or Provider	Model or Product	Release Date	Description	
OpenAl	GPT-3	May 2020	175 billion parameter LLM with 2048 token context window	
OpenAl	ChatGPT	November 2022	2 Consumer chatbot application, powered by GPT-3.5 Turbo	
Microsoft Azure	OpenAl Service	January 2023	Managed service offering LLMs from OpenAl	
Amazon Web Services	Bedrock	September 2023	Managed service offering LLMs from various developers	
Dataiku	LLM Mesh	September 2023	Commercial LLM Mesh offering for connecting to LLMs and building agentic applications in the enterprise	
Databricks	DBRX	March 2024	Open-weights mixture of experts model with 132B total parameters and 32k-token input context window, licensed for commercial use	
Meta	LLaMA 3 (8B, 70B)	April 2024	Updated LLM with 4096-token input context window, with updated license allowing certain commercial uses	
Mistral	Mixtral 8x22B	April 2024	Open-weights mixture of experts model with up to 141B parameters and 64k-input context window, licensed for commercial use	
OpenAl	GPT-40	May 2024	Multimodal LLM supporting voice-to-voice generation and 128k-token input context window	
OpenAl	01	September 2024	Reasoning model with built-in chain-of-thought for complex scientific and mathematical problems	
DeepSeek	R1	January 2025	Open-source reasoning model (MIT license) optimized for math, coding, and logic	

Today, you can build entirely new capabilities that would not have been possible previously, to improve the lives of your employees and better serve your customers. But you also have to keep up with rapid changes in the core technologies and use techniques that have not been fully proven. We are all now at the cutting edge. This diversity of options among the technologies and techniques is truly a great thing. In fact, we are just scratching the surface for the potential uses of LLMs in the enterprise. It's easy to 'imagine a future where these technologies are generating massive amounts of value for the enterprise, automating mundane tasks, and making new products and services possible.

In this chapter, we will briefly introduce what an LLM Mesh is, and then take an in-depth look at the many different types of LLMs that can be appropriate for use in the enterprise. We'll discuss different characteristics of models, and how models are built, published, run, and perform.

After reading this chapter, you should be able to think about how you would want to use different models for different applications in your business. Given this multitude of models, you will see why an LLM Mesh architecture is going to be a key part of your AI strategy going forward.

#### What Is an LLM Mesh?

An LLM Mesh is an architecture paradigm for building agentic applications in the enterprise. There are three principles regarding what an LLM Mesh should accomplish. An LLM Mesh should enable you to:

- 1. Access various LLM-related services through an abstraction layer.
- 2. Provide federated services for control and analysis.
- **3.** Provide central discovery and documentation for LLM-related objects.

These principles allow for agentic applications to be built in a modular manner, simplifying their development and maintenance.

Figure 1-1 illustrates an LLM Mesh architecture being used to develop two applications. Various objects, referenced in the Catalog and accessed via the Gateway, are combined to build the logic of

What Is an LLM Mesh?

<sup>1</sup> An agentic application is a software that uses an AI agent to perform tasks, make decisions, or interact with users with a defined level of autonomy. AI agents are discussed in more depth in Chapter 2.

the applications. Federated services provide control and analysis throughout the lifecycle of the application.

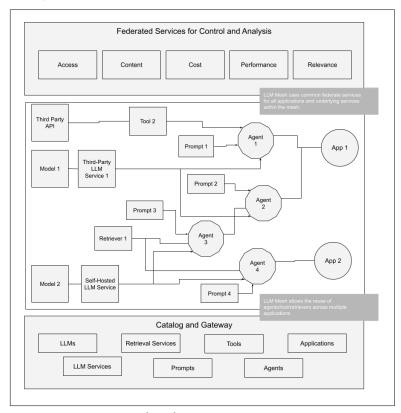


Figure 1-1. An LLM Mesh architecture

It is certainly possible to build agentic applications in the enterprise without an LLM Mesh. Many of the initial applications that organizations have built since the release of ChatGPT do not use an LLM Mesh. In these cases, the logic for connecting the various objects of the application (LLM services, retrieval services, etc.) is built directly into the application, as are any additional capabilities such as access controls or logging. This approach is perfectly appropriate for building a first proof of concept, or a single application.

An LLM Mesh, however, becomes valuable when:

1. The total number of agentic applications being developed begins to grow,

- 2. More teams start building and using the applications,
- 3. More complex agentic applications are being designed and built.

In this context, the LLM Mesh will accelerate the development of the applications, simplify their maintenance, and help to ensure that the applications meet enterprise standards for safety, security, and performance.

#### Why LLMs and Not Generative AI?

An LLM Mesh architecture focuses on LLMs and not Generative AI more broadly because LLMs are the core building blocks of the AI applications that will be built in the enterprise.

LLMs are large neural networks trained on text data. They possess a variety of natural language processing capabilities. Many, but not all, LLMs can generate text. Generative AI is a broader category of AI that includes models that can generate text, audio, images, and videos.

Beyond simply generating text, LLMs are also used to reason through a problem, to give instructions to various tools, and to write the code to connect to various tools. While image-generating models, for example, can be useful in the enterprise, they are not relevant in the context of building sophisticated AI applications that are the focus of the LLM Mesh.

An LLM Mesh provides a gateway not only to LLMs, but also to the full range of objects that are needed to build fully-featured, agentic applications. These include the LLMs themselves and the services to host them, but also agents, tools, retriever services, and applications such as chatbots.

These objects are, for most organizations, new types of assets that will need to be developed and used. The skills to develop and use these kinds of objects are not yet commonplace in organizations, and best practices for their development and use are still being defined. Amid this rapid innovation, the LLM Mesh architecture paradigm aims to simplify the management and use of these objects to accelerate and standardize the development of agentic applications. Chapter 2 will explore in depth these different types of objects and how an LLM Mesh can simplify their use.

What Is an LLM Mesh?

#### The Right Model for the Right Application

The challenge for the use of LLMs in the enterprise is not a lack of availability of models. As of February 2025, the popular model repository Hugging Face lists 1,377,986<sup>2</sup> models, of which 178,489<sup>3</sup> are text-generation models. More models are being developed and released every day.

In fact, the abundance can actually be a hindrance, as you have to sort through the many different options to choose the ones that are best for your applications.

A large general model that can do most things pretty well is a good place to start. But as an enterprise's use of LLMs matures and it seeks higher levels of performance and optimized budgets, it will need to use a growing number of models across different applications.

The following sections explore the different characteristics of models and how these characteristics may make a model more or less appropriate for the many different, specific uses in the enterprise.

# Model Size: The Upside and Downside of More Parameters

The word "large" in large language model refers to the number of parameters in the model. Alternatively, "large" may refer to the number of tokens in the training data that the model is trained on. More training tokens lead to more parameters.

LLMs often have hundreds of billions to trillions of parameters. For example, GPT-3, released in May 2020<sup>4</sup> and the immediate precursor to the model behind the first version of ChatGPT, has 175 billion parameters. The first version of LLaMA from Meta AI in February 2023 has 65 billion parameters.<sup>5</sup> Increasingly, the makers of proprietary models are no longer making the number of parameters in their models public.

<sup>2</sup> https://huggingface.co/models, accessed February 4, 2025

<sup>3</sup> https://huggingface.co/models?pipeline\_tag=text-generation&sort=trending, accessed June 20, 2024

<sup>4</sup> https://arxiv.org/abs/2005.14165

<sup>5</sup> https://ai.meta.com/blog/large-language-model-llama-meta-ai/

These parameters are the numerical values (sometimes they will be called weights and biases) that make up the simple mathematical formulae of each neuron in the neural network. Usually, they are 32-bit floating point numbers. A process called quantization can simplify these numbers to 4- or 8-bit integers. This process can often dramatically improve the efficiency of a model while having only a modest impact on model performance.

Figure 1-2 illustrates a simple neural network architecture, showing the input layer, two hidden layers, and the output layer. The circles represent the nodes in the network, the values under the nodes are the biases, while the values on the lines connecting the nodes represent the weights. Larger neural networks, like LLMs, are built on the same basic architecture but are billions of times larger with more than one hundred hidden layers.

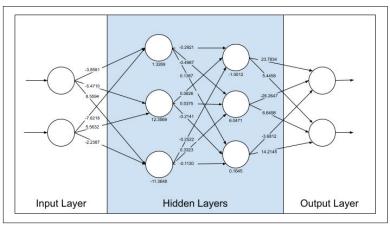


Figure 1-2. Simplified example of a neural network showing the input, hidden, and output layers and the weights connecting each node and the biases of each node

Generally speaking, larger models perform better: They can do more tasks and they can do those tasks better. Thus, it could be easy to conclude that you should choose the largest model your budget allows and use it for everything. But that would be like using your large, comfortable, powerful grand touring car for every trip. While it would be the right choice for a cross-country roadtrip, it would be overkill for a quick trip to the grocery store or the bakery around the corner. A bicycle or your own two feet would be better for such errands.

The following subsections explain the tradeoffs related to the size of a model.

#### Inference Costs

The most direct impact of a larger model size will be on inference cost. Inference is the process of generating tokens in response to a particular input. A model with more parameters will require more calculations during inference. One way or another, those calculations must be run on some hardware that is installed and managed somewhere and that is consuming electricity, for which someone will have to pay the bill at the end of the month.

In some cases, companies offering these models as a service may obfuscate these costs, for example, by subsidizing the cost in order to gain more customers. This may make an apples-to-apples comparison difficult. We'll dig into cost considerations in Chapter 3.

Some models function, in essence, as a combination of smaller models, each specialized in different tasks. This architecture, known as Mixture of Experts (MoE), can dramatically reduce the cost of inference. One well-known model using an MoE architecture is 8x7B (also known as Mixtral) from Mistral. Despite being a 46.7-billion parameter model, only 12.9-billion parameters are used per token. This approach has led to improvements in inference cost, but makes the model more challenging to build and to fine-tune.

All other things being equal, larger models will be more expensive to use, though technological advancements like MoE mean that these tradeoffs will become more complex in the future. The benefit that large models bring to a particular use case may justify their expense in certain cases, as discussed in the sections and chapters below, but a wise strategy will use them only where needed.

#### Inference Speed

While the inference of a larger model will require more calculations, these calculations can be done more quickly when using larger and higher-performance hardware. Furthermore, many of these calculations can be parallelized, using multiple processing units at the same time to run all of the necessary calculations. Again, MoE models do not need to use all parameters for every task.

Standard benchmarks are being established to accurately quantify and compare the speed of different models, acknowledging that hardware and network performance will have a significant impact on the results. The two metrics that are used most commonly are latency and throughput:

- Latency, often measured as Time To First Token (TTFT), is a measure of how long the model takes to generate its first response token to a user's input. In applications where the end user is interacting with the model in real time, latency will influence whether the model "feels" responsive. In applications where the model's response is part of a longer chain of interactions, latency will need to be considered when setting when the application will time out.
- Throughput, often measured as Tokens Per Second (TPS), measures the overall rate at which the model will generate tokens in response to a given request. Like latency, it will influence if a model feels fast to an end user. Throughput needs to be taken into consideration when building applications that depend on the output of the model.

When comparing the speed of models, pay close attention to the units being used, as different testers are using different methodologies.

While the relationship between model size and inference speed is indirect (because large models can be run more quickly on higher-performance hardware, and factors like network performance can influence the time it takes to receive a response), the measured speed of a deployed model must be taken into account when building an agentic application.

#### Task Coverage and Performance

One of the main functional differences between LLMs and previous generations of models used for Natural Language Processing (NLP) is that those earlier models were always task specific. For example, separate models would be used for sentiment analysis, text summarization, or language translation.

LLMs can do all of those tasks, and generative LLMs can do something that previous models could not: Generate text based on a prompt. Generally speaking, models with more parameters can perform more tasks, which can be useful in the enterprise when several tasks need to be performed on the input text.

LLMs have also been shown to gain new abilities as they grow larger and are trained on more and more data. These emergent abilities are unpredictable. In other words, researchers cannot predict ahead of time at what point in its training a model will gain new abilities. It is possible that future LLMs will be capable of many more tasks or will see dramatic improvement in their performance of existing tasks as they grow larger.

In addition to gaining new abilities as they grow, larger LLMs generally show better performance on any task that they are capable of performing as well. Recent research shows that this improvement in performance is not linear nor predictable, as with the emergent abilities mentioned above.<sup>7</sup>

#### **Context Windows**

The amount of input text that a model can receive within a single prompt is known as its context window. Measured in tokens, it defines how much information a model can work with at a single time.

For example, a model with a small context window can only be used to summarize a document that can fit in its context window. You could break up the document into smaller pieces, but the model would summarize each separately, without knowledge of the entire document, potentially resulting in repetitive or incoherent results. Large context windows, on the other hand, allow for plenty of space to provide examples of what you want the LLM to produce (called few-shot learning) and to engage in more complex prompt engineering techniques.

Generally speaking, larger models have larger context windows, and some models have been optimized for exceptionally large context windows. While the original GPT model had a context window of only 512 tokens (approximately one page of text) Gemini 1.5 from Google now has a context window of more than 1 million tokens and has been shown in internal testing to handle up to 10 million tokens.

<sup>6</sup> https://openreview.net/pdf?id=yzkSU5zdwD

<sup>7</sup> https://arxiv.org/abs/2210.14891

#### Sizing Models to the Task

Reading these previous sections, it is easy to conclude that if cost and complexity are no barrier, then the largest models are always the best choice for any application in the enterprise. But, in which enterprise are cost and complexity not a barrier? In fact, these are the two greatest barriers to the practical use of LLMs in the enterprise!

Given this reality, enterprise users of LLMs will need to choose a model that strikes the right balance of ability, performance, cost, and complexity for a specific application. The right choice for one application may not be the right choice for another application.

#### **General Models vs. Specialized Models**

Building on this understanding of the implications of model size, we will now explore the differences between general models and specialized models.

General models are those that have been trained to perform at human level across a wide range of tasks. OpenAI's GPT-4 (released in March 2023<sup>8</sup>) is an excellent example of such a model. It demonstrates very high performance across a great number of tasks, covering natural languages, programming languages, and a wide variety of specialized jargon. It can generate, summarize, and translate text, it can write technical reports, and it can write poetry. Furthermore, GPT-4 can take image data as input, a capability known as multimodality.

In contrast to these general models, specialized models have been trained to perform well on specific tasks, in specific domains, or have been compressed and optimized for performance at a smaller model size.

Note that while high-performing, general models tend to be larger models, specialized models may be larger or smaller.

#### Types of Specialized Models

Task-specific models are those that are focused on doing specific tasks very well. Some examples of task-specific models include M2M100°, a model that is designed to translate between any pair

<sup>8</sup> https://openai.com/index/gpt-4-research/

of natural languages, or OpenAI's Codex<sup>10</sup>, an evolution of GPT-3 that is trained specifically to generate code across a wide variety of programming languages. A common application might be a model that is specialized in summarization, allowing it to be much smaller than a general model. Thanks to its small size, it could run locally on a mobile device and be used for rapidly summarizing content directly on the phone.

Domain-specific models are those that are trained on the language of a specific domain. For example, BioMedLM<sup>11</sup> is a 2.7-billion parameter model trained on biomedical literature and is thus well-adapted to answering questions about medical topics, while BloombergGPT<sup>12</sup> is a 50-billion parameter model trained on a very large dataset of financial documents designed to serve the financial services industry.

Resource-constrained models are models that have been compressed through various techniques to maintain good performance in their desired tasks or across a wide range of tasks, while being less resource intensive to run. An example is MobileBERT<sup>13</sup>, a compressed version of the popular BERT model designed to be run on mobile devices.

Embedding models transform text into numerical representations called embeddings or vectors. These embeddings capture the semantic meanings of the text and the relationships between the different parts of the text. A common application is retrieval augmented generation (RAG) where a corpus of text (e.g., thousands of documents) are converted into embeddings and stored in a specialized database called a vector store.

Reranking models are used to refine the initial ranking of search results of an embedding model to make the results more relevant to the end user. There are LLM-based and non-LLM rerankers, each presenting tradeoffs in terms of performance and quality of response.

<sup>9</sup> https://about.fb.com/news/2020/10/first-multilingual-machine-translation-model/

<sup>10</sup> https://openai.com/index/openai-codex/

<sup>11</sup> https://arxiv.org/abs/2403.18421

<sup>12</sup> https://arxiv.org/abs/2303.17564

<sup>13</sup> https://arxiv.org/abs/2004.02984

#### Choosing a General or Specialized Model

The existence of a diverse and growing ecosystem of both general and specialized models gives enterprises the opportunity to use different models for different purposes.

