

From language models to successful products

# LLMs IN PRODUCTION

Christopher Brousseau  
Matthew Sharp



MEAP



MANNING

From language models to successful products

# LLMs IN PRODUCTION

Christopher Brousseau  
Matthew Sharp



MEAP

 MANNING

# LLMs in Production

1. [welcome](#)
2. [1 Word's Awakening: Why Large Language Models Have Captured](#)
3. [2 Large Language Models: A deep dive into language modeling](#)
4. [3 Large Language Model Operations: Building a platform for LLM](#)
5. [4 Data Engineering for Large Language Models: Setting up for suc](#)
6. [5 Training Large Language Models: How to generate the generator](#)
7. [6 Large Language Models in Production: A Practical Guide](#)

# welcome

Thank you for purchasing the MEAP for *LLMs in Production*. This book is a labor of love born from our excitement and experience working with them professionally, and we're very grateful that you're interested in sharing this journey with us.

Machine Learning is everywhere, and it's becoming increasingly difficult to keep up. At the same time, the fundamentals of actually taking a ML project from idea to working product haven't really changed. When asked about when we were going to learn how to productionize the models so that people could use them, one of our deep learning teachers in college actually said, "That isn't academically interesting, so we won't be doing it." We posit that it's not only academically interesting, it's necessary for any Data Scientist, Data Engineer, and Machine Learning Engineer to have a complete picture to boost their own work.

We aim to provide you the reader with the tools to see the larger picture in production, from as many perspectives as possible to really help you at every step of the ML lifecycle. As such, we will need your help to make that happen, because we don't have your perspective. Give us feedback and help make this book something that can immediately help take any ML professional to the next level. This isn't a history or a linguistics book, but it also isn't a computer science or math book. We want you to get as much out of this book as possible. We have quite a lot of ground to cover, so be prepared to research concepts that interest you further. You'll want to have some experience with Python, and any experience you have with ML as a field will help.

This book is divided into several parts, the first of which focuses on ML foundations. We'll see different types of language models and how to solve for actual language and how to perform traditional MLOps.

The second part will go more in-depth on actual LLMs, how you can help them be ready for production from how you acquire your dataset to how to

train them and compensate for their size.

The third and final part of the book will go through several projects building and using LLMs with various goals on various hardware. We trust that it will be useful to see the application of the concepts we've gone over.

Please reach out to us with your thoughts, either through comments in the [liveBook Discussion forum](#) or in the GitHub repo accompanying the book.

— Christopher Brousseau and Matthew Sharp

**In this book**

[welcome](#) [1 Word's Awakening: Why Large Language Models Have Captured Attention](#) [2 Large Language Models: A deep dive into language modeling](#) [3 Large Language Model Operations: Building a platform for LLMs](#) [4 Data Engineering for Large Language Models: Setting up for success](#) [5 Training Large Language Models: How to generate the generator](#) [6 Large Language Models in Production: A Practical Guide](#)

# 1 Word's Awakening: Why Large Language Models Have Captured Attention

## This chapter covers

- What Large Language Models are and what they can and cannot do
- When you should deploy your own Large Language Models and when should not
- Large Language Model myths and the truths that lie behind them

"Any sufficiently advanced technology is indistinguishable from magic."

- *Arthur C. Clarke*

The year is 1450. A sleepy corner of Mainz, Germany, unknowingly stands on the precipice of a monumental era. In Humbrechthof, a nondescript workshop shrouded in the town's shadows pulsates with anticipation. It is here that Johannes Gutenberg, a goldsmith and innovator, sweats and labors amidst the scents of oil, metal, and determination, silently birthing a revolution. In the late hours of the night, the peace is broken intermittently by the rhythmic hammering of metal on metal. In the lamp-lit heart of the workshop stands Gutenberg's decade-long labor of love—a contraption unparalleled in design and purpose.