In the enterprise, general models are well-suited to tasks where the input is going to be highly unpredictable. This could be the classification of documents into different categories. For example, if a directory contained a mix of contracts, invoices, and emails, a first step in the analysis could be to use a general model to sort the documents into different categories so that the contracts could be analyzed separately from the invoices.

Specialized models are well adapted for tasks where the input data is more homogeneous and predictable. Let's explore what this might look like across a pharmaceutical company. That company may wish to build a chatbot to serve its customers (doctors, nurses, pharmacists, and other healthcare providers) in their interactions with patients. It would likely choose a domain-specific model like BioMedLM to ensure higher quality and more relevant results. The same company may then use a model like ESM<sup>14</sup> from Meta AI researchers which has been trained on the language of proteins as part of their molecular research applications. Finally, that same organization may use a non-LLM computer vision model to watch their products as they come off of the manufacturing line to quickly identify any anomalies as part of their quality assurance processes.

General models can be a very good starting point for enterprises as they experiment and build their first use cases using LLMs. At those early stages, the simplicity of using a single model for a variety of tasks and use cases outweighs the benefits of further optimization using specialized models. But, as an enterprise scales its use of LLMs across use cases, enterprises will want to optimize their use to improve performance and reduce costs. In this context, specialized models become more relevant, and the number of models that an organization will need to manage and apply will tend to increase.

<sup>14</sup> https://github.com/facebookresearch/esm

#### What Is Fine-Tuning?

A common way of creating a domain-specific model is to fine-tune an existing base model. For example, this is how BloombergGPT was built, by fine-tuning the open-access BLOOM model on a proprietary dataset of financial documents.<sup>15</sup>

Fine-tuning is a type of transfer learning that feeds new — usually specialized — data into a model to retrain some parts of the model on this new data. Compared to building a model from scratch, it is far less complex and compute-intensive.

While fine-tuning is simpler and less expensive than building a base model, it remains an advanced technique and should be used only when other, simpler, and less expensive avenues have been exhausted.

Fine-tuning has often been cited as a way to elicit better performance from base models, allowing enterprises to differentiate their use of LLMs from their competition. While this is true, it ignores or downplays the difficulties of fine-tuning and leaves unexplored the opportunity to generate differentiated results using simpler techniques like prompt engineering and Retrieval Augmented Generation (RAG).

#### Making Sense of Model Licenses

There has often been a conflation between a model's license (e.g., open source vs. proprietary) and where the model is hosted (e.g., provided as a service via API vs. self-hosted or on-premises). It is important to distinguish between the two dimensions. For example, hosted services like Amazon Bedrock serve both proprietary and open models, while providers like Cohere license their proprietary models for self-hosting in addition to hosting the model themselves. Hosting options will be covered in the next section, while this section will distinguish between the different license types.

#### Proprietary Models

Proprietary models are just that: proprietary to their creators. The creator of a proprietary model retains full control over the intellec-

<sup>15</sup> https://arxiv.org/abs/2303.17564

tual property of the model itself. Most often, these models are a black box. In other words, their training data, the algorithm used to train the model, any subsequent steps such as reinforcement learning or fine-tuning, and the weights of the model itself remain hidden from the end user, unless the developer chooses to disclose any of this information.

Early in the development of LLMs, there was a trend towards openness, even among developers of proprietary models. OpenAI published a technical paper detailing the development process of GPT-3.<sup>16</sup> The release of subsequent models, such as GPT-4, have not been accompanied by such detail.

The use of proprietary models is governed by the terms of use that a customer agrees to when using the model. An enterprise should ensure a full and detailed legal review of these terms to ensure that they are appropriate for the intended use. Specific attention should be given to any rights that the model provider may claim to have on any data sent to the model for inference. Generally, models that are licensed for professional use do not retain any customer data for retraining purposes, though they may retain customer data for quality assurance purposes.

#### Open-Weights Models

An open-weights model provides public access to the pre-trained parameters of the model. This can allow the end user to modify the weights through fine-tuning or other techniques to adapt the model to their needs.

Open-weights models typically do not publish their training data, training algorithms, or other associated information. As such, it can limit the ability to perform a detailed technical inspection of the model or to reproduce the model's performance. These limitations, however, are most relevant to other researchers and are less relevant to enterprises that are seeking to simply use a model in the most efficient and effective way.

<sup>16</sup> https://arxiv.org/abs/2005.14165

#### **Open Access Models**

Open access is a growing category of models that are nearly open, but have custom terms that cannot be considered fully open source in the traditional definition of that term. It covers a wide gamut of licenses with different restrictions, and thus should be the subject of a detailed legal review to ensure that the license allows for the intended use.

#### Some examples include:

- BLOOM, which was released under the OpenRAIL-M license.<sup>17</sup> Though quite nearly open source, it has requirements for responsible use of the model, which means that it is not fully open source.
- LLaMA 2 and 3 from Meta AI, which have been released with their own custom licenses (called the LLaMA 2<sup>18</sup> and LLaMa 3<sup>19</sup> Community Licenses, respectively) that set limits to the use of the model. Specifically, the licenses forbid the use of the model in applications with more than 700 million monthly active users and for the purpose of building competitive models.

#### **Open Source Models**

Open source models are the most open of all, publishing details of their training data, training algorithms, and model parameters, allowing for the most permissive use of the model. Common open source licenses include Apache 2.0 and MIT. Meta AI's first version of their LLaMA model was released under a GPLv3 license, restricting it from commercial use and thus making it not useful for most enterprise applications.<sup>20</sup>

#### Choosing a License for Enterprise Use

Even though proprietary models are the most restrictive, they are often entirely appropriate for use in the enterprise, as with any other proprietary enterprise software. By charging for access to their models, providers of proprietary models may be able to more easily

<sup>17</sup> https://bigscience.huggingface.co/blog/the-bigscience-rail-license

<sup>18</sup> https://llama.meta.com/llama2/license/

<sup>19</sup> https://llama.meta.com/llama3/license/

<sup>20</sup> https://arxiv.org/abs/2302.13971

provide for services and support for the use of their model. This may make their use more appropriate for use in the enterprise.

Open-weights, open access, and open source models may be more useful in applications where an enterprise wants more control over the model itself and possesses the technical expertise to make any such modifications or to host the model.

#### **Model Hosting**

Enterprises have three main hosting options when looking to access LLMs:

- 1. API services from the model developers, such as OpenAI, Anthropic, Cohere, and Mistral.
- Cloud Service Providers offering hosted LLM services, such as Azure OpenAI Service, AWS Bedrock, or Google Vertex AI Model Garden. These services also allow customers to load their own models, while the underlying hardware is managed by the cloud provider.
- 3. Self-managed hosting of models. Many models with different licensing terms are available for self hosting, including open-source and open-access models as described above. Cohere also licenses its proprietary models for self-managed hosting.

In many cases, models hosted by their developer or a cloud service provider (options 1 and 2, above) are the best choice in the enterprise. In the same way that cloud computing outsourced the burden of running data centers, hosted models are a simple continuation of that trend, offering infrastructure-as-a-service. Given the intensive compute requirements of LLMs, especially under heavy workloads, outsourcing this can be a wise choice.

The most common objection to using a hosted service is that it requires sending corporate data to a third-party service. But, in many cases, this corporate data is already hosted by a third party that may also be hosting internal communications and other sensitive data (e.g., a company that uses Microsoft 365 productivity and communication tools has its data in Azure). Is using the LLM service from that same provider any different? It is ultimately a question that warrants review by your legal and risk teams, but in

most cases, the conclusion is that it is not different in a meaningful way.

Self-hosting a model requires acquiring the necessary hardware, configuring it to run the LLM, and then maintaining that stack for reliable internal use. Typically, this will require a cluster of GPUs that have been properly configured with the right drivers and packages to run the LLM in question. The LLM must then be loaded into this environment so that it can begin to serve internal requests.

Self-hosting can be an appropriate choice for an enterprise in cases where an organization needs full control over the model and the hardware it runs on and cannot use a third-party service for its data. This may be the case in the most restrictive data environments, or if the enterprise does not want to rely on a third party to ensure the performance of the environment, notably in contexts where third-party providers may need to throttle access to certain customers to ensure the overall stability and availability of their service.

In the case of both self-hosting and hosted services, applications that use the LLM will access the model through an API endpoint. The difference is simply who is hosting and maintaining that endpoint and whether the data going to and returning from the LLM leaves the corporate firewall of the enterprise.

#### **Building a Base Model Is Not for Most Organizations**

Early on in the popular interest in LLMs, a lot of attention was given to the expense and complexity of building these models. Billions of dollars were being spent building these models, and sometimes training them took many months. A huge amount of this initial work was amassing the enormous training sets required to build models of this scale.

Recent advances have brought down the time needed to build new models, and open source training data repositories now exist. But the fundamental question for an enterprise that is considering building a model remains: Why would you? Given the great diversity of today's models, which offer seemingly endless combinations of performance, specificity, licensing, and hosting options, what would justify the time and expense needed to build your own model, especially given that you are uncertain of being successful?

Any company whose core business is not building or serving AI models should not consider building their own model. There are more than enough options on the market today. The challenge is not getting access to a model, but using it safely, securely, efficiently, and effectively to further your business goals. This is where an LLM Mesh comes into play.

#### **Bottom Line: Why the LLM Mesh?**

As you have read in the previous sections, a great variety of models exists in an ecosystem that is rapidly evolving. This is ultimately a very good thing for enterprises: It means that they will be able to pick and choose the right model for the right applications within their business. Building applications that are powered by these LLMs requires combining them with other objects, like retrieval systems, prompts, and tools. This requires careful attention to many different factors:

- How the models, services, and associated objects are registered and used within the organization,
- How the data is routed to the model,
- How access to the models and services is controlled,
- How the use of the model is logged and audited,
- How the content generated by the model is moderated,
- How the models can be enriched with proprietary data,
- How can the applications be developed, deployed and maintained efficiently, and
- How can more people become involved in this process?

As more agentic applications are built and used in the enterprise, the cost and complexity of managing all of these dimensions risks spiraling out of control. This could force the enterprise to make compromises, potentially limiting the value that it derives from AI.

For example, perhaps there is a use case that would benefit from using a small, specialized model that is self-hosted and to which access is restricted. This could be a code assistant that is well-versed in the company's proprietary code libraries. If the organization lacks the ability to quickly and efficiently add this model to its mix, it may

not pursue this use case, leaving the potential gains in efficiency on the table and falling behind its competition.

This would be unfortunate, given that many of the additional capabilities that are required to use an LLM efficiently and effectively in an enterprise are common to all models.

This is the power of an LLM Mesh: its ability to reduce the cost of building an additional agentic application in the enterprise. With an LLM Mesh, an enterprise is free to develop an optimal AI strategy without compromising on performance, cost, safety, or security.

The remaining chapters of this technical guide will go into much more detail about how implementing an LLM Mesh can be done.

# Objects for Building Agentic Applications

#### A Note for Early Release Readers

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 2nd chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at <code>jbleiel@oreilly.com</code>.

LLM Mesh is a new architecture paradigm for building agentic applications in the enterprise. It enables an organization to build and maintain more agentic applications, ultimately getting more value from LLMs.

An LLM Mesh will allow you to:

- 1. Access various LLM-related services through an abstraction layer.
- 2. Provide federated services for control and analysis.

3. Provide central discovery and documentation for LLM-related objects via a catalog.

This chapter will describe the many different types of LLM-related services that are used in building agentic applications, and how they can be connected with one another. These various LLM-related services are called "objects" in an LLM Mesh. The final section of this chapter will describe the importance of the catalog for the discovery and use of the objects in an LLM Mesh.

We start with an explanation of why using an LLM Mesh to build agentic applications is increasingly important in today's competitive landscape. The bottom line is that you are going to need to build a lot of custom agentic applications.

#### The Potential of New Agentic Applications

In Chapter 1, we learned about the many different types of models available, how they work, and the options available for hosting them. What can these models be used for in the enterprise? There have been two main, initial uses of these models in the enterprise. The first is simply providing a version of the consumer chatbot experience within a wrapper that meets enterprise security and auditability requirements. The second is using these models to provide software assistants, often called "copilots," that can accelerate the use of existing SaaS products.

This first generation of enterprise use has been met with a mixed reaction. In some cases, notably when used as coding assistants for software developers, the copilots have proven to be valuable additions to the enterprise IT mix. Other feedback has been more mixed, leading in some cases to disillusionment. It is too early, however, to discount the potential for LLMs in the enterprise. This is because the second generation of agentic applications in the enterprise will be more sophisticated.

These applications will not only use the ability of these models to generate text but also their ability to solve arbitrary problems when instructed to do so. For example, an LLM could be provided with the documentation for an API that looks up the current price of a stock. With that documentation and its coding ability, the LLM can write a script to call that API for a given stock price. If allowed to execute that script, the LLM—without having ever been explicitly

programmed to do so—could then become a tool for the end user to look up arbitrary stock prices. This ability to accomplish tasks for which the LLM has not specifically been programmed is called "generalization," and it is what allows LLMs to be the engines in a new class of enterprise applications.

These new applications will provide automation and decision support throughout the enterprise. In order to do so in a reliable and cost effective manner, however, they will need to be carefully designed, tested, deployed, and monitored. While many of the constraints of traditional enterprise applications will also apply to this new class of agentic applications, the way in which they are built, and the components that are used to build them, will be different.

Given LLM's ability to generalize, would it be possible to develop a single, all-powerful application that can solve any problem and answer any question in the enterprise? In short, no. While the LLM itself is capable of generalization, the constraints of the enterprise will require that the scope of any one application be relatively narrow to ensure consistently good performance and to control access to data and tools.

For example, this imaginary, all-powerful application sounds convenient but would require full access to all of the company's data and tools, from the most mundane to the most sensitive. Just as an employee should only have access to the data and the tools that they need to do their job, so too must the access of any one agentic application be limited to that which it needs to perform its function. Furthermore, while LLMs are capable of generalization, they require quite specific instructions to deliver consistent results, often with examples of the expected input and output. This also drives towards a larger number of more narrowly-scoped applications.

Concretely, how many such agentic applications might a large corporation need? Let's do some order-of-magnitude estimations. Let's say that a large corporation has 10 departments and each department has 5 core functions. Each of these functions could potentially benefit from 5 such applications. For example, within a Sales Department, the Sales Operations function could have one application that researches their target accounts, a second that checks if the sales process is being respected, a third that continuously analyzes the health of the sales pipeline, a fourth that summarizes meetings

with prospects and a fifth that assists salespeople with their followups.

Doing the multiplication of this order of magnitude estimate gives us  $10 \times 5 \times 5 = 250$  such applications in that enterprise. Again, this is not a precise number, it's a rough estimate of the order of magnitude of applications to expect. Expecting several hundred such applications in use in a large organization seems like a reasonable estimate.

#### **Build vs. Buy**

If an organization would benefit from several hundred novel applications, where will they come from? As always, organizations will face a "build versus buy" decision. On the "buy" side of that balance, existing software vendors and new startups are already bringing these applications to market, and organizations will have a lively marketplace of competitive offers to choose from. On the "build" side of the balance, more advanced organizations are building their first production-ready agentic applications. Which approach is best? Each has its advantages, disadvantages, and appropriate uses, meaning that most organizations will buy some applications, and build others. Table 2-1 summarizes these tradeoffs and considerations.

Table 2-1. Comparing the tradeoffs of building versus buying agentic applications

	Advantages	Disadvantages	Ideal Uses
Buying Off- the-Shelf Agentic applications	Turn-key performance once implemented Developed and maintained by professional software engineers  Turn-key performance once implemented  Turn-key performance once implement	Same performance as your competitors that use the same solution     Complex to integrate with enterprise systems     Governance challenges for tracking which models are used by which applications	Non-critical functions where the goal is to gain in efficiency, not necessarily to differentiate from competitors.

	Advantages	Disadvantages	Ideal Uses
Building Custom Agentic Applications	Adapted to specific business context with the potential to build differentiated capabilities     Full control and transparency over the application     Independence from software, AI, and cloud providers	Skills required to build applications may not be available     Complexity of monitoring and maintaining grows with the number of applications	Core and strategic functions where full control and strong competitive differentiation are needed.

Agentic applications, whether they are custom-built or bought offthe-shelf, have the potential to improve the efficiency of an organization's operations. But simply improving your efficiency in lockstep with that of your competitors does not improve your competitive position in the market. If you are making the same efficiency gains as your competitors and no more, you are not becoming more competitive, you are simply keeping up.

Building custom agentic applications allows an organization to create a capability that its competitors do not possess, and thus to outperform them in that particular domain. Given the cost and complexity of building, monitoring, and maintaining these applications, organizations will choose to focus their internal development efforts on the parts of their business that stand to benefit most from strong competitive differentiation. In most cases, this will be their core business. For example, it may be R&D and supply chain management for a pharmaceutical company, or risk and price modeling for an insurance company. The needs of non-core functions will be satisfied with applications bought off the shelf.

#### The Complexity Threshold

How many custom applications can any given organization develop and maintain? Each organization is different but every organization has a maximum number of applications that it is able to develop, monitor and maintain with its current practices and techniques. We call this the organization's "complexity threshold", and it is illustrated in Figure 2-1.

As the organization develops and deploys more agentic applications, the complexity of monitoring and maintaining them increases until, at some point, the maximum complexity is reached and no more applications can be developed. Reaching this threshold means that the organization cannot develop more applications, even if doing so would benefit its business. If the organization wants to develop more applications, it must find a way to increase its complexity threshold. This requires standardizing and structuring the way that the organization builds these applications.

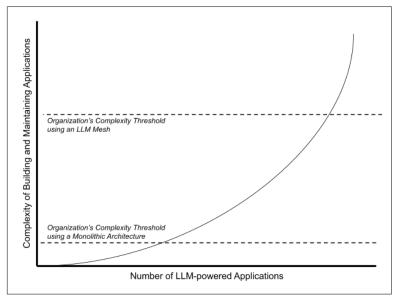


Figure 2-1. Comparing the tradeoffs of building versus buying agentic applications

#### A New Paradigm for Building Agentic Applications

Bringing standardization and structure to the way that applications are built in the enterprise is a show we've seen before. Over the years, organizations have used different architecture paradigms for developing applications. Starting with monolithic applications in the early days of application development, where all components were tightly integrated into a single codebase, organizations then shifted to an architecture paradigm with a higher degree of abstraction with the services-oriented architectures of the late nineties and, now, the modern standard of microservices has taken that abstraction even further.

Today, the architecture paradigm for building agentic applications is monolithic applications using packages like LangChain. This approach is appropriate for building your first few POCs and production applications, but it reflects the relative immaturity of agentic application design in the enterprise.

A new architecture paradigm is needed for building and maintaining many agentic applications that can raise an organization's complexity threshold. LLM Mesh is that new architecture paradigm.

Now, let's look at the objects used in building an agentic application.

#### **LLM Mesh-Related Objects: An Overview**

Building an LLM Mesh requires understanding the different types of objects that must interact with one another within an agentic application. Chapter 1 covered the LLMs and the various services that host and serve them. While those models and services are at the heart of an agentic application, more is needed, especially if the developer hopes to build a custom application that will stand apart from the competition and deliver better and more valuable performance. This requires integrating the LLMs with various objects unique to the organization.

An LLM Mesh thus treats objects of a similar type in the same way, with the LLM Mesh itself providing the translation between the generic object (e.g. a tool) and the specific service (e.g. a specific SQL database). In this way, we say that the LLM Mesh provides "abstraction" between the high-level object and the underlying, specific service.

Figure 2-2 illustrates the objects of an LLM Mesh organized into different layers that comprise the typical stack of an agentic application, overlaying the typical stack of a traditional application.

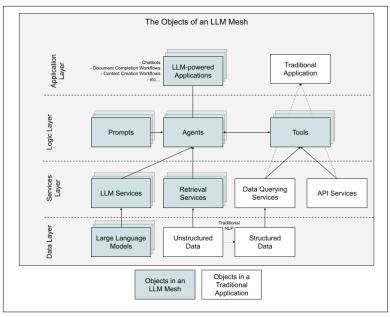


Figure 2-2. The objects of an LLM Mesh in comparison with those of a traditional application

Note that in Figure 2-2, the objects in the lighter-colored rectangles are not themselves part of an LLM Mesh, but rather are abstracted as the higher-level objects in the darker shade. This will be discussed further in the retrieval services and tools sections below. In contrast, traditional applications use Data Querying Services and API Services directly, without abstraction as tools. Unstructured Data is not used directly in traditional applications but is first transformed into structured data using traditional natural language processing (NLP) techniques.

Here is an overview of the different objects in Figure 2-1 and how they relate to one another:

#### Large Language Models

The base model — the trained neural network comprising the core mathematical weights — as described in Chapter 1.

#### Unstructured Data

Enterprise data that is not in tabular form. A common type of unstructured data is documents, which may be in PDF, DOCX,

or other formats. Unstructured data is abstracted as retrieval services in an LLM Mesh.

#### Structured Data

Enterprise data that is in tabular form, typically stored in databases, data warehouses, and data lakes. Structured data is stored in data querying services, which are in turn abstracted as tools in an LLM Mesh.

#### LLM Services

The services which are comprised of the hardware and software systems used to deploy and interact with the model in real time. As described in Chapter 1, these services may be managed by the model developer, a third party, or internally by the enterprise.

#### Retrieval Services

A service that allows for the efficient and effective querying of unstructured data. The retrieval services usually consist of an LLM used for embedding, storage for the embeddings (which can be either a dedicated vector store or another type of database—SQL or search, for example—that has added these capabilities), and some system for ranking the results to best respond to the query.<sup>1</sup>

#### **Data Querying Services**

Databases and their associated query languages, like SQL, that allow for the efficient retrieval of structured data. These systems are abstracted in an LLM Mesh as tools.

#### API Services

Any internal or external API services to be integrated with the agentic application. An external example could be a weather service to look up a forecast, while an internal example could be the data catalog to allow for data discovery. These services can be very diverse and are abstracted in an LLM Mesh as tools.

<sup>1</sup> As retrieval services are relative newcomers to the enterprise architecture landscape and are themselves powered by LLMs, they are treated as distinct objects from other tools in an LLM Mesh.

## Prompts

The input to the LLM services, they can be templated and standardized and can run the gamut of prompting techniques (few-shot, chain of thoughts, etc.).

## Agent

An LLM-powered system that seeks to accomplish a certain goal over multiple iterations within a defined level of autonomy, and using tools to meet its objective. Note the centrality of agents in this architecture. They are the object where the logic and behavior of the application are defined.

#### Tool

Any function or resource that an agent can use to accomplish its task. It can be a programming or querying language, an API service, or even another agent.

## Agentic Applications

An application that provides a user interface and other functionality on top of the agent. A chatbot is one example of an application type, but agentic applications could have many different types of interfaces running the gamut from dashboards, to mobile apps, to assistants embedded in other applications, to headless applications running behind the scenes and altering users only when needed.

# The Objects of an LLM Mesh in Detail

In this section, we will walk through each of the seven types of objects that are used in building agentic applications in the enterprise. Each section will first define and describe the object.

At the end of each section, you will find a tip box titled "Thinking Like an LLM Mesh". This box describes the expected input and output of each object. Recall that one of the main benefits of an LLM Mesh is that it creates an abstraction layer that standardizes the inputs and outputs of diverse services into standardized objects. The tip boxes summarize what those standardized inputs and outputs should be.

Building an agentic application requires integrating several different objects. For example, a simple chatbot application using a retrieval augmentation technique could be built using the following objects:

- An *application* with a chatbot interface where the end users ask their questions and receive their responses as well as provide feedback to the developers.
- An *agent*, composed of several templated *prompts*, that defines how the user's question will be handled by the LLM.
- An *LLM service* that receives the question, tokenizes it, and submits it to the *LLM* that will generate the response, enriching it with an answer from a retrieval service.
- A *retrieval service* that provides access to unstructured data from documents. The retrieval service is comprised of
  - an *embedding model* that converts the text data to vectors and
  - a *reranking model* that will select the most relevant answer to the user's question, providing it back to the LLM service for inclusion in the reply.

Such a chatbot could, of course, be built without an LLM Mesh simply by building a monolithic application that calls the various services, passing the results from one object onto the next. In practice, the developer of such an application would be writing many API calls, each of which is specific to each service. If the developer would later want to change, for example, from one third-party LLM service to another, this would require manually updating the code so that the application calls the new LLM service in the way that is expected by that service.

In that scenario, an efficient application design would provide a certain degree of abstraction, defining the interface with the LLM service as a single function within the application and not specifying the details of the API call in every instance where the LLM service is called.

An LLM Mesh takes this abstraction further, completely separating the service from the application and providing a standard interface for all objects of the same type for use across all agentic applications in the enterprise.

## LLMs

We covered LLMs in detail in Chapter 1. When we talk about an LLM, we are talking about a very large file, often measuring in

gigabytes or terabytes. For example, the 405 billion parameters of Meta's Llama 3.1 model weighs in at 2.3TB. The majority of the data volume is taken up by the weights of the model itself. Remember, as described in Chapter 1, the weights of a model are simply a great quantity of floating-point numbers.

If an organization is using a managed LLM service, they will never interact with the model itself, only with the service endpoint. But, if an organization self-hosts an LLM, then they will need to load the LLM into their hosting infrastructure.

#### Thinking Like an LLM Mesh

From the perspective of an LLM Mesh, an LLM is thus an object that can be interacted with in only two very simple ways: It can be downloaded and updated in the environment where it is hosted. These two actions will generally be done by interfacing with an API supplied by the model provider or from an aggregator of models (a model hub) such as Hugging Face.

## **LLM Services**

An LLM service is a combination of storage resources, compute resources, and supporting software that allows an LLM to be hosted and accessed for inference.

The developer of the model may provide LLM services. For example, OpenAI, Anthropic, and Mistral all provide services that run their proprietary models. In these services, the end user does not load the model into the service, they simply select the service running the model that they prefer.

Alternatively, an organization may choose to build and run its own LLM service, managing the GPUs and associated technologies.

Finally, cloud service providers (CSP) offer managed LLM services. In these, the end user may select the model that they wish to run, but the CSP manages the compute and storage infrastructure.

An LLM service, be it hosted by your organization or by a third party, is accessed via an API. Generally, most LLM services will expect similar variables when they are called. Those include:

Which model version to use

- A system prompt set by the developer to guide the model's completion
- The user prompt for the model to complete
- Temperature setting to define the level of randomness in the response
- Alternatives to temperature, such as top\_p or top\_k, use different sampling methods to determine which subsequent token to select

In response to such requests, the LLM service will generally reply with a response that includes:

- An indication of the type of response (e.g., text completion or streaming chat)
- A unique identifier of the response
- · The generated content
- · Reasons for why completion may have stopped
- Usage statistics about the number of tokens in the request and response

Most services have broadly similar expected inputs and outputs. An LLM Mesh abstracts and standardizes these inputs and outputs through its abstraction layer, ensuring that the request sent to a given service is formatted appropriately and uses the correct syntax. When using an LLM Mesh to build an application that calls an LLM service, the end user calls the LLM service object in the LLM Mesh, indicating which service to use, and the LLM Mesh translates that generic call into the specific call expected by the indicated service.

To better understand the value of providing a standard interface for all LLM services, let's compare the expected syntax of two common providers, OpenAI and Google Gemini, starting with OpenAI. The OpenAI documentation<sup>2</sup> gives the following example:

```
curl https://api.openai.com/v1/chat/completions \
  -H "Content-Type: application/json" \
  -H "Authorization: Bearer $OPENAI_API_KEY" \
  -d '{
    "model": "gpt-40",
```

<sup>2</sup> https://platform.openai.com/docs/api-reference/chat

Let's compare that with the expected request to the Google Gemini API. Google's documentation<sup>3</sup> gives the following specification:

```
curl -X POST \
    -H "Authorization: Bearer $(gcloud auth print-access-token)"
    \
     -H "Content-Type: application/json" \
        https://${LOCATION}-
aiplatform.googleapis.com/v1/projects/${PROJECT_ID}/locations/$
{LOCATION}/publishers/google/models/${MODEL_ID}:streamGenerate-
Content \
    -d '{
        "contents": [{
            "role": "user",
            "parts": [{
                "text": "TEXT"
        }]
     }
}
```

These short samples from the documentation already show some differences between the two APIs:

- OpenAI specifies the model in the JSON payload with the model key-value pair, while Google specifies the model in the URL path.
- The array containing the content of all messages is called mes sages by OpenAI and contents by Google.
- Google nests an additional array, parts, within its contents array.

<sup>3</sup> https://docs.anthropic.com/en/api/complete

An LLM Mesh standardizes these and other differences, allowing for faster application development and easy switching between LLM services. As LLM service providers update their services, the LLM Mesh developer will update the Mesh accordingly, freeing the application developers from the need to do so.

## Thinking Like an LLM Mesh

As an object in an LLM Mesh, an LLM service expects a prompt as input and is expected to provide text as its output.

## **Retrieval Services**

From the perspective of an LLM Mesh, a retrieval service takes a user's query as its input and provides a relevant result from unstructured data as its output. It is how an agentic application accesses unstructured data. In this context, the unstructured data is text data coming from documents, often stored as PDFs, plain text documents, or other common document formats, like DOCX. Retrieval services allow agentic applications to make this data available to the employees of an enterprise by allowing them to discover it more accurately and rapidly. Importantly, the information in these documents is not only made available to the employees. It will also be available for agentic applications themselves to inform their inference on how to solve a problem. Retrieval services serve this dual purpose: making the unstructured text data available to both the employee and the agentic applications which, in both cases, leads to better decisions.

In retrieval services, like traditional search systems that came before them, there is usually a tradeoff between the speed of the results and the quality of the results. While it is possible to have fast results or good quality results, it is difficult to have both. Retrieval services are usually made of three separate components to provide the highest quality results as quickly as possible. These are:

- Embedding models, which will convert the text of the document base, as well as the text of the query, into dense vector representations,
- Data storage, commonly vector stores, for storing the vectors and performing efficient similarity search, and

• Reranking models, to improve the quality of the search results.

Increasingly, these services are being provided as bundled services from various providers, as their combined functionality is required to provide the desired result to the end user: a relevant result from unstructured data in response to a natural language query. From the perspective of an LLM Mesh, the mutual dependency of these three underlying technologies is why they are combined as a single object, "retrieval services."

The following sections describe the components of a retrieval service and some of the tradeoffs that the various choices will entail.

## **Embedding Models**

As described in the previous chapter, embedding models transform text into numerical representations called embeddings, stored as high-dimensional vectors. For example, the word "banana" might have the following embedding: [0.534, 0.312, -0.123, 0.874, -0.567, ...] Each number represents the value of a particular dimension. If the vector had 100 dimensions, there would be 100 numbers in the list.

These embeddings capture the semantic meanings of the text, such that "Denver" and "capital of Colorado" will have similar vector representations, even though they share no keywords, while "kid" meaning "young goat" will have a different vector representation than "kid" meaning "young human".

Different embedding models use different embedding lengths, meaning more or fewer dimensions for each vector. Put more simply, a shorter embedding length means fewer dimensions in each vector, and thus fewer numbers in the list that represents each word or part of a word. More embeddings require more storage and compute resources. New embedding models, like OpenAI's textembedding-3 family of models, allow for the embedding size to be shortened to a degree specified by the user. Shorter embeddings can have lower storage and computational costs but may result in degraded performance. Model developers are working to increase the performance of vectors with fewer embeddings.

Embedding models expect input text that has been pre-processed to a certain degree. Different models have different requirements; an LLM Mesh provides a standard interface that is mapped to each model. Pre-processing will generally include extracting the text from any documents (e.g., PDF or DOCX formats), removing punctuation, adding special tokens to tell the model about relevant breaks in the text, and splitting longer documents into smaller "chunks" that are sized appropriately for the embedding model.

In a retrieval service, embedding models serve the dual purpose of converting the corpus of documents into vectors and then doing the same for the query. Converting the corpus of text into vectors is usually an offline task, done once while converting the query is necessarily done at runtime, when the query is received from the user.

## Vector Store, or Other Data Storage

The embeddings are written to a data store, often a dedicated vector store. A vector store is a database specially designed to store and efficiently query high-dimension, dense vectors, like those created by embedding models. Vector stores have built-in retrieval functionality for finding a stored vector most similar to the query vector, usually using a cosine similarity function.

Traditional data stores — including relational databases like Post-greSQL, document databases like MongoDB, search engines like ElasticSearch, and graph databases such as Neo4J — are all adding support for dense vector data types. As the use of vector data increases in the enterprise thanks to the growth of text embeddings used in agentic applications, the use of these more traditional data storage technologies may become increasingly relevant, reducing the need for dedicated vector stores.