This is no ordinary invention. Craftsmanship and creativity transform an assortment of moveable metal types, individually-cast characters born painstakingly into a matrix. The flickering light dances off the metallic insignias. The air pulsates with the anticipation of a breakthrough and the heady sweetness of oil-based ink, an innovation from Gutenberg himself. In the stillness of the moment, the master printer squares his shoulders and with unparalleled finesse, lays down a crisp sheet of parchment beneath the ink-loaded matrix and allows his invention to press firmly and stamp finely print

onto the page. The room adjusts to the symphony of silence, bated breaths hanging heavily in the air. As the press is lifted, it creaks under its own weight, each screech akin to a war cry announcing an exciting new world.

With a flurry of motion, Gutenberg pulls from the press the first printed page and slams it flat onto the wooden table carefully examining each character which are all as bold and magnificent as the creator's vision. The room drinks in the sight, absolutely spellbound. A mere sheet of parchment has become a testament to transformation. As the night gives way to day, he looks upon his workshop with invigorated pride. His legacy is born, echoing in the annals of history and forever changing the way information would take wings. Johannes Gutenberg, now the man of the millennium, emerges from the shadows, an inventor who dared to dream. His name synonymous with the Printing Press which is not just a groundbreaking invention, but the catalyst of a modern world.

As news of Gutenberg's achievement begins to flutter across the continent, scholars from vast disciplines are yet to appreciate the extraordinary tool at their disposal. Knowledge and learning, once coveted treasures, are now within the reach of the common man. There were varied and mixed opinions surrounding that newfound access.

“In our time, thanks to the talent and industry of those from the Rhine, books have emerged in lavish numbers. A book that once would’ve belonged only to the rich - nay, to a king - can now be seen under a modest roof. [...] There is nothing nowadays that our children [...] fail to know.” - Sebastian Brant

“Scholarly effort is in decline everywhere as never before. Indeed, cleverness is shunned at home and abroad. What does reading offer to pupils except tears? It is rare, worthless when it is offered for sale, and devoid of wit.” - Egbert of Liege

People have had various opinions on books throughout history. One thing we can agree on, living in a time when virtual printing presses exist and books are ubiquitous: the printing press changed history. While we weren't actually there when Gutenberg printed the first page using his printing press, we have watched many play with large language models (LLMs) for the first time. The astonishment on their face as they see it respond to their first prompt.

The excitement they have when challenging it with a difficult question only to see it respond as if it was an expert in the field. The light bulb moment when they realize they can use this to simplify their life or make themselves wealthy. I imagine this wave of emotions were, but a fraction of those felt by Johannes Gutenberg. Being able to rapidly generate text and accelerate communication has always been valuable.

## **1.1 Large Language Models accelerating communication**

Every job has some level of communication. Oftentimes this communication is shallow, bureaucratic, or political. I've often warned students and mentees that every job has its paperwork. Something that used to be a passion can easily be killed by the day to day tedium and menial work that comes with it when it becomes a job. In fact, when we talk about our professions we often talk them up, trying to improve our social standing, so you'll rarely get the full truth. You won't hear about the boring parts and the day-to-day grind gets conveniently forgotten.

However, envision a world where we reduce the burden of monotonous work. A place where police officers no longer have to waste hours of each day filling out reports and could instead devote that time to community outreach programs. Or a world where teachers, no longer working late into the night grading homework and preparing lesson plans, instead being able to think about and prepare customized lessons for individual students. Or even a world where lawyers would no longer be stuck combing through legal documents for days, instead being free to take on charity cases for causes that inspire them? When the communication burden, the paperwork burden and the accounting burden, are taken away, the job becomes more akin to what we sell it as.

For this, LLMs are the most promising technology to come along since, well, the printing press. For starters, they have completely upended the role and relationship between humans and computers, transforming what we believed they were capable of. They have already passed medical exams[\[1\]](#), the bar exam[\[2\]](#), and multiple theory of mind tests. They've passed both Google and



Amazon coding interviews[\[3\]](#). They've gotten scores of at least 1410 out of 1600 on the SAT. One of the most impressive to the authors is that GPT-4 has even passed the Advanced Sommelier exam—which makes us wonder how they got past the practical wine tasting portion. Indeed, their unprecedented accomplishments are coming at breakneck speed and often make us mere mortals feeling a bit queasy and uneasy. What do you do with a technology that seems to be able to do anything?