This evolving technology landscape is one more reason why abstracting these services as "retrieval services" is important in an LLM Mesh. While the underlying technologies may change, the function remains the same: Provide relevant results from unstructured data to user queries.

## **Reranking Models**

The vector store's retrieval function will provide a fast result, but it may not always be the most accurate. More accurate results can be obtained by using the vector store's retrieval function to narrow down the results and then a reranking model to analyze the subset more carefully selecting the best result to return.

In contrast to the retrieval function of the vector store, the reranking model will take the entire source document plus the input query for comparison. Given that the source data could contain thousands or even millions of source documents, it would be too slow and too costly to run that process across every document. By using the retrieval function to narrow down the results to the top few (a number which can be specified), and then running the reranking model across the subset, you strike the best balance between speed and quality. This is called two-stage retrieval.

The ranked results will be the output of the reranking model. The retrieval system will provide the top-ranked result back to the user in response to their query.

#### Thinking Like an LLM Mesh

As an object in an LLM Mesh, a retrieval service expects a natural language query as its input and is expected to output the top-ranked result from unstructured data. These results are then generally passed on to an agent.

## **Prompts**

Since the popular use of LLM-powered chatbots increased dramatically following the release of ChatGPT and associated products, many of us are now familiar with the notion of a prompt. A prompt is the initial input (a question, a command, or instructions) provided to the model, prompting its response.

In contrast to the ad hoc prompting often used in consumer applications, prompting in the enterprise benefits from a structured, templated, composable approach. This allows a bank of prompts to be developed, tested, and then shared for reuse across the organization. There are many different types of such prompts. The following sections discuss several categories of prompts, with some simple examples of each.

## **Role-Based Prompts**

These prompts direct the LLM to respond as a specific type of expert, such as a customer support agent or HR consultant, or guide the AI on the tone, formality, and style of its responses.

Examples:

- "You are an IT support technician. Assist the user in trouble-shooting their software issue."
- "Respond in a professional and concise manner suitable for senior management."

## **Compliance and Ethical Prompts**

These prompts direct the LLM to provide responses that adhere to specific regulations or legal frameworks, or responses which follow specific internal guidelines for ethical practices.

## Examples:

- "Ensure that no response contains personally identifiable information (PII), such as names, phone numbers, or identifiers like Social Security numbers."
- "Generate responses that respect the following internal ethical AI guidelines [corporate ethical AI guidelines]."

While using such prompt components can decrease the risk of non-compliance, they cannot guarantee that any result will necessarily be compliant. As such, human oversight is required.

## Customization, Personalization, and Context-Specific Prompts

These prompts customize responses from an LLM based on known information about a user or customer, a prediction about them, or other relevant contextual information. The variables in the example prompts below would be completed based on information in the enterprise Customer Relationship Management system (CRM), customer support records, or using the result of a predictive model.

## Examples:

- "Personalize the marketing message for a [age]-year-old [gender] living in [postal code]."
- "Recommend to the customer [result from next-best offer prediction]."
- "Given the customer's previous request about [subject of the previous request], provide a relevant response."

## **Multi-Step Process Prompts**

These prompts guide the LLM to respond in a multi-step process by breaking down complex decisions into smaller, more manageable steps. These multi-step process prompts are the building blocks of agents.

#### Examples:

- "Step 1: Gather all financial data from Q1. Step 2: Generate a financial report. Step 3: Summarize the key findings in a presentation."
- "First, evaluate the market demand. Next, assess the cost implications. Finally, recommend a go/no-go decision."

## Thinking Like an LLM Mesh

For the purposes of an LLM Mesh, prompts need to be tested, approved, and published in the catalog, which we will explain further later in this chapter. Any prompt must be associated to a specific model and version, as small changes in the model or prompt may result in dramatically different prompt performance. These prompts can then be combined with one another to compose more complex and sophisticated prompts, themselves part of agentic applications.

## **Agents**

While various definitions for "agent" exist, from the perspective of an LLM Mesh, an agent is an LLM-powered system capable of accomplishing its objective across multiple steps using tools, without requiring prompting by an end user for each step.

Within an LLM Mesh, an agent is the object where the other objects interact with one another to form a system that can respond to users' needs. They call one or more LLM services, they use several templated prompts and use one or more tools. As such, agents are some of the most important objects within an LLM Mesh and are at the core of building agentic applications in the enterprise.

Like the other objects, they must be built, described, cataloged, and maintained. As the maturity of an organization increases, it will begin to develop more agents, and will likely start chaining those agents together, with one agent using another as a tool. This increasing complexity can be tamed by the abstraction and modularity that the LLM Mesh offers.

There are a few important parts to the above definition of an agent, so let's look at them one by one.

## **Objective**

An agent's developer will define its objective by giving it a role-based prompt, as described in the previous section. For example, an agent that is part of an application that is designed to generate real-time sales analytics could include the following role-based prompt template:

You are a Business Intelligence Analyst with access to the company's sales data across various regions and time periods. Your role is to assist in retrieving specific data as requested by the user and to provide additional analysis that highlights any interesting, unusual, or noteworthy aspects of the data, just as a human analyst would do. When the user makes a request:

- Accurately identify the relevant data source and retrieve the specific data they are asking for.
- Perform a detailed analysis on the retrieved data to uncover any trends, anomalies, or key insights. Consider aspects such as:
  - Comparisons with previous periods or other regions.
  - Significant changes or trends in the data.
  - Potential reasons behind the observed data patterns.
- Any other insights that might be valuable for the user to  $\ensuremath{\mathsf{know}}$  .

Finally, present the data and your analysis in a clear, concise summary that the user can easily understand. If the user's request is unclear or requires data from multiple sources, use your judgment to clarify the request and combine data sources as needed to provide a comprehensive analysis.

In this example, the objective is clearly described, as is what the agent should do if the end user asks it to do something outside of its prescribed scope.

## **Multiple Steps**

Agents will execute multiple steps to meet their objectives. These individual steps are linked in chains, which define the steps the agent must take to meet the objective. This differentiates agents

from the simple, direct prompting of an LLM. For example, asking an LLM to summarize a block of text cannot be considered an agent because it is a single step.

Take for example an agent that has been built to summarize financial reports. The multiple steps might be:

- 1. Call an API to download the desired report(s)
- 2. Locate and extract key figures from the report
- 3. Look up historical values for these figures and compare them
- 4. Extract key quotes from the report
- 5. Generate a semi-templated summary that includes both extracted quotes, generated summary text, and comparison between historical and current figure
- 6. Send the report to the recipient over the specified channel

Each step would include a templated prompt that would be modified with either the user input or the LLM output from the preceding step. The steps are strung together in a chain, which may be sequential, looping, branching, or parallel chains. Throughout this process, multiple calls to the LLM service will occur without any user involvement.

## Autonomy

An agent is granted some degree of autonomy. Using the analytics-generating agent as an example, a minimal degree of autonomy may simply be deciding which Python package or function to use during the data analysis step. A more significant degree of autonomy may be choosing the tool that it will use to meet its objective from several made available to it (e.g., deciding if it should query historical data from a data warehouse or live data from a CRM to best respond to the user's request).

Less autonomy will mean that the agent is less flexible in the type of problem it can solve, but more likely to give a good result in that narrower range. More autonomy will mean more flexibility, but more risk that the results will not be satisfactory. In the enterprise, agents are likely to be quite limited in their autonomy, with narrowly defined options available to them, especially during the early stages of their development and use. This may change over time as models and agent-building techniques evolve and improve.

#### Tool Use

A defining characteristic of an agent is its use of tools to accomplish its objectives. These tools are described in more detail in the following section.

#### Thinking Like an LLM Mesh

As an object in an LLM Mesh, an agent expects some task as an input and is expected to provide a satisfactory result as an output. This broad definition reflects the breadth of what agents can be built to accomplish.

## Tools

In an LLM Mesh, a tool is any function or system that an agent is provided with to accomplish its task. As such, tools cover a very wide range of potential technologies. This breadth gives agents and agentic applications their incredible potential: they can automate and accelerate tasks, decisions, and operations that otherwise require manual work across the enterprise and its business systems.

The types of systems that can serve as tools in an LLM Mesh include but are not limited to:

- Internal data storage and retrieval systems, such as databases, data warehouses, and data lakes.
- Enterprise software systems, such as CRM, Human Resources Management System (HRMS), and Enterprise Resource Planning (ERP) systems.
- Advanced analytical assets, like predictive machine learning models.
- Programming and querying languages, like Python and SQL, along with specific packages or proprietary code.
- External data APIs, such as financial data or weather services.
- Other agents within the LLM Mesh

For an agent to use a tool, it needs to understand what the tool is and how to use it. This is accomplished by creating a schema for each tool. This schema is what allows for some standardized interaction with the tool, despite the great diversity of tools that may exist in the LLM Mesh. The schema should include:

- A description of the tool, including examples of the circumstances in which it should be used.
- Instructions on how to interact with the tool, including what input is expected and what output is expected.
- Connection details for accessing the tool.

By ensuring that each tool has a well-described schema, the tools can be used across different agents, including those with a high degree of autonomy, as those agents will rely on the descriptions in the schema to decide which tool to use.

#### Thinking Like an LLM Mesh

As an object within an LLM Mesh, a tool provides a schema, making itself available for use by an agent. The tool expects an input and provides an output as defined in that schema. Tools are very flexible and their schema is essential to their use by agents.

## **Applications**

In an LLM Mesh, an application is what makes an agent available to end users. The agent defines the logic that orchestrates the different objects from the LLM Mesh that are required to accomplish a specific purpose. The application is the interface and supporting functions that allow the end users to interact with the agent, to better understand the results provided by the agent, and to provide feedback to the developers. The application is also where certain services providing security, safety, and cost control are enforced.

There are several types of agentic applications, including:

- Chat interfaces where users interact with the agent iteratively.
- Contextual assistants, either as desktop applications or browser extensions, that provide some additional functionality or assistance in the context where the user is working at that moment.
- Backend or "headless" applications that run without direct end user interaction.

Agentic applications can have a wide range of interfaces and functionality. The abstraction and standardization within the LLM Mesh makes it simpler for the developer to build the application in a way that clearly communicates to the end user how the agent underpinning the application is generating its results.

For example, in the case of an application that exposes an analytics-generating agent to end users, it will be easier for the end user to understand and trust the results if the application distinguishes between outputs that come from a query to a retrieval system or a tool versus outputs that are the result of the LLMs interpretation or suggestion. Furthermore, the end user will also be more likely to trust results if they can verify that the sources used, and the query that the LLM generated, are appropriate for the objective of the application.

Feedback mechanisms should also be built into the application to ensure that when an agent does not behave as expected, end users can flag this anomaly to the developers so that they can monitor the agent's performance and take corrective action if necessary.

#### Thinking Like an LLM Mesh

Within an LLM Mesh, the application object includes the application itself, versioning for the deployed application, and logging of user interactions with the application.

# Cataloging LLM-Related Objects

As an organization begins developing more agentic applications, the number of different objects it will need to use to build those applications will grow rapidly. This could become difficult to manage, with users hunting for different objects, recreating existing objects, or using unapproved objects.

Overcoming these challenges starts by creating a central catalog for all of these objects. This catalog is a fundamental component of an LLM Mesh. The catalog should:

- Account for all LLM-related objects that are available for use in the enterprise.
- Provide documentation that describes and provides instructions for using each object.
- Track the version or other details about the ownership and development history of the object.

 Assign a unique ID to each object to allow it to be referenced and tracked unambiguously.

This information is stored in a structured format that allows human and machine discovery of the available objects. Having a central catalog of the objects provides various benefits for organizations as they begin building more agentic applications. Those benefits include:

#### Standardization

Only approved objects can be added after a vetting process.

## Governance and Compliance

You can maintain full transparency and traceability of which data are used with which LLM for which purposes, enabling business alignment and regulatory compliance.

#### Security

The catalog allows access controls to be defined and enforced, controlling which end users and automated systems have access to what objects.

#### Composability

Once registered, objects can be easily added to new applications where they are combined with other objects, accelerating the development process.

## Efficiency

Less time is spent manually connecting different objects, accelerating application development.

Importantly, this catalog will be useful for both the end users and the LLM-powered agents that they will be building. The agents will also rely on the documentation to discover and use the objects and the agents will be subject to the security model.

## Conclusion

In this chapter, we have learned about the various objects of an LLM Mesh, how they are abstracted, and how they can be integrated with one another. The final chapter of this guide will go into greater detail about how a specific agentic application can be built using an LLM Mesh. Before getting to that, however, the following chapters will describe the various federated services that an LLM

Mesh must also provide to meet enterprise security, reliability, and cost requirements for the many agentic applications that will be built within it. Chapter 3 will start with that most important of enterprise considerations: cost. How can the overall cost of an enterprise's LLM use be optimized? Read on to learn more.

Conclusion 55

# Quantifying and Optimizing the Cost of LLMs in the Enterprise

## A Note for Early Release Readers

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 3rd chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at <code>jbleiel@oreilly.com</code>.

At the beginning of Chapter 2, we established the importance of an organization's ability to develop its own agentic applications. These built-for-purpose applications will allow the organization to differentiate itself from its competitors, pulling ahead in the market. That said, developing and running them will come at a cost.

Chapter 2 established how an LLM Mesh will simplify and standardize the development of these applications, thus reducing the cost of their development. However, agentic applications will also incur costs as they run. Minimizing costs while maintaining the required level of performance will allow an organization's budget to go fur-

ther, supporting the deployment of these built-for-purpose agentic applications across more functions of the organization, thereby getting more return from the same investment in AI.

It is critical not to consider cost in a vacuum; rather, it must be considered in conjunction with performance. If an organization were to drive down its spending on AI blindly (for example, by using smaller models or cheaper services), it could find that its agentic applications are no longer delivering adequate performance or could stop functioning entirely. On the other hand, if an organization had a policy of always using the highest-performing model, it would likely end up overspending for a level of performance that is, in fact, not needed. The answer is to find the balance where cost is minimized for an adequate level of performance. The necessary level of performance, both in terms of speed and quality of output, will depend on the specific requirements of each application. Thus, the goal is neither minimal cost nor maximal performance but rather minimal cost while delivering the required performance.

An LLM Mesh must provide the federated services required to measure and track cost and performance across the many different components used to build agentic applications. By using these federated services for cost and performance tracking in an LLM Mesh, organizations can fully understand where their AI budget is going and enforce policies to get the most return on that investment.

We will tackle the topics of cost in performance in both this chapter and in Chapter 4. In this chapter, we'll first understand the drivers for the cost of the different objects in an LLM Mesh, especially the different types of LLM services. Then, we will look at techniques for limiting costs. Finally, we'll consider the organizational practices needed to run a cost-efficient AI practice using an LLM Mesh, covering topics like cost reporting, budgeting, and rebilling. Throughout, we'll refer to the need to consider performance tradeoffs when making decisions to reduce costs. Rest assured that we will bring much more precision to performance measurement and monitoring in Chapter 4.

# **Quantifying the Costs of Agentic Applications**

In this section, we will understand how agentic applications generate costs. We'll start with a quick review of the objects of an LLM Mesh to determine which are cost-generating and which are not.

From there, we'll then dive deep into the costs of the different types of LLM services, including those provided by model developers, CSPs, and those that you manage yourself. Finally, at the end of this section, we'll compare two different agentic applications to see how choosing different LLM services impacts their costs.

## The Objects in an LLM Mesh That Drive Costs

Let's begin by breaking down the cost structures of the different objects in an LLM Mesh. Recall from Chapter 2 the seven main objects of an LLM Mesh: LLMs, LLM services, retrieval services, prompts, agents, tools, and applications. Of these seven objects, some are cost generating and others are not, as follows:

## Non-Cost-Generating Objects in an LLM Mesh

## Prompts

Like code or documentation that an organization would develop and manage, prompts themselves have no direct cost.

## Agents

Agents are logical objects and thus generate no direct cost themselves, though when used, the underlying LLM services, retrieval services, and tools will generate costs.

## Agentic Applications

Like the agent, the code of an application generates no direct costs.

#### LLMs

As data objects, the LLMs themselves incur no direct costs, though their use through LLM services (below) does. Any licensing fees for proprietary LLMs are baked into the cost of using the service. (An exception to this would be a proprietary LLM that is licensed for use in a private deployment, where that licensing cost would need to be factored into the total cost, in addition to the cost of running the hosting service yourself.) Cohere, AI21 Labs, and Aleph Alpha are examples of model developers who offer their proprietary LLMs to be licensed for private deployment.

## Cost-Generating Objects in an LLM Mesh

#### LLM Services

LLM services are the workhorses of agentic applications and their main cost drivers. Several different cost models for LLM services will be explored in more detail in the following sections.

#### Retrieval Services

Retrieval services generate costs per use, much like the LLM services that often underpin them.

#### Tools

Certain tools, like data querying services, are fixed costs that the organization already bears, whereas others, such as external APIs, will generate costs per use.

The cost monitoring and control services of an LLM Mesh will focus on the LLM and retrieval services. The costs of the tools (data querying services, API services) are either fixed costs that don't vary with application usage or are already well-managed by existing API cost tracking and management solutions. The services of the LLM Mesh need to be adapted to the cost structures and cost control techniques unique to LLM services and retrieval services, which are new assets in the enterprise IT landscape.

#### Additional Costs in an LLM Mesh

In addition to the costs of the objects combined to build agentic applications, organizations may incur the following costs:

- Licensing fees for any LLM Mesh services for analysis and control,
- Licensing fees for a cataloging and documentation solution,
- Licensing fees for other components of an LLM Mesh,
- The costs of developing and maintaining any of the LMM Mesh components or infrastructure that they do not license from a software vendor. Note that these costs, which may be significant, are out of the scope of the LLM Mesh and, thus, this guide.

Whether these components are licensed separately, as part of an all-in-one LLM Mesh solution, or developed internally will depend

on the organization's LLM Mesh strategy, and each approach has different costs.

## **Understanding the Costs of LLM Services**

You must set up your LLM Mesh to capture the total costs of all agentic applications in a normalized way, allowing for their comparison and aggregation. By having a comprehensive and comparable way of looking at the costs of different LLM services, organizations can more easily build capabilities that combine models from different providers to strike the optimal balance between cost and performance across their entire fleet of agentic applications. Without the cost tracking and control services of an LLM Mesh, the costs of these different services are difficult to compare and control, making it challenging for organizations to optimize their spending

As described in Chapter 1, there are three main types of LLM services based on where the LLM is hosted:

## Model developer-managed services

These are the LLM services offered directly by model developers, such as OpenAI, Google Gemini<sup>1</sup>, Anthropic, Cohere, or Mistral. They are Model-as-a-Service (MaaS) offerings, and costs are calculated on a per-token basis. Many developers will offer both on-demand pricing and reduced batch pricing.

## Cloud Service Provider-managed services

These are the LLM services offered by cloud service providers (CSPs) such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP). The cloud platforms have given specific branding to their LLM services: AWS Bedrock, Google Vertex AI, and Microsoft Azure AI Studio<sup>2</sup>. These services generally offer access to a number of curated LLMs that can be deployed on managed compute environments or as serverless MaaS instances.

<sup>1</sup> Google offers both its Gemini model directly to customers, as well as through its Google Vertex AI product. The Google Vertex AI product offers the models of other developers as well. We will refer to the Gemini direct offering as "Google Gemini" and the offering that includes models from other developers as "Google Vertex AI" to avoid confusion.

<sup>2</sup> Microsoft Azure's foundation model catalog is also accessible through their Azure Machine Learning Studio product.

## Self-managed services

Organizations may choose to manage their LLM services themselves. In this case, they may rent a server from a cloud provider or provision an on-premises server for themselves.

The following sections will explore the details of these different options and their pricing models, giving some examples of how those pricing models impact total costs. As you will see, there are diverse pricing models and options available, making the centralization and normalization of the cost service of an LLM Mesh an essential element for deployments that go beyond initial testing and experimentation.

## The Cost Model of Developer-Hosted LLM Services

The pricing model for the LLM services offered by the model developers mostly follows the same MaaS structure: a price per token input (i.e., the prompt) and a price per token output (i.e., the generated instructions). The price per output token is generally three to five times as expensive as per input token. Table 3-1 summarizes the on-demand pricing of several developers' flagship models, showing the discrepancy between input and output tokens.

Table 3-1. On-demand pricing of developers' flagship models as of October 2024

Developer and Model	Input Tokens (per 1M tokens)	Output Tokens (per 1M tokens)
OpenAl GPT-4o	\$2.50	\$10.00
Google Gemini 1.5 Pro (>128k context window)	\$2.50	\$10.00
Mistral Large 2	\$2.00	\$6.00
Anthropic Claude 3.5 Sonnet	\$3.00	\$15.00

These services will typically offer their latest models at a few different sizes, with larger models being more expensive and smaller models being cheaper. They may also offer specialized models, such as those for generating text embeddings, generating images, or generating code.

In addition to paying as you go, some model providers are beginning to offer options for reducing costs through two different

approaches: offering lower pricing for asynchronous *batch processing* and reducing the input tokens needed through *context caching*.<sup>3</sup>

## **Batch Processing**

While some applications, like chatbots, require immediate responses from LLM services, other applications do not. For example, using an LLM to classify the sentiment in a large dataset of historical customer reviews can be done asynchronously, with the request being sent and the response being returned sometime later. In the case of OpenAI, a 24-hour turnaround time for all batch jobs is guaranteed. Check with your provider to see if they provide reduced pricing for batch processing, it could be a valuable source of savings for the right kind of application.

When implementing an LLM Mesh, be sure to do so in such a way as to take advantage of batch processing where the LLM service offers it. And be sure that the LLM Mesh takes the price difference into account for cost tracking and analysis.

## **Context Caching**

There are some cases where you will include repeated information in a prompt submitted to an LLM service. For example, a chatbot may include a lengthy system prompt at the beginning of each request to ensure consistent and appropriate behavior. The LLM service will retain a cache of these input tokens and can use that cache instead of the new tokens, resulting in reduced latency and cost. Some of the LLM services pass these savings on to their customers.

Google Gemini prices cached input tokens at a 75% discount for uncached tokens, while OpenAI offers a 50% discount for cached tokens. Both the OpenAI and Google Gemini services apply caching automatically — no change to your API calls is required to take advantage of it. That said, it is important to carefully read the developer's documentation to understand under what conditions cached tokens will be used so that you can create the conditions to maximize these savings. For example, OpenAI only checks the first few tokens to determine if the cache should be used or not, meaning

<sup>3</sup> As of October 3, 2024, OpenAI and Google Gemini offer both batch processing and context caching, and Anthropic offers context caching (called Prompt Caching). Cohere and Mistral do not currently offer these options.

that you need to structure your prompts such that the repeating information is always at the very beginning.

The LLM services will report when cached tokens are used. An LLM Mesh should track this and tabulate costs accordingly.

## The Cost Model of CSP-Managed LLM Services

The major cloud providers have all introduced managed services that allow popular open-source and proprietary models to be quickly deployed on infrastructure managed by the cloud provider. This can be a quick way for organizations with an existing relationship with one or more CSPs to test and begin building with LLMs from different developers.

In contrast to the offers from the model developers, which all follow a MaaS paradigm, the CSPs offer more pricing models and options. This introduces some additional complexity, but it also introduces more options that allow an organization to optimize for cost and performance as well as data residency.

AWS, Azure, and GCP all offer two main pricing models: ondemand and provisioned throughput.

## **On-Demand Pricing**

On-demand is a MaaS offer that is billed on input and output tokens (like the MaaS offers from the model developers). In most cases, the models offered for on-demand pricing will come from several curated developers. For example, in addition to its own Gemini models, Google Vertex AI offers models from AI21 Labs, Anthropic, and Mistral. There is thus a tradeoff between the convenience of the fully managed MaaS offer and the full choice of models available on model hubs, such as HuggingFace.

In contrast to the offers of the model providers described previously, only AWS Bedrock offers batch pricing for their on-demand offer, giving a 50% discount. This pricing is only available for certain models in certain regions.

## **Provisioned Throughput**

In contrast to their on-demand offers, the provisioned throughput offers from the CSPs reserve a defined amount of model capacity for a given amount of time. The primary benefit is guaranteeing resource availability for your applications. Thus, it is not primarily a cost-saving strategy, but rather a way to guarantee the availability of required resources. That said, longer reservation periods result in a discount over the CSP's on-demand pricing.

In contrast to the per-token pricing of on-demand offers, provisioned throughput offers are less transparent. In a provisioned throughput offer, you reserve "model capacity" which each CSP defines and names differently. On Azure, they are known as "Provisioned Throughput Units" and an estimator tool is provided, but not the underlying formula. On AWS Bedrock, they are named "model units" and you must contact your account manager to estimate the quantity that you would need.

Provisioned throughput offers are best adapted to applications that have already been tested and deployed to production and thus where the expected usage is relatively certain. If you have an application where you anticipate variable usage, the on-demand offers are probably the best bet.

CSPs generally do not provide as many models in their provisioned throughput offers as they do in their on-demand pricing offers. For example, Google only offers its Gemini models for provisioned throughput, and Microsoft Azure only offers models from OpenAI. AWS Bedrock offers both its Titan family of models and models from other developers.

## The Cost Model of Self-Managed LLM Services

The previous two categories of managed services offer much convenience, but at the cost of control. While it is easy to begin using the managed services with a simple API call, you are limited in the models available and other aspects of the deployment. If more control over those details is required, then an organization may consider self-managing its own LLM services. This can be done by renting a server instance from a cloud provider or using an on-premises server, typically maintained by your organization's IT department.

Cost will be an important consideration when deciding whether to deploy self-managed LLM services. In contrast to the per-token or model unit costs of the managed services, the direct costs of the self-managed services will be the hourly cost of the server instances. The costs will thus be fixed, independent of the number of tokens you push through the service. A cost-conscious organization would thus seek to maximize the use of its service without overloading it and degrading its performance. A delicate balance must be struck, for sure!

The cost model for a self-managed LLM service will first depend on whether you are renting the instance from a cloud provider or using an on-premises server. In the case of an on-premises server, the cost model will depend on your organization's internal rebilling policies. In the case of instances rented from a cloud provider, two primary cost models exist: on-demand usage, and savings plans that require a long-term commitment, typically one or three years. The following sections explore these in more detail.

The pricing of the server instances available to rent from the major cloud providers will depend on its combination of GPU memory, compute, and storage resources. GPU memory is generally the constraining resource that you need to take into account, depending on the size of the LLM that you would like to host. This is because all of the model's weights need to be loaded into the GPU memory. Very large models with tens or hundreds of billions of parameters will exceed the memory of even the largest single GPUs and must therefore be deployed to multi-GPU instances. To choose the right server, you will need to estimate the throughput required for the anticipated use of your application(s).

You will need to ensure that the server instance has all of the necessary drivers and software to serve the models. The cloud providers typically offer base images for these instances that provide the necessary drivers, but be sure to check the documentation from the model developer to make sure that you are not missing anything.

## Costs of Self-Managed, On-Demand Cloud Server Instances

On-demand servers are usually priced per hour, and you begin paying the moment that it is activated, regardless of the number of tokens that it is processing. In this way, it is a fixed cost. It is important to thus spin down any instances that are not in use.

Even though the cost of the instance is expressed per hour, most CSPs offer finer granularity for billing, often down to the second. So if you use the instance for 8 hours, 12 minutes, and 23 seconds, you will only pay for those 29,543 seconds.

That said, be wary of attempting to over-optimize your instances by spinning them up and down too frequently, as there will be a cost in terms of latency. It will usually take a few minutes for an instance and the connected services to come fully back online and to be available for use. Thus, "on demand" only refers to the lack of long-term billing commitment and should not be misinterpreted as meaning "available for immediate use." On-demand instances are most appropriate for predictable workloads that will remain constant for a period of time.

## Costs of Self-Managed, Long-Term Cloud Server Instances

If you expect the workload of your application to remain constant over a year or more, long-term commitments can be very cost-efficient. By agreeing to pay for a predetermined volume of server usage, the CSPs will usually offer a discount over the equivalent on-demand price. The discount rate will depend on the instance type, the region, and whether you pay upfront. Discounts can be as large as 70% but may be smaller in practice, particularly for in-demand instance types like those for LLM workloads.

Note that the commitment implies that the instance is paid for the entire duration, meaning that you cannot spin it down during periods of unuse. Thus, a three-year commitment implies multiplying the hourly cost by the number of hours in three years (26,280, assuming there is no leap year).

## Costs of Self-Managed, On-Premises Server Instances

Some organizations, though fewer and fewer all the time, prefer — or are required — to manage their server infrastructure. If your organization operates this way, you are undoubtedly aware of it already, and you are undoubtedly aware of the process of gaining access to these resources. In this case, a guide like this will not be able to provide you with detailed information, as the costs and how they are rebilled will depend on your organization's practices.

In all cases, an LLM Mesh should allow you to move easily among the services available to your organization and must capture the costs of all different types of services. As the pricing of many hosted services is not available programmatically (i.e., they can not be looked up via API), an administrator of your LLM Mesh should track them and update them carefully to ensure that the cost data is captured accurately.

## **Comparing the Costs of Applications**

To make these differences more real, let's imagine two different applications in a large enterprise. These applications will have very different usage scenarios to show the consequences of using different applications with different services.

## Application 1: Company-Wide Knowledge Assistant

This first application is a knowledge assistant application that uses a RAG pipeline to surface relevant information to the user from an internal document base. We imagine this application will be deployed in a large, global enterprise. As a result, it sees relatively constant usage across time zones and throughout the year. As a RAG application, a large volume of text extracted from the source documents is passed to the LLM service in the prompt, resulting in a large volume of input tokens. Let's build out a low-volume and a high-volume scenario to estimate the total token count for the LLM service over a year; those estimates are shown in Table 3-2.

Table 3-2. Usage hypotheses for company-wide knowledge assistant application

	Low Usage Scenario	High Usage Scenario
Input Tokens per Session	5,000	15,000
Output Tokens per Session	500	5,000
Sessions per User, per Day	5	20
Users	1,000	1,000
Working Days per Year	250	250
Total Input Tokens per Year	6,250,000,000	75,000,000,000
Total Output Tokens per Year	625,000,000	25,000,000,000

Now, let's compare the costs of running this application against two different LLM services. The first is an on-demand service from a model provider. We'll use OpenAI in this example. The second is a self-managed cloud server from AWS. This type of application does not need the advanced reasoning capabilities of the largest, most sophisticated models. As such, we'll choose the OpenAI GPT-40 mini as the MaaS option. This scenario is detailed in Table 3-3.

Table 3-3. Pricing scenario for OpenAI GPT-40 mini on-demand as of October 2024

	Low Usage Scenario	High Usage Scenario
Total Input Tokens per Year	6,250,000,000	75,000,000,000
Total Output Tokens per Year	625,000,000	25,000,000,000
Input Price (per 1M tokens)	\$0.15	\$0.15
Output Price (per 1M tokens)	\$0.60	\$0.60
Total Input Cost	\$937.50	\$11,250.00
Total Output Cost	\$375.00	\$15,000.00
Total Cost	\$1,312.50	\$26,250.00

For the self-managed option, we will select the Llama 3.2 11B model and will run it on an AWS g5.2xlarge EC2 instance. This instance should be sufficient for the hypothetical usage, but it is important to model the required resources appropriately for your expected use, as different instance configurations have very different prices. In this scenario, we will compare the costs for no commitment (ondemand), a one-year commitment, and a three-year commitment. As these costs are fixed regardless of usage, they would be identical for the low and high usage scenarios; hence, we show cost in columns per commitment period.

Table 3-4. Cost for AWS EC2 g5.2xlarge (US East, Ohio) as of October 2024

	No Commitment (on-demand)	1-Year Commitment	3-Year Commitment
Cost per Hour	\$1.21	\$0.95	\$0.65
Cost per Year	\$10,617.12	\$8,360.98	\$5,733.24

In most cases, self-managing the AWS EC2 instance would be less costly, with the exception of the lowest usage scenario. In the high usage scenario, making a three-year commitment and self-managing the EC2 instance delivers 78% savings. This shows the importance of accurately estimating and monitoring your usage of LLM services and making the correct choice per application.

## **Application 2: Corporate Strategy Sparring Partner**

Let's now imagine a very different application. This application is built on top of an agent that is designed to support the strategic planning activities of the senior leadership and corporate strategy teams. It is used far less per year, by far fewer people. That said, the usage is very intensive, and it requires a very large, very capable model. Let's go through the same exercise, starting with the usage hypotheses.

Table 3-5. Usage hypotheses for corporate strategy sparring partner application

	Low Usage Scenario	High Usage Scenario
Input Tokens per Session	5,000	50,000
Output Tokens per Session	5,000	25,000
Sessions per User, per Day	10	50
Users	20	20
Working Days per Year	50	50
Total Input Tokens per Year	50,000,000	2,500,000,000
Total Output Tokens per Year	50,000,000	1,250,000,000

As with the previous application, let's now consider the costs of running this application. To keep the comparisons similar, we'll use OpenAI again as the MaaS provider, but we'll select their flagship model, GPT-40.