Passing tests is fun and all, but not exactly all that helpful unless our aim was to build the most expensive cheating machine ever, which we promise there's better use of our time. What LLMs are good at is language, particularly, helping us improve and automate communication. This allows us to transform common bitter experiences into easy enjoyable experiences. For starters, imagine entering your home where you have your very own personal JARVIS, as if stepping into the shoes of Iron Man, an AI-powered assistant that adds an unparalleled dynamic to your routine. While not quite to the same artificial general intelligence (AGI) levels as those portrayed by JARVIS in the Marvel movies, LLMs are powering new user experiences from improving customer support to helping you shop for a loved one's birthday. It knows to ask you about the person, learn about their interests and who they are, find out your budget, and then make specialized recommendations. While many of these assistants are being put to good work, many others are simply just chatbots that users can talk to and entertain themselves—which is important because even our imaginary friends are too busy these days. Jokes aside, these can create amazing experiences allowing you to meet your favorite fictional characters like Harry Potter, Sherlock Holmes, Anakin Skywalker or even Iron Man.

What we're sure many readers are interested in though is programming assistants, because we all know Googling everything is actually one of the worst user experiences. Being able to write a few objectives in plain English and see a copilot write the code for you is exhilarating. I've personally used these tools to help me remember syntax, simplify and clean code, write tests, and learn a new programming language.

Video Gaming is another interesting field where we can expect LLMs to create a lot of innovation. Not only do they help the programmers create the

game, but they are allowing designers to create more immersive experiences. For example, talking to NPCs (non-player characters) will have more depth and intriguing dialogue. Picture games like Animal Crossing or Stardew Valley having near infinite quests and conversations.

Consider other industries, like Education where there just doesn't seem to be enough teachers to go around meaning our kids aren't getting the one on one attention they need. An LLM assistant can help save the teacher time doing manual chores as well as serve as a private tutor for the kids that are struggling. Corporate is looking into LLMs for talk-to-your-data jobs, tasks helping employees understand quarterly reports and data tables, essentially giving everyone their own personal analyst. Sales and Marketing divisions are guaranteed to take advantage of this marvelous innovation, for better or worse. The state of Search Engine Optimization (SEO) will change a lot too since currently it is mostly a game of generating content to hopefully make websites more popular, which is now super easy.

The above list is just a few of the common examples I've seen where companies are interested in using them. People are using them for personal reasons too. Writing music, poetry, and even books, translating languages, summarizing legal documents or emails, and even using them for free therapy—which yes, is an awful idea since they are still dreadful at this. Just a personal preference, but we wouldn't try to save a buck when our sanity is on the line. Of course, this leads us to the fact that people are already using them for darker purposes like cheating, scams, and fake news to skew elections. At this point, the list has become rather large and varied, but we've only begun to scratch the surface of the possible. Really, since LLMs help us with communication often, it's better to think, “What can't they do?” than “What can they do?”.

Or better yet, “What shouldn't they do?” Well, as a technology there are certain restrictions and constraints, for example, LLMs are kind of slow. Of course, slow is a relative term, but responsive times are often measured in seconds not milliseconds. We'll dive deeper into this in Chapter 3, but as an example, we probably won't see them being used in autocomplete tasks anytime soon which require blazingly fast inference in order to be useful. After all, autocomplete needs to be able to predict the word or phrase faster

than someone types. In a similar fashion, LLMs are large complex systems, we don't need them for such a simple problem anyway. Hitting an autocomplete problem with an LLM isn't just hitting the nail with a sledgehammer, it's hitting it with a full on wrecking ball. And just like it's more expensive to rent a wrecking ball than buy a hammer, an LLM will cost you more to operate. There are a lot of similar tasks where we should consider the complexity of the problem we are trying to solve.

There are also many complex problems that are often poorly solved with LLMs such as predicting the future. No, we don't mean with mystic arts, but forecasting problems, acts like predicting the weather or when high tide will hit on the ocean shore. These are actually problems we've solved, but we don't necessarily have good ways to communicate how these have been solved. They are expressed through combinations of math solutions like Fourier transforms and harmonic analysis or through black box ML models. There are many problems that fit into this category, like outlier prediction, calculus, or finding the end of the roll of tape.

You also probably want to avoid using them for highly risky projects. LLMs aren't infallible and make mistakes often. To increase creativity we often allow for a bit of randomness in LLMs which means you can ask an LLM the same question and get different answers. That's risky. You can remove this randomness by doing what's called turning down the temperature, but that might make it useless depending on your needs. For example, you might decide to use an LLM to categorize investment options as good or bad, but do you want it to then make actual investment decisions based on its output? Not without oversight, unless your goal is to create a meme video.

Ultimately an LLM is just a model, it can't be held accountable for losing your money, and really it didn't lose your money you did by choosing to use it. Similar risky problems could include filling out tax forms or getting medical advice. While an LLM could do these things, it won't protect you from heavy penalties in an IRS audit like hiring a certified CPA would. If you take bad medical advice from an LLM, there's no doctor you could sue for malpractice. However, in all of these examples, the LLM could potentially largely help the practitioner better perform their job roles both reducing errors and improving speed.

## When to use an LLM

Use them for:

- Generating content
- Question and answering services
- ChatBots and AI assistants
- Text-to-something problems (diffusion, txt2img, txt23d, txt2vid, etc.)
- Talk-to-your-data applications
- Anything that involves communication

Avoid using them for:

- Latency-sensitive workloads
- Simple projects
- Problems we don't solve with words but with math or algorithm - forecasting, outlier prediction, calculus, etc.
- Critical evaluations
- High-risk projects

Language is not just a medium people use to communicate. It is the tool that made humans apex predators and gives every individual self-definition to their community. Every aspect of human existence, from arguing with your parents to graduating from college to reading this book is pervaded by our language. Language models are learning to harness one of the fundamental aspects of being human and have the ability, when used responsibly, to help us with each and every one of those tasks. They have the potential to unlock dimensions of understanding both of ourselves and others if we responsibly teach them how.

LLMs have captured the world's attention since their potential allows imaginations to run wild. LLMs promise so much, but... where are all these solutions? Where are the video games that give us immersive experiences? Why don't our kids have personal AI tutors yet? Why am I not Iron Man with my own personal assistant yet? These are the deep and profound questions that motivated us to write this book. Particularly that last one, it keeps me up at night. So while LLMs can do amazing things, not enough people know how to actually turn them into a product and that's what we aim to share in

this book.



This isn't just a Machine Learning Operations book. There are a lot of gotchas and pitfalls involved with making an LLM work in production because LLMs just don't work like traditional software solutions. To turn an LLM into a product that can interact coherently with your users will require an entire team and diverse set of skills. Depending on your use case you may need to train or finetune and then deploy your own model, or you may just need to access one from a vendor through an API.

Regardless of which LLM you use, if you want to take full advantage of the technology and build the best user experience, you will need to understand how they work. Not just on the math/tech side either, but on the soft side for how to make it a good experience for your users. In this book we'll be covering everything you need to make LLMs work in production. We'll talk about the best tools and infrastructure, how to maximize their utility with

prompt engineering and other best practices like controlling costs. LLMs could be one step towards a greater equality, so if you are thinking, “I don’t feel like the person this book is for,” please reconsider. This book is for the whole team and anyone who will be interacting with LLMs in the future.

We’re going to hit on a practical level everything that you’ll need for collecting and creating a dataset, training or finetuning an LLM on consumer or industrial hardware, and deploying that model in various ways for customers to interact with. While we aren’t going to cover too much theory, we will be covering the process from end to end with real-world examples. At the end of this book you will know how to deploy LLMs with some viable experience to back it up.

## **1.2 Navigating between the Build and Buy decision with Large Language Models**

If you bought this book, you are likely already convinced of the overwhelming potential LLMs can have in your life and in your organization. Buying this book then is the first step to turning your dreams into a reality, because none of it is possible until we know how to put these models into production. After all, if you talk to any entrepreneur or investor out there they will tell you good ideas are a dime a dozen, what matters is execution and actually manifesting those ideas. What we need to do is get these models into production, where they are readily available to do actual work for you.