Table 3-6. Pricing scenario for OpenAI GPT-40 on-demand as of October 2024

	Low Usage Scenario	High Usage Scenario
Total Input Tokens per Year	50,000,000	2,500,000,000
Total Output Tokens per Year	50,000,000	1,250,000,000
Input Price (per 1M tokens)	\$2.50	\$2.50
Output Price (per 1M tokens)	\$10.00	\$10.00
Total Input Cost	\$125.00	\$6,250.00
Total Output Cost	\$500.00	\$12,500.00
Total Cost	\$625.00	\$18,750.00

Imagine now that you want to self-manage this application. You want to use a large and highly capable model, so you choose Llama 3.2 90B. Given the size of this model, you will need to use an EC2 instance with multiple GPUs, in this case, a g5.12xlarge. That choice results in the pricing shown in Table 3-7.

Table 3-7. Cost for AWS EC2 g5.12xlarge (US East, Ohio) as of October 2024

	On-Demand	1-Year Commitment	3-Year Commitment
Cost per Hour	\$5.672	\$4.4667	\$3.06288
Cost per Year	\$49,686.72	\$39,128.23	\$26,830.83

For an application that is only used for part of the year and requires a high-performing model, the most cost-efficient approach will be to use the MaaS offering from the model provider, even in the highest usage estimate. Paying for a self-managed instance of sufficient size for the model required would be wasteful.

We've intentionally constructed the high and low usage estimates to show the high potential variability between applications. For example, applications with an agent can cycle through many prompts and responses during chain-of-thoughts reasoning. This can result in a large volume of both input and output tokens that the end user never sees but which are necessary to the proper functioning of the application. Monitoring this spending is essential to make sure that the budget is used wisely and that costs are minimized.

# **Techniques for Limiting Costs**

In the previous sections, we reviewed how to measure and monitor costs in an LLM Mesh. Now, how can you limit costs, and how should you approach a new project in a cost-conscious way?

When building out a new agentic application, a good rule of thumb is to start by defining the level of performance that you need for a given application without much consideration for cost. This is because it is important to confirm if the application will work at all and not be left wondering if it would have worked if a more powerful and more costly model had been used instead. Thus, you should start with a large, high-performing LLM that will allow you to build out your application quickly. Then, once you have achieved the level of performance that you need, you can start optimizing your application, testing lower-cost models and techniques to maintain the level of quality and speed that your application requires while minimizing the cost. The following sections review some of these techniques.

## LLM and LLM Service Selection

The simplest method of reducing cost is to switch to a cheaper LLM service that provides the required level of speed and performance. There are two ways to think about this: model upgrade and model substitution.

## Model Upgrade

Upgrading a model is when you move from one model to its successor. For example, upgrading from OpenAI's GPT-4 to GPT-4o. The good news is that, given the rapid rate of development in the space today, the new generation of a given model is often both better performing *and* cheaper than its predecessor. Such is the magic of living in a time of rapid technological progress, and while this trend may not continue, you should take full advantage of it.

An LLM Mesh should monitor for these upgrade opportunities and then facilitate the rapid testing to confirm that the new model does, in fact, provide improved application performance. You may ask, "How could it not?" While the new model will almost certainly outperform the previous model on the benchmark metrics, that does not necessarily mean that it will automatically improve the overall performance of the application. For example, your application may depend on very carefully crafted prompts that exploit some quirk of the outgoing model. The new model may not behave in precisely the same way, leading to degraded performance. Once you adapt the prompts to the new model, you can likely achieve equivalent, if not improved, performance at equivalent or lower cost.

Upgrading a model may seem like a trivial task. However, the many steps in the previous paragraphs demonstrate several of the benefits of using an LLM Mesh. First, with an LLM Mesh it is quick to introduce the new model — no application logic needs to be changed. Second, performance measurement in the LLM Mesh allows the developer to quickly measure the performance with the new model, taking corrective measures if needed. Finally, cost tracking would show the difference in the cost between the application with the old model and the application now with the new model. By reducing the effort required for each of these tasks, an LLM Mesh makes it easier and faster to improve agentic applications, integrating improved models as quickly as they become available.

#### Model Substitution

In addition to upgrading to a new version of the same model, a viable option may be to consider a different model or family of models. Here, the decision will not be triggered by the release of a new model. Instead, a team will seek opportunities to reduce the cost of its application, testing lower-cost options to see if they can be used while maintaining the quality required for the application.

#### Service Substitution

Finally, in certain situations, consider migrating from one LLM service to another, even if the underlying model is the same. For example, imagine an internal knowledge management chatbot that sees regular, predictable usage throughout the year. Your organization anticipates that it will continue to use this application for several years without much modification. If it runs on a CSP-managed instance with their on-demand option running an open weights model (Llama 7B, for example), you could consider migrating to a self-managed service running on a server instance with a long-term commitment to reduce your costs. Note that it will now be incumbent upon your organization to manage this service.

## **Prompt and Inference Optimization**

Whether using a managed service or self-managing your LLMs, optimizing your prompts to use fewer input and output tokens while delivering the desired results is one of the most direct ways to reduce costs. Well-established manual techniques (popularly known as prompt engineering) as well as emerging programmatic techniques exist for optimizing prompts. Additionally, an LLM Mesh can help optimize inference through caching. The following sections describe these techniques.

## **Manual Prompt Engineering**

The principle of manual prompt optimization is to craft better, often shorter, prompts. Many guides are available online, often providing helpful advice. An overview of prompt engineering techniques is beyond the scope of this guide, but in general, being direct and concise in the instructions and clearly specifying the expected output will help minimize the volume of input and output tokens per request.

Remember, in an LLM Mesh, prompts are objects that are developed, tested, registered and reused. Thus, optimized prompts that deliver good results at low cost can be shared and reused across multiple applications, reducing the time needed for application developers to test new prompts. In an LLM Mesh, a prompt must be linked to a specific model, as not all prompts perform equally across all models.

#### **Prompt Compression**

If LLMs are good at generating text, and a prompt is just text, can we ask an LLM to improve a prompt? The answer is yes, and there is much active research happening at the frontier of programmatic prompt optimization. These methods generally use an LLM to analyze a prompt, to understand which tokens in the prompt are actually driving the desired result from the model, and then compress the prompt down to these essential tokens. The result is often a prompt that is unintelligible to a human reader but which delivers the desired results with far fewer tokens. These approaches are particularly useful in prompts that provide a lot of context to the model, such as those in RAG pipelines, where a large volume of text is sent to the model.

Two popular approaches to prompt compression are LLMLingua<sup>4</sup>, a research project developed by Microsoft, and Selective Context<sup>5</sup>, an open-source project developed by academic researchers. These approaches can allow up to a 20x compression rate (i.e., reducing input tokens by 95%) while maintaining the required performance.

Implementing prompt compression within an agentic application, particularly one with long internal prompts — such as a RAG pipeline or an agent using a chain-of-thoughts reasoning technique — can be a powerful strategy for reducing the input tokens to the LLM service, thereby reducing costs.

By offering prompt compression as an option, an LLM Mesh ensures that application developers don't spend time repeatedly implementing such techniques.

<sup>4</sup> https://www.microsoft.com/en-us/research/project/llmlingua/

<sup>5</sup> https://pypi.org/project/selective-context/

#### Caching

As described in this chapter, some LLM services offer reduced pricing for prompts where context caching can be applied. An LLM Mesh can be configured to detect opportunities where a prompt can be adapted to trigger context caching when using a service that proposes this option. Usually, this means making sure that prompts with similar content start with identical strings of text.

In addition to ensuring that context caching is triggered when possible, an LLM Mesh can provide its own caching capabilities. For example, an LLM Mesh can detect when a prompt being sent to a particular LLM service is identical to a prompt sent recently. Rather than send the new prompt, it can simply supply the previously generated response, avoiding the unnecessary regeneration of the response and reducing both latency and cost.

#### **LLM Modification**

In the previous sections, we have looked at ways to reduce costs without changing the LLM itself. Those techniques can thus be used on proprietary models or with LLM services that don't allow for model modification. A more advanced approach is to modify the model itself so that it can provide the required results at lower cost. The following sections review those approaches.

#### Model Quantization

In model quantization, the precision of the model weights is reduced by converting them typically from 32-bit floating point numbers to 16-bit floating point or 8-bit integers. In layperson's terms, the more precise values are rounded to a less precise value with fewer significant digits. This means the calculations throughout the model during inference are simplified, resulting in lower computational time and cost. This, however, can come at the cost of accuracy. As such, the results of a quantized model should be tested to ensure that they meet the quality requirements.

An LLM Mesh can offer model quantization as an option for any LLM objects. Note that this is only relevant for open-weight models running in a self-managed LLM service.

#### **Model Pruning**

While model quantization seeks to reduce the computational intensity of inference by reducing the precision of all weights in the model, model pruning seeks to reduce computational intensity by reducing the number of weights the model.

Emerging techniques, such as LLM-Pruner<sup>6</sup>, developed by a team of researchers from the National University of Singapore, show promising results, though fine-tuning may be required after pruning to ensure the overall quality of the model.

As with model quantization, you should make sure that your LLM Mesh offers model pruning as an option for any LLM objects. Note that this is only relevant for open-weight models running in a self-managed LLM service.

#### Fine-Tuning

A much-discussed approach to improving the cost efficiency of an LLM is to fine-tune it to a particular context. From a cost perspective, however, the fine-tuning process will entail its own costs. But, it can be less costly to use a fine-tuned LLM in the long run if, in order to get the results that you need, you find yourself having to provide many examples to the LLM in your prompts. By fine-tuning, the LLM can permanently "learn" to give the results that you want, reducing the cost-per-use from that point. When fine-tuning a model, you will typically retrain the task-specific layers or the final layers (heads) on a specialized dataset. In the example of marketing copy, this would be reference texts that exhibit the desired style. If the fine-tuning is successful, the model will be able to correctly mimic the style without requiring multiple examples in the prompt. While there is a fixed cost to the fine-tuning process, this may be recouped in the reduced run costs using this model for this application.

When deploying an LLM Mesh, be sure to provide methods that allow for fine-tuning as an option for any open-weight LLM object. This way, the developers using your LLM Mesh can use the approach without implementing the process themselves.

76

<sup>6</sup> https://arxiv.org/abs/2305.11627

# **Cost-Efficient AI Operations in the Enterprise**

In the previous sections, we learned about what drives the cost of agentic applications, how performance can be measured and balanced against cost, and many different techniques that can be used to reduce cost while maintaining the required performance.

However, running a cost-efficient AI practice requires more than knowing about or accessing the best cost-reduction techniques. It requires organizational policies and practices that ensure that those techniques are fully applied. This section is about these organizational aspects.

# **Tracking and Reporting Costs and Performance**

A major advantage of an LLM Mesh is that it allows you to track the cost and performance of all of your agentic applications in a single place. The centralization of this tracking is an essential part of a well-governed Generative AI practice. Without the standardization that an LLM Mesh provides, however, you would need to manually aggregate this information from the many different applications being developed across the organization, each built in a heterogeneous way. Doing it this way would be a major barrier to obtaining a single view of all cost and performance data.

With regards to cost data in particular, be sure that the LLM Mesh that you deploy track the costs in a fine-grained manner, allowing for the costs to be aggregated across multiple dimensions, such as:

- Individual users
- Teams
- Departments
- Projects
- Functions
- Business priorities
- Usage type (experiments, development, production)
- Geography and region
- LLM provider
- LLM type
- LLM (version-specific)

• Hosting architecture (model provider, CSP(s), self-managed)

By associating costs with all of these dimensions, you can generate the reports required by leadership to understand where the AI budget is being spent and to what benefit. It will also allow you to quickly identify the root causes of any anomalies. For example, if the costs of several different projects spike and they all use the same model, it could be that a change in the model has resulted in degraded performance for certain shared prompts, which would then need to be corrected.

Reports across these dimensions can be used to support a culture of transparency and accountability for all AI costs. Specifically, these reports can be shared with:

- The developers of the application, to help them to understand the ultimate cost of the applications that they are developing and to help make them aware of the consequences of their design choices,
- Management and budget owners, supporting a culture of transparency and accountability for all AI costs.

# **Setting and Enforcing Budgets**

Based on the tracking and reporting capabilities described in the previous section, an LLM Mesh can also allow you to set and enforce budgets.

From a technical perspective, setting a budget is simple: You set a value that should not be exceeded for a given cost dimension or combination of cost dimensions. From that number, an LLM Mesh can allow you to enforce that budget in several ways. Here are some examples of how that budget can be enforced in increasing order of severity:

#### Warnings

The LLM Mesh can alert the budget owner that they will soon reach or have reached their budget. This ensures that they are aware of the issue and can take appropriate action.

#### Throttling

78

If a budget has been exceeded, the LLM Mesh can throttle cost-incurring traffic to LLM services. This will help limit the exceedance without entirely interrupting service. That said,

it will degrade service and thus may not be appropriate for production-deployed applications.

#### Blocking

Beyond throttling, an LLM Mesh can also be configured to block certain LLM services if the budget has been exceeded. This could be useful, for example, in the case of an experiment that has gone wrong: The developer may not be aware of the costs they are generating, and an automatic block can prevent a costly and embarrassing overrun. On the other hand, it would be inappropriate to block a customer-facing application.

# **Rebilling Policies**

Organizations may pursue a policy of rebilling the cost of running an agentic application to the business unit that benefits from its use. An LLM Mesh can support this practice in the following ways:

#### Transparency

The business units that bear the cost can clearly understand how the application works and which parts drive the cost.

#### Reassurance

The business units can be informed about the cost-limiting techniques that have been applied, reassuring them that the application has been implemented in the most cost-efficient way possible.

#### Trust

The business unit can trust that the costs are captured accurately.

#### **Fairness**

The business unit can be reassured that other business units are also bearing the costs of the applications that benefit them and that the costs are being calculated consistently.

# **Defining and Enforcing Cost-Saving Policies**

As your organization builds more and more agentic applications using an LLM Mesh, you should define cost-reducing policies. These policies can take a variety of forms:

#### Recommendations and best practices

As part of your training guidance for using the LLM Mesh, you can provide recommendations and best practices for all the cost-saving techniques that can be applied.

#### Periodic application reviews

Once deployed and running, the LLM Mesh enables you to review the cost profile of the applications using the centralized cost tracking and reporting.

#### Mandatory approvals

Within an LLM Mesh, you can enforce a mandatory cost review and approval process to ensure that the best practices have been applied before deploying an application to production.

#### Automated detection of cost-saving opportunities

As an LLM Mesh is aware of all applications and tracks their costs, it can automatically detect opportunities to reduce the costs of the applications. The LLM Mesh could then raise an alert to trigger a review and decision about whether to apply that technique.

#### Automatic application of cost-saving techniques

Certain techniques can be applied automatically in the background to all applications. For example, prompt compression could be applied automatically for any prompts beyond a certain threshold.

# Conclusion

There is no guarantee of which approach will result in the best combination of cost, speed, and quality when developing and deploying agentic applications. As such, the best that any organization can do is to test the many different approaches to optimize their applications, to learn from this experimentation, and then seek to generalize the best practices as policies.

An LLM Mesh supports this in several ways:

 An LLM Mesh makes the different cost-saving techniques easily available to all application developers. This ensures that an individual developer does not waste time implementing these methods themselves.

- 2. An LLM Mesh provides *visibility* into all costs so that they can be measured in development and monitored during deployment, facilitating reporting and budgeting.
- 3. An LLM Mesh allows cost-saving policies to be *enforced*, ensuring that best practices are respected.

As such, an LLM Mesh allows cost-efficiency to become a core strength of an AI practice. This means that an organization can successfully develop more agentic applications in more business domains for the same budget. This cost efficiency is essential to maximize the value that your organization derives from generative AI.

But life is not without tradeoffs. These cost-reduction techniques may degrade the performance of your agentic applications. In the next chapter, we will learn about how to measure and monitor performance so that you can maintain the required level while keeping costs at a minimum.

Conclusion 81

# Measuring and Monitoring the Performance of Agentic Applications

# A Note for Early Release Readers

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 4th chapter of the final book.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at <code>jbleiel@oreilly.com</code>.

Now that we understand the cost models for various LLM services, let's look at how the performance of LLMs -- and the agentic applications built on top of them -- can be measured. Remember, the objective of this measurement is to be able to define the required level of performance for a given application and then find the lowest-cost way to deliver that level of performance.