There’s no getting around it and no need to sugar coat it either, deploying LLMs into production is hard. Often anything worth pursuing is. In this book we aim to teach you everything you need to know to do it as well as give you some practical hands-on experience. But because it is so hard, it is mighty tempting to take a shortcut. Large corporations like OpenAI and Google have some great offerings of models to choose from, why not just buy them? Let’s start by considering what they offer and if this may be a good choice, and then we’ll take a look at the other side of the coin where these offerings tend to fall flat.

### **1.2.1 Buying: the beaten path**

There are many great reasons to simply just buy access to an LLM. First and foremost is the speed and flexibility accessing an API provides. Working with an API is an incredibly easy and cheap way to build a prototype and get your hands dirty quickly. In fact, so easy you can see in listing 1.1 that it only takes a few lines of Python code to start connecting to OpenAI's API and start using LLMs. Sure, there's a lot that's possible, but it would be a bad idea to heavily invest in LLMs only to find out they happen to fail in your specific domain. Working with an API allows you to fail fast. Building a prototype application to prove the concept and launching it with an API is a great place to get started.

**Listing 1.1 A simple app calling OpenAI's API**

```
import os
import openai

# Load your API key from an environment variable
openai.api_key = os.getenv("OPENAI_API_KEY")

chat_completion = openai.ChatCompletion.create(
    model="gpt-3.5-turbo",
    messages=[{"role": "user", "content": "Hello world"}],
)
```

Oftentimes, buying access to a model can give you a competitive edge. In many cases it could very well be that the best model on the market is built by a company specializing in the domain who are using specialized datasets they have spent a fortune to curate. While you could try to compete and build your own, it may better serve your purposes to simply buy access to their model instead. Ultimately, whoever has the better domain specific data to finetune on is likely going to win, and that might not be you if this is a side project for your company. Curating data can be expensive after all. It can save you a lot of work to go ahead and buy it.

Which leads to the next point, buying is a quick way to access expertise and support. For example, OpenAI has spent a lot of time making their models safe with plenty of filtering and controls to prevent the misuse of their LLMs. They've already encountered and covered a lot of the edge cases, so you don't have to. Buying access to their model also gives you access to the system they've built around it.

Not to mention, the LLM itself is only half the problem in deploying them to production. There's still an entire application you need to build on top of it. Sometimes buying OpenAI's model has thrived over its competitors in not a small part due to their UX and some tricks like making the tokens look like they're being typed. We'll take you through how you can start solving for the UX in your use case, along with some ways you can prototype to give you a major head start in this area.

## **1.2.2 Building: the path less traveled**

Using an API is easy, and in most cases, likely the best choice. However, there are many reasons why you should aim to own this technology and learn how to deploy it yourself instead. While this path might be harder, we'll teach you how to do it. Let's dive into several of those reasons starting with the most obvious: Control.

### **Control**

One of the first companies to truly adopt LLMs as a core technology is a small video game company called Latitude. Latitude specializes in Dungeon and Dragons like role playing games utilizing LLM chatbots, and they have faced challenges when working with them. This shouldn't come off as criticizing this company for their missteps, as they contributed to our collective learning experience and were pioneers forging a new path. Nonetheless, their story is a captivating and intriguing one, like a train wreck that we personally couldn't help but keep watching.

Latitude's first release was a game called AI Dungeon. At inception, it utilized OpenAI's GPT2 to create an interactive and dynamic storytelling experience. It quickly garnered a large gathering of players, who of course started to use it inappropriately. When OpenAI gave Latitude access to GPT3 it promised an upgrade to the gaming experience, instead, what it got was a nightmare.[\[4\]](#)

You see, with GPT3 they added Reinforcement Learning from Human Feedback (RLHF) which greatly helps improve functionality, but this also meant OpenAI contractors were now looking at the prompts. That's the



human feedback part. And these workers weren't too thrilled to read the filth the game was creating. OpenAI's reps were quick to give Latitude an ultimatum. Either they needed to start censoring the players or they'd remove their access to the model—which would have essentially killed the game and the company. With no other option, they quickly added some filters but the filtering system was too much a band-aid, a buggy and glitchy mess. Players were upset at how bad the system was and unnerved to realize Latitude's developers were reading their stories, completely oblivious to the fact that OpenAI was already doing such. It was a PR nightmare. And it wasn't over.