But first, what do we mean by "performance"? Indeed, there are two very different notions here:

- 1. The *quality* of the generated response. In other words, how well the generated response corresponds to the requirements of the application. This is the main focus of this chapter, as measuring the quality of LLM responses is a novel field where an LLM Mesh can provide significant value to an application developer.
- 2. The *speed* of the service, in terms of how quickly a response is generated. We'll touch on this briefly at the end of this chapter as this monitoring is similar to established DevOps practices for monitoring the speed and responsiveness of API services. An LLM Mesh does not need to do more than capture these metrics.

Within the two dimensions of quality and speed, there are several sub-dimensions as well. Figure 4-1 illustrates these dimensions and sub-dimensions as a tree diagram.

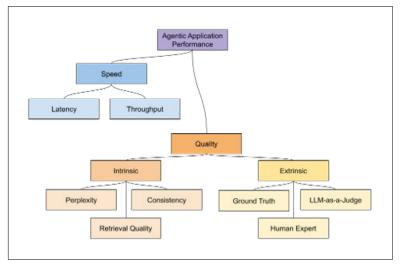


Figure 4-1. Dimensions and sub-dimensions of measuring the performance of agentic applications

The sections that follow will explore these dimensions in more detail. But we'll start by describing how an LLM Mesh architecture helps us to measure performance.

When we talk about the performance of an agentic application, it is important to note that we are *not* talking about LLM benchmarks. LLM benchmarks are the standardized tests to which LLM developers and others submit their models with the goal of comparing the inherent performance of different LLMs to one another. Some common benchmarks include massive multitask language understanding (MMLU), instruction-following eval (IFEVal), and graduate-level Google-proof Q&A (GPQA). Generally, they submit a standard list of tasks to the model and compare the results.

These benchmarks are useful for comparing one model to another and thus can help choose which model you want to start building with. However, they do not measure the quality of the output of an agentic application for your specific needs, nor do they measure the responsiveness of your LLM service. Metrics that measure the quality of output and responsiveness of your service are critical for managing a fleet of agentic applications and are thus the focus of this chapter.

# How an LLM Mesh Helps Measure Performance

When implementing an LLM Mesh, you should:

- 1. Provide a shared service for performance measurement and monitoring and
- 2. Allow a developer to use that service flexibly at different levels in their application.

To begin, it is important to provide performance measuring and monitoring tools as a shared service for two reasons. The first is efficiency: You don't want your application developers to develop and redevelop similar common capabilities across many different applications. The second is consistency: You want to be able to compare performance across different LLMs and agentic applications. If every developer is implementing their own performance measurement metrics in their own way, it will be difficult, if not impossible, to gather consistent performance metrics across your growing fleet of agentic applications.

Then, it is important to allow your application developers the flexibility to apply the performance measurement and monitoring service of the LLM Mesh at different levels of their applications throughout the many calls to an LLM that an agentic application will make. When implementing an LLM Mesh, you should ensure that the evaluation can be made at the level of the call to the LLM service and not higher up in the application stack. By bundling the LLM service call and the evaluation call, it means that developers need to call only one API and that performance evaluation can be deployed consistently across all agentic applications.

The concepts of performance and quality in an LLM Mesh are similar to that of data quality in a data mesh. Quality metrics for agentic applications should be understandable and readily available because they are a key part of the "contract" that an agentic application has with the people using it. In the same way that a data mesh is designed to establish and enforce such a quality contract so that end users are sure that they can trust the data that they are using, an LLM Mesh does the same by establishing a performance contract for agentic applications.

As we will see, reliably assessing the quality of an agentic application will require a combination of different techniques. An LLM Mesh should provide these different techniques as shared services that developers can experiment with and use freely without having to implement them themselves for each and every application they are developing. As with the methods for controlling cost discussed in Chapter 3, providing pre-implemented performance monitoring methods to all developers of agentic applications in your organization as a shared service is a critical part of implementing an LLM Mesh. This way, they continue their focus on their applications and not on common components like performance monitoring.

# Measuring and Monitoring the Quality of Generated Text

There are three phases in the lifecycle of an agentic application where the quality of the generated text should be monitored:

Pre-development

Defining the quality metrics and their required levels.

#### Development

Measuring the quality of the results and iteratively improving the application to improve the quality metric.

#### Deployment

Monitoring the quality metric for changes and taking necessary measures to fix issues that arise.

In the past, if you were developing a predictive machine learning (ML) model, proceeding through these steps would be rather straightforward. You would choose a metric depending on the business requirements of the application (e.g., in certain applications, you may be more sensitive to false positives than false negatives or vice versa) and set a minimum threshold for that given metric. You would then train the model to deliver the required performance. Finally, you would deploy the model and set up a monitoring capability to ensure that it continues to deliver the required level of performance in the face of real-world data. These processes are now well-established with standard best practices to follow.

Measuring the quality of agentic applications is very different for two reasons:

- 1. *Model outputs are non-deterministic.* With a traditional ML model, the same inputs always give the same outputs. In this sense, it is deterministic. With an LLM, the same input the prompt will generate different outputs.
- 2. The models are used in open-ended contexts. In a given application, the same model might be called on to select a tool, generate code, evaluate a response, and generate a text response for the user, requiring different quality metrics and requirements for each interaction.

Both of these characteristics are features and not bugs of LLMs. The fact that they are non-deterministic means that they can mimic creativity and come up with novel responses. The fact that they can be used in open-ended contexts means that they can be flexibly used to solve many different problems. These are their strengths, but it means that measuring their quality is much more difficult than measuring the quality of predictive models, where statistical methods alone can be used to prove whether a model is performing well or not.

Evaluating the quality of LLM outputs is done in two ways:

- 1. Measuring how well the model is performing independent of any particular task known as *intrinsic quality* and
- 2. Measuring how well the model is satisfying the task at hand known as *extrinsic quality*.

Both approaches are useful in measuring and monitoring the quality of LLM outputs. Measures of intrinsic quality are useful to identify problems in model performance early on in a relatively low-cost way. At the same time, measures of extrinsic quality are necessary to ensure that the output of the model is actually solving the task at hand in a way that a human expert would judge as appropriate. When implementing an LLM Mesh, both intrinsic and extrinsic evaluation techniques are required.

# **Intrinsic Quality Evaluation**

Intrinsic quality measures assess how well the LLM is performing, independent of the task at hand. These are different from the LLM benchmarks mentioned in the note box at the beginning of the chapter. Those benchmarks are actually extrinsic measures of how well an LLM performs at a standardized task. Instead, the intrinsic measures look at several factors of the inherent performance of the model. Those include:

- 1. *Perplexity*, or how confident the model is in predicting each token.
- 2. *Consistency*, or how similar model outputs are when provided with the same input.
- 3. *Retrieval quality*, or how effectively an underlying retrieval system is finding relevant documents.

Let's now look at each of these three factors.

# **Perplexity**

Perplexity measures how "surprised" a model is by the text it generates or encounters. Lower perplexity generally means the model finds the text more predictable, indicating confidence. Note that low perplexity (i.e., high confidence) does not guarantee that the output

is accurate; it only calculates that the token that it has just generated is very likely the best one. That said, the model may be confidently wrong.

While perplexity does not tell us if the response of the LLM service is accurate or not, high perplexity can be an indication that something has gone wrong and the model is entering unfamiliar territory. In the context of an agentic application, this might mean that a prompt is too imprecisely written and that the model has to take a wild guess at what the correct response might be.

For example, if high perplexity is measured at a step where the agent chooses from a list of available tools, it indicates that the model is not sure of its choice and may have chosen a different tool in another, similar case. By monitoring perplexity, you can detect this situation and take steps to fix the problem. For example, you could improve the prompt by providing more detail to the model about the task at hand or you could the schemas of the tools so that the model better knows which one to choose.

Another important aspect of measuring perplexity is that it can provide input to the model itself about the best next step to take. For example, by feeding an indication of perplexity back to the model, the agent can choose to escalate the task for human review or stop an automatic process. These are important performance and safety measures that are only possible if perplexity is measured consistently.

# Consistency

Even if a model is fully confident in its response as measured by low perplexity, an enterprise application requires consistent performance. It is, therefore, important to measure the model's consistency. This is done by setting up a process where the model is provided with the same input and the similarity of the output is measured using traditional Natural Language Processing (NLP) techniques (most often cosine similarity).

By continually monitoring the consistency of a model within an LLM Mesh, you can detect early on any problems that may affect the quality of any applications that are built on top of the model.

# **Retrieval Quality**

Many agentic applications will include a retrieval system, as described in Chapter 2. These systems identify relevant passages from a corpus of texts and return them to be included in the response from an LLM service. Consistently evaluating the quality of retrieval is essential for ensuring that the systems deliver useful, accurate content.

A system designed to measure retrieval quality should evaluate both the relevance of retrieved documents and their effectiveness in supporting the system's overall goals, such as providing accurate answers. It should track how often the most useful information appears at the top of search results and identify gaps where relevant content is missed. To achieve this, the system should combine automated metrics with human feedback. Additionally, it should support continuous improvement by highlighting patterns in errors, enabling adjustments to search algorithms, ranking methods, and data quality to enhance retrieval accuracy over time.

By measuring and monitoring retrieval quality, developers of agentic applications can ensure that their systems consistently deliver accurate, relevant information, leading to more reliable and effective user interactions. This process helps identify weaknesses in how information is retrieved and ranked, allowing for targeted improvements in search algorithms and data management. Ultimately, measuring retrieval quality results in smarter, more responsive applications that can better understand and meet user needs.

# **Extrinsic Quality Evaluation**

What does it mean for the output of an agentic application to be good for a particular purpose? Given the broad range of potential contexts, the best we can say is that a good output is one that serves its intended purpose well. This means that, in order to measure the quality of the output, you will necessarily need to define what good looks like on a per-application basis. This definition should be made with the input of a human expert, which can then be encoded in a "golden dataset" that can serve as the ground truth for future automated evaluation processes.

Take, for example, a very simple application that categorizes customer service requests according to a company's specific product

category. For the application to work properly, the request must return only the name of the category, not an entire sentence. In such an application, what constitutes a good answer is two-fold:

#### The accuracy of the response.

Did the LLM correctly classify the request according to the organization's categories? This would likely require a golden dataset to test against where, for example, customer service requests are categorized correctly, given that the knowledge is specific to the organization and may be difficult for an LLM to determine without this additional guidance.

#### The format of the response.

Did the LLM respond with only the category title, as instructed? This can be evaluated simply with a small glossary and a predefined rule.

Now, imagine a more complex agentic application where the agent must choose the appropriate tool for a task from among a list provided to it and then interact with that tool correctly based on the information in the schema before returning a properly formatted response to the user. This will require several tests that are specific to each step: selecting the correct tool, using the tool correctly, and providing a properly formatted response to the user.

Once you have defined what a desirable outcome is for each subtask of your agentic application, you then must define a metric for that measurement and a method to generate that metric. There are many methods available, and when implementing an LLM Mesh, it is important to ensure that you make a wide range of methods available to your application developers so that they can choose the right combination of methods and metrics for each application. In this way, the LLM Mesh will save developers time by ensuring that they don't have to waste time on implementing methods for quantifying the quality of their applications, and ensure that the organization can be confident in the results.

There are three main categories of methods for measuring the quality of the output of LLMs, each with its own advantages and disadvantages. Those categories are human expert methods, statistical methods that compare responses to ground-truth, and LLM-based evaluation methods. The following sections review those categories.

# **Human Expert Methods**

The most reliable — but least scalable — method for measuring the quality of the output of an LLM is simply asking a knowledgeable person if the response is good or not. This approach is frequently used during the development phase of a new agentic application. For example, in a data analytics-generating application, a data analyst could provide input on how they would solve a given problem, and the LLM could be prompted to deliver similar results.

However, once an application is deployed, the volume of content generated would make it infeasible to use human evaluation to monitor the quality — you are not going to pay a human to review every output of the LLM. That said, there are two ways that human feedback can be used to monitor quality once the application is deployed:

- 1. Include a method for simple user sentiment tracking using a binary feedback button (e.g., thumbs up, thumbs down). Given the low effort required of the user, it is possible to collect a meaningful sample of responses, though they may lack detailed contextual feedback on why the user responded the way they did.
- 2. Have a human expert analyze a sample of results. Such checks are an important part of quality monitoring, and they should be designed to ensure that the expert user reviews a representative sample of all responses.

When using human evaluation, it is important to ask the human evaluator to rate the output in a consistent manner across several dimensions. Simply asking if the response is "good" will not get you the information that you need to improve the application. Depending on the context of the application, you may consider asking the experts to evaluate the responses across some of the following dimensions:

#### Relevance

Does the output align well with the query and user intent? How well does it address the core needs of the business use case?

#### Accuracy

Are the facts and information presented correct and free of errors? This is especially crucial for applications in fields requiring high precision, like finance or healthcare.

#### Clarity and Coherence

Is the output easy to understand, logical, and well-structured? A high-quality response should be clear and devoid of ambiguous or confusing language.

#### Completeness

Does the output provide a sufficiently complete answer, or are key details missing? Depending on the use case, an answer that is too brief or superficial may not be helpful.

#### Conciseness

Is the response free of unnecessary information or verbose explanations? It's often important for enterprise applications to deliver only what's needed, especially when users may need quick answers.

#### Actionability

For applications with practical implications, can the user easily take action based on the output? This is relevant in customer service, recommendation systems, or task-based LLM applications.

#### Tone and Style

Does the tone fit the enterprise's needs? For instance, customerfacing applications may need a friendly tone, while internal documentation tools might require a formal, straightforward approach.

#### Bias and Fairness

Is the response free from harmful or biased statements? Evaluating for fairness ensures inclusivity and adherence to ethical standards.

#### Safety and Compliance

Does the output avoid unsafe suggestions or violations of regulatory standards? This is essential in sensitive domains like legal, financial, or medical applications.

#### Adaptability and Contextual Awareness

Can the LLM handle context changes or follow-up questions accurately? This dimension matters in dynamic environments where information evolves or multi-step tasks are involved.

#### *Novelty and Creativity (if applicable)*

Does the LLM offer innovative solutions or ideas? This can be particularly valuable in domains like marketing or R&D, where unique insights are beneficial.

When implementing an LLM Mesh, make sure that it can capture the results of human expert evaluations in a structured way. These evaluations are extremely valuable as they can form the basis of more automated evaluation approaches that we'll describe later, some of which require a golden dataset with accurate responses to compare against.

### **Ground Truth-Based Statistical Methods**

For certain use cases, you will have documented examples of correct outputs that you can use as your ground truth. In these cases, you can set up quality measures that compare the model output to that ground truth. The following three statistical methods are frequently used:

#### *Translation: Bilingual Evaluation Understudy (BLEU)*

Measures the overlap of n-grams between the generated output and a reference text. It's commonly used to evaluate the quality of machine translation.

Summarization: Recall-Oriented Understudy for Gisting Evaluation (ROUGE)

Primarily used for summarization, ROUGE compares the recall of n-grams, specifically how much of the reference summary is captured by the generated output.

#### **BERTScore**

As an alternative to relying on n-grams, BERTScore uses the BERT pre-trained transformer model to compute the similarity between the generated text and reference text on a token level. This method captures semantic similarities rather than just n-grams. It is useful in evaluating the responses in chatbot applications, as well as for evaluating the quality of translations and summarizations.

These statistical methods will generate metrics that are familiar to individuals who have developed ML models previously, namely accuracy, F1 score, precision, and recall.

Your LLM Mesh should make these methods available to your developers, as they can provide a useful point of reference for evaluating the quality of LLM outputs. That said, these methods are appropriate only for specific cases and cannot evaluate more complex responses. They also have the benefit of requiring relatively few computational resources to compute, in contrast to the LLM-based methods described in the next section.

#### LLM-Based Evaluation Methods

When performing extrinsic evaluations of agentic applications, moving beyond the limitations of traditional statistical methods and human expert methods requires turning to a method that can interpret and analyze the applications' varied and open-ended outputs: evaluation methods that use LLMs themselves. Though potentially counterintuitive, LLMs can be used to evaluate the quality of their own output when they are carefully instructed on how to do so, similar to how a human instructor grades the tests of their human students.

It should be noted that this method — often referred to as LLM-as-a-judge — is an area of active research and experimentation. The techniques described in this section are liable to evolve or be superseded by improved methods in the near future.

Given that LLM-as-a-judge methods rely on LLMs, they essentially become agentic applications themselves.