OpenAI decided the game studio wasn't doing enough, stonewalled, they were forced to increase their safeguards, and so they started banning players. Here's the twist, the reason so many of these stories turned to smut, was because the model had a preference for erotica. It would often unexpectedly transform harmless storylines into inappropriately risqué situations causing the player to be ejected and barred from the game. OpenAI was acting the paragon of purity, but it was their model that was the problem. Which led to one of the most ironic and unjust problems in gaming history: players were getting banned for what the game did.

So there they were, a young game studio just trying to make a fun game stuck between upset customers and a tech giant that pushed all the blame and responsibility onto them. If the company had more control over the technology they could have gone after a real solution, like fixing the model. Instead of having to throw makeup on a pig.

In this example, control may come off as your ability to finetune your model and OpenAI now offers finetuning capabilities, but there are many fine-grained decisions that are still lost by using a service instead of rolling your own solution. For example, what training methodologies are used, what regions the model is deployed too, or what infrastructure it runs on. Control is also important for any customer or internal-facing tool. You don't want to have a code generator accidentally output code that's copyrighted or that creates a legal situation for your company. You also don't want your customer-facing LLM to output factually incorrect information about your company or its processes.

Control is your ability to direct and manage the operations, processes, and

resources in a way that aligns with your goals, objectives, and values. If a model ends up becoming central to your product offering and the vendor unexpectedly raises their prices there's little you can do but pay it. If the vendor decides their model should give more liberal or conservative answers that no longer align with your values, you are just as stuck.

The more central a technology is to your business plan, the more important it is to control it. This is why McDonald's owns the real estate for its franchises and why Google, Microsoft, and Amazon all own their own cloud networks. Or even why so many entrepreneurs build online stores through Shopify versus just using other platforms like Etsy or Amazon Marketplace. Ultimately control is the first thing that's lost when you buy someone else's product. Keeping it will give you more options to solve future problems and will also give you a competitive edge.

## **Competitive edge**

One of the most valuable aspects of deploying your own models is the competitive edge it gives you over your competition. Customization - train the model to be the best at one thing. For example, after the release of Bidirectional Encoder Representations from Transformers (BERT) in 2017, which is a transformer model architecture you could use to train your own model, there was a surge of researchers and businesses testing this newfound technology on their own data to worldwide success. At the time of writing, if you search the Hugging Face Hub for "BERT," there are more than 13.7k models returned, all that people trained individually for their own purposes aiming to be the best model for their task.

One of my personal experiences in this area was training SlovenBERTcina. After aggregating the largest (at-the-time) monolingual Slovak language dataset by scraping the Slovak National Corpus with permission, along with a whole bunch of other resources like the OSCAR project and the Europarl corpus. It never set any computational records, and has never appeared in any model reviews or generated partnerships for the company I worked for. It did, however, outperform every other model on the market on the tasks it trained on.

Chances are, neither you nor your company need AGI to generate relevant insights from your data, and in fact if you invented an actual self-aware AGI and planned to only ever use it to crunch some numbers, analyze data and generate visuals for PowerPoint slides once a week, that would definitely be reason enough for the AGI to eradicate humans. More than likely, you need exactly what I did when I made SlovenBERTcina, a large language model that performs the two to three tasks you need better than any other model on the market and doesn't also share your data with Microsoft or other potential competitors. While some data is required to keep secret for security or legal reasons, a lot of data should be guarded simply because they are trade secrets.

There are hundreds of open source LLMs both for general intelligence, and for foundational expertise on a specific task. We'll hit some of our favorites in Chapter 4. Taking one of these open source alternatives and training it on your data to create a model that is the best in the world at that task will ensure you have a competitive edge in your market. It will also allow you to deploy the model your way and integrate it into your system to make the most impact.

## **Integrate anywhere**

Let's say you want to deploy an LLM as part of a choose your own adventure style game that uses a device's GPS location to determine story plots. You know your users are often going to go on adventures into the mountains, out at sea, and generally to locations where they are likely to experience poor service and lack of internet access. Hitting an API just isn't going to work. Now, don't get me wrong, deploying LLMs onto edge devices like in this scenario is still an exploratory subject, but it is possible, and we will be showing you how in chapter 10. Relying upon an API service is just not going to work for immersive experiences.