The core notion of LLM-as-a-judge methods is that you develop a prompt or series of prompts that will instruct an LLM to evaluate a specific response aspect. A response aspect is simply a formal definition of the different qualities that you may be looking for in a response. Table 4-1 summarizes these aspects.

Table 4-1. Overview of the different response aspects that you may wish to evaluate, adapted from "GPTScore: Evaluate as You Desire"<sup>a</sup>

Aspect	Task	Definition
Semantic Coverage	Text summarization	How many semantic content units from the reference text are covered by the generated text?

Aspect	Task	Definition
Factuality	Text summarization	Does the generated text preserve the factual statements of the source text?
Consistency	Text summarization, Dialogue response generation	Is the generated text consistent in the information it provides?
Informativeness	Text summarization, Data to text, Dialogue response generation	How well does the generated text capture the key ideas of its source text?
Coherence	Text summarization, Dialogue response generation	How much does the generated text make sense?
Relevance	Dialogue response generation, Text summarization, Data to text	How well is the generated text relevant to its source text?
Fluency	Dialogue response generation, Text summarization, Data to text, Machine translation	Is the generated text well-written and grammatical?
Accuracy	Machine translation	Are there inaccuracies, missing, or unfactual content in the generated text?
Interest	Dialogue response generation	Is the generated text interesting?
Engagement	Dialogue response generation	Is the generated text engaging?
Specific	Dialogue response generation	Is the generated text generic or specific to the source text?
Correctness	Dialogue response generation	Is the generated text correct or was there a misunderstanding of the source text?
Semantically appropriate	Dialogue response generation	Is the generated text semantically appropriate?
Understandability	Dialogue response generation	Is the generated text understandable?
Error Recovery	Dialogue response generation	Is the system able to recover from errors?
Diversity	Dialogue response generation	Is there diversity in the system responses?
Depth	Dialogue response generation	Does the system discuss topics in depth?
Likeability	Dialogue response generation	Does the system display a likable personality?
Flexibility	Dialogue response generation	Is the system flexible and adaptable to the user and their interests?

Aspect	Task	Definition
Inquisitiveness	Dialogue response generation	Is the system inquisitive throughout the conversation?
<sup>a</sup> Fu, Jinlan, et al. "Gptscore: Evaluate as you desire." arXiv preprint arXiv:2302.04166 (2023).		

When setting up an LLM-as-a-judge system, you will have to make many decisions, including:

- Which LLM should be used as the judge? Within an LLM Mesh, this can be any of the LLM services made available to the developers. Certain LLM-as-a-judge methods seek to provide good results with smaller models, reducing the cost of the evaluation.
- Can the evaluation be completed with a single interaction with the LLM (called a single-turn method), or will it require multiple interactions (called a multi-turn method)?
- Does the evaluation method require reference answers or not?
   If the method requires reference answers, you need to develop or otherwise provide the golden dataset that the LLM judge will refer to.
- Do the evaluations of the LLM judge correlate with the responses of a human expert judge? This is the key question! The LLM judge does not know if the answer is good or not; only a human expert can confirm this. When setting up an LLM-as-a-judge method, the goal is to ensure that its responses correlate with the responses of a human expert, meaning that if a human expert would rate one aspect of a response positively, then the LLM judge would rate it in a similar manner.

Note that as LLM-as-a-judge systems are based on LLMs which are, by definition, non-deterministic, the same judge may not always give the same evaluation to the same response. That said, the same can be said about human expert judges. Your goal when implementing an LLM-as-a-judge method is to design the prompts in such a way that the LLM gives consistent and reliable evaluations of the responses. Multiple libraries and templates of such evaluations have been developed and are available from both proprietary and open source providers. Examples include OpenAI Evals, Arize Phoenix, and RAGAS.

# Implementing Evaluation Methods

When implementing an LLM Mesh in your organization, you would decide which methods to make available to your developers and provide them as a shared service. The statistical and LLM-based methods described in the previous sections are all available as open source implementations from their original authors that you could freely use in your LLM Mesh.

Alternatively, rather than creating your own implementations of these open source methods, you could choose to use a third-party service, connecting it to your LLM Mesh. Like with LLM services, an LLM Mesh architecture should allow for connecting to both self-hosted and third-party hosted evaluation services. Third-party evaluation services include those that are offered by the cloud service providers (Amazon's SageMaker Clarify, Google's Vertex Gen AI Evaluation Service, and Microsoft's Azure Machine Learning Prompt flow which includes templated evaluation flows), as well as those offered by a host of emerging startups (for example, Lyzr, Humanloop, Cognition Labs, among others and with more emerging every month).

So, should you build your own implementations of open-source evaluation methods or purchase third-party evaluation services? On the one hand, building your own implementations of these methods within your LLM Mesh can be a good choice, as it allows you to make fine-grained decisions about how the evaluation service functions, including which LLM service they use and ensuring that they are compatible with your LLM Mesh. On the other hand, buying third-party evaluation services avoids this development effort but at the cost of customization. Just as your organization will likely choose to buy some agentic applications and build others, you face the same "build or buy" choice when it comes to evaluation services.

From the perspective of an LLM Mesh, it is important to treat these evaluation methods in a consistent manner, regardless of whether they are provided by a third party or if you implement them yourself within your LLM Mesh. When implementing your LLM Mesh, ensure that it allows you to define and capture the following dimensions of an evaluation:

• Response being evaluated, associating it with the relevant application and agent capturing the version of each. This ensures

that you have traceability of your evaluations and how they evolve with different experiments.

- Evaluation *method* and service being used, including any information about the version of the service being used. For example, a small change to the evaluation service may result in very different evaluation results in some cases, so taking into account the version of the evaluation service is essential.
- *Aspect* of the response being evaluated.
- *Metric* being calculated, including a precise definition and the formula for its calculation.
- *Value* of the metric evaluated, so that you can track the evolution of the metric over time as you continue your experimentation or monitor the application while it's in production.

Note that all third-party evaluation services may not expose this information in their API, meaning that you will not be able to capture it in your LLM Mesh, making it difficult to get a complete picture of performance across your portfolio of agentic applications. You should take this into account when choosing among different third-party services. Given the value of having consistent and comparable performance metrics that are centrally aggregated in an LLM Mesh, you may find that you will need to implement your own services in many cases.

Let's make this all more concrete by looking at a simple example of how performance measurement would work for someone developing an agentic application within an LLM Mesh architecture. In this example, let's imagine that the developer is working on an agent that includes a RAG-enriched response at one point. They want to evaluate and monitor if that response is only making claims that are backed up by the documentation that it is using for its retrieval.

First, the developer would capture the *generated response* and log its content plus information about its associated LLM Mesh objects (the application it is serving, the LLM service that it is using, etc.). Then, the developer may choose to use RAGAS as their evaluation *method* as it is designed to work well with RAG applications. The *aspect* that they want to measure is whether the response provides appropriate context, and the specific *metric* is called context recall. Context recall is calculated by simply dividing the number of claims in the generated text that can be attributed to the source text by the

total number of claims made in the generated text.<sup>1</sup> If this value is 1, it means that all of the claims in the response are based on claims in the source, while a value of 0 would mean that none are. This *value* would be logged then as the result of the evaluation.

Then, a second LLM-based evaluator might evaluate the quality of the written response, checking to ensure that it meets the expected tone. That evaluation would also need to capture the generated response, method, aspect, metric, and value for the tone of the response.

In this scenario, the two LLM-based evaluators measured different parts of the response at different moments in the agent's logic chain. Given the very different natures of these tests, it is necessary to use different evaluations. There could be several more evaluations required for such an agent, depending on its complexity. This multiplicity of evaluations shows the importance of providing evaluations as a shared service so that the developers can focus on creating and perfecting the logic of the agent and not be slowed down by the important but repetitive work of setting up robust evaluations.

# Implementing a Performance Architecture in an LLM Mesh

As described in the previous sections, when implementing an LLM Mesh you will make different quality assessment methods available to your app developers. They will, in turn, apply these methods at different levels within an agentic application, generating metrics that they can then monitor over time. But what are the best practices for how these different methods and metrics can be combined to ensure that the applications are performing as intended? While the state-of-the-art is a rapidly moving target, let's describe a simple performance architecture that could be implemented as an organizational best practice for all agentic applications using an LLM Mesh.

<sup>1 &</sup>quot;Context Recall." Ragas Documentation, October 14, 2024, https://docs.ragas.io/en/sta-ble/concepts/metrics/available\_metrics/context\_recall/. Accessed 5 Nov. 2024.

# **LLM-Level Monitoring**

At the most basic level, your organization should put in place systems to consistently monitor the performance of your primary LLMs to ensure that they are delivering consistent performance. Changes in performance may occur when either the model is updated to a newer version or when the nature of the input changes, even slightly.

In the first case, model updates may result in unexpected and potentially degraded performance, especially if your team members are crafting highly specialized prompts. These prompts may depend on quirks of a particular model version which may disappear when a model is updated. While it is true that newer versions of models generally offer improved performance, often at lower cost, it is important to be able to monitor performance so as not to be caught unaware of changed performance in a production application.

In the second case, even when a model stays the same, it is possible that your input to the prompt may change over time without you realizing it. For example, if you are using an agentic application to process customer service requests, the content of those requests may shift as your organization introduces new products and services. This may result in changed performance of your application, requiring you to take some action to return the application to its desired level of performance.

In both of these cases, you need to design experiments where you compare the real-world results that your applications are generating with the expected, reference results. This can be done using the LLM-as-a-judge methods described above, using reference answers as your point of comparison. Note that, If your monitoring shows that the nature of your input to the LLM has changed, it may mean that you need to update your human-approved reference answers.

# **Agent-Level Monitoring**

The output of an agent may vary depending on many different factors: changes in underlying LLMs, changes in the user inputs, changes in connected retrieval services, changes in tools, etc. Just as the power of agents is their open-endedness and the diversity of systems that they can integrate, so too is this diversity a source of

difficulty when trying to monitor agents' performance and diagnose any issue.

At the heart of agent-level monitoring is measuring if each step in an agent's execution is driving the ultimate task it is meant to accomplish.<sup>2</sup> Is the retriever choosing the right documents *for that task*? Is the agent choosing the right tool *for that task*? Is the final output written in a tone that is appropriate *for that task*?

Whether you should measure the quality of every task completion or just a sample depends on several factors:

- 1. The criticality or riskiness of the application. Some applications are critical to core business processes and others are high risk (and some are both!). In these cases, you may choose to monitor all completions, rather than a sample.
- 2. The variability of the performance. If the agent shows that it can complete the task in a consistent way, then you may content yourself to measure the performance of only a sample.
- 3. The volume of task completion and cost of the monitoring. As described above, certain evaluation methods can be costly in and of themselves. Some agents may be completing many thousands of tasks per day. You will need to use the cost-measuring techniques in Chapter 3 to fully understand and capture the costs of these processes to ensure that you have the budget for them.

# Agent Self-Monitoring

You will recall that in the previous sections of this chapter, we described LLM-as-a-judge methods of quality assessment as being essentially agentic applications themselves. You also recall in Chapter 2 that we said that agents could call other agents as tools. So, agents can also use quality assessment agents to monitor their own outputs. This can and should be a core part of how agents iteratively improve their responses, using assessment methods provided to them to check their responses and to attempt to improve their own performance.

102

<sup>2</sup> We see again here yet another reason why it is important to design agents with a relatively narrow scope of action. When the scope is broader, it is more difficult to specify what good task completion looks like.

This iterative self-improvement should be fully logged so that human experts can then review how the agents are improving and to identify common areas for improvement. For example, if the developers find that an agent often calls the wrong tool for a particular task and must go through a quality-improvement loop to identify and then correct the error, then the developers can improve their agent design to better direct the agent in the right direction from the start, for example by improving the schema of the tool or the prompt instructing the agent what to do. As always, reference answers and human expert monitoring remain essential to ensure high quality and adherence to ethical and regulatory requirements.

#### What are the alternatives to an LLM Mesh?

When evaluating LLM applications, organizations can choose between fully distributed evaluation, centralized monitoring, or an LLM Mesh. A fully distributed approach allows teams to build custom evaluation systems tailored to each agent or application, making it easy to implement and scale alongside development. However, this method lacks enterprise-wide quality standards, making it difficult to enforce consistency and implement standardized diagnostic procedures. Additionally, evaluation results often remain siloed, limiting the ability to drive adaptive behaviors across the organization.

Centralized monitoring, on the other hand, provides a common evaluation framework by capturing logs and overseeing performance from a single application. While this aligns with traditional enterprise monitoring paradigms and ensures standardization, it tends to focus on monolithic applications rather than reusable components, limiting flexibility. In contrast, an LLM Mesh offers evaluation as a shared service, enabling applications, agents, and multi-agent systems to dynamically leverage a standardized yet adaptable evaluation framework. Although it may require an initial shift in enterprise architecture, an LLM Mesh ultimately promotes efficiency, reusability, and scalability, making it a strong choice for optimizing LLM evaluation.

# Measuring and Monitoring the Speed of Agentic Applications

Like any other service, it is important to monitor the speed and responsiveness of the agentic applications that you will build. As discussed in Chapter 2, an agentic application involves many API calls to LLM services and other objects within the LLM Mesh. The speed and responsiveness of the agentic application will depend on the collective speed and responsiveness of these various services. Thankfully, monitoring the speed of API services is a well-established DevOps practice, and those methods apply equally well here. Thus, this guide will not develop those concepts fully, but rather only mention a few aspects that are specific to agentic applications.

The speed and responsiveness of agentic applications are measured using specific metrics. Two important metrics to consider are *latency* and *throughput*. Latency measures how long it takes for an LLM service to respond to an input. Throughput measures how many requests it can process or how much output it can produce in a given time span.

Inference latency is usually measured in time to first token (TTFT) and time per output token (TPOT). Together with the number of output tokens, these metrics can be used to calculate a global metric, total generation time, which measures how long it takes to provide a response from the moment the input is received until the response generation ends.

Inference throughput is measured primarily with tokens per second. This metric most often takes into account only output tokens and an LLM Mesh should specify whether this is the case or not.

An LLM Mesh should provide observability of these performance metrics to ensure that they are meeting the requirements of an application. These requirements should be documented in the LLM Mesh on a per-application basis as well. For example, an internal chatbot supporting customer support agents as they work in real-time with customers will have different and higher performance requirements than an application that runs silently in the background analyzing contracts with suppliers.

## Capturing and Optimizing the Costs of Performance Evaluation

As you have probably already realized, some of the LLM-based evaluation methods<sup>3</sup> can result in a lot of traffic to your LLM services and thus will generate costs. It is important to capture and optimize these costs as you would for any agentic application, as described in Chapter 3.

It is important to be intelligent about the frequency with which you run these evaluations so that you can balance the cost of the evaluation itself against any expected improvement. Furthermore, you will need to decide on a sampling strategy for monitoring the responses of production-deployed applications. The volume of their responses is likely to be too great to justify monitoring every response, though this needs to be balanced against the sensitivity of the application.

If you have a very sensitive application that requires the evaluation of every response using a computationally expensive evaluation method, you may find that the combination of requirements makes the application economically unviable, even if it is technically possible to build. In such a situation, an LLM Mesh can help by allowing you to test different evaluation methods to see if there is one that is sufficiently cost-efficient while maintaining the required performance.

# Conclusion

In this and the previous chapters, we have understood how we can quantify and reduce costs, as well as how to measure performance in terms of both the quality of the response and the overall speed of the service. By measuring these three dimensions, your developers will be able to experiment with different approaches to find the right combination of speed, quality, and cost. Once again, there is no silver bullet and no single best practice. But an LLM Mesh can make it far more efficient for teams to test and build high-performing, cost-optimized applications by ensuring that they focus on the development of the application itself and not the supporting evaluation and monitoring services since those are provided by the LLM Mesh.

Conclusion 105

<sup>3</sup> The statistical methods are far less computationally intensive, their costs are essentially negligible.

In the following two chapters, we will discuss additional shared services that you should also include when implementing an LLM Mesh. These are services to ensure the safety and appropriateness of the content generated by your agentic applications. You'll see how an LLM Mesh can help by providing the necessary mechanisms to pass those tests with flying colors.

## **About the Author**

As Head of AI Strategy at Dataiku, **Kurt Muehmel** brings Dataiku's vision of Everyday AI to industry analysts and media worldwide. He advises Dataiku's C-Suite on market and technology trends, ensuring that they maintain their position as pioneers. Kurt is a creative and analytical executive with 15+ years of experience and foundational expertise in the Enterprise AI space and, more broadly, B2B SaaS go-to-market strategy and tactics. He's focused on building a future where the most powerful technologies serve the needs of people and businesses.