Similarly, using third party LLM and hitting an API adds integration and latency issues, requiring you to send data over the wire and wait for a response. APIs are great, but they are always slow and not always reliable. When latency is important to a project it's much better to serve the service in-house. The previous section on Competitive Edge discussed two projects with edge computing as a priority, however many more exist. LLAMA.cpp and

ALPACA.cpp are two of the first of such projects, and this space is innovating quicker than any others. Quantization into 4-bit, Low-Rank Adaptation, and Parameter Efficient Finetuning are all methodologies recently created just to meet these needs, and we'll be going over each of these starting in Chapter 3.

When my team first started integrating with ChatGPT's API, it was both an awe-inspiring and humbling experience. Awe-inspiring because it allowed us to quickly build some valuable tools. Humbling because as one engineer joked to me, "When you hit the end point you will get 503 errors, sometimes you get a text response as if the model was generating text, but I think that's a bug." Serving an LLM in a production environment, trying to meet the needs of so many clients, is no easy feat. However, deploying a model that's integrated into your system allows you more control of the process affording higher availability and maintainability than you can currently find on the market. This of course also allows you to better control costs.

## **Costs**

Considering costs is always important because it plays a pivotal role in making informed decisions and ensuring the financial health of a project or an organization. It helps you manage budgets efficiently and make sure that resources are allocated appropriately. Keeping costs under control allows you to maintain the viability and sustainability of your endeavors in the long run.

Additionally, considering costs is crucial for risk management. When you understand the different cost aspects, you can identify potential risks and exert better control over them. This way, you can avoid unnecessary expenditures and ensure that your projects are more resilient to unexpected changes in the market or industry.

Finally, cost considerations are important for maintaining transparency and accountability. By monitoring and disclosing costs, organizations demonstrate their commitment to ethical and efficient operations to stakeholders, clients, and employees. This transparency can improve an organization's reputation and help build trust.

All of these apply as you consider building versus buying LLMs. It may seem immediately less costly to buy, as the most costly service widely used on the market currently is only \$20 USD per month. Compared to an EC2 instance on AWS, just running that same model for inference (not even training) could run you up a bill for about \$250k USD per year. This is where building has done its quickest innovation, however. If all you need is an LLM for a proof of concept, any of the projects mentioned in the Competitive Edge section will allow you to create a demo for only the cost of electricity to run on the computer you are demoing on, and spell out training easily enough to allow for significantly reduced costs to train a model on your own data, as low as \$100 (yes, that's the real number) for a model with 20 billion parameters. Another benefit is knowing that if you build your own, your cost will never go up, like it very much will when paying for a service.

## **Security and privacy**

Consider the following case. You are a military staff member in charge of maintenance for the nuclear warheads in your arsenal. All the documentation is kept in a hefty manual. There's so much information required to outline all the safety requirements and maintenance protocols cadets are known to forget important information despite their best efforts. They often cut the wires before first removing the fuse<sup>[5]</sup>. You decide to fine-tune an LLM model to be a personal assistant, giving directions and helping condense all that information giving soldiers exactly what they need when they need it. It's probably not a good idea to upload those manuals to another company—understatement of the century—you're going to want to train something locally that's kept secure and private.

This scenario may sound farfetched, but when speaking to an expert working in analytics for a police department they echoed this exact concern. Talking with them, they expressed how cool ChatGPT is even having their whole team take a prompt engineering class to better take advantage of it, but lamented that there was no way for his team to use it for their most valuable work—the sort of work that literally saves lives—without exposing sensitive data and conversations. Anyone in similar shoes should be eager to learn how to deploy a model safely and securely.

You don't have to be in the army or police force to handle sensitive data. Every company has important intellectual property and trade secrets that are best to keep a secret. Having worked in the semiconductor, healthcare, and finance industries we can tell you first hand, paranoia and corporate espionage are part of the culture in these industries. Because of this Samsung and other industry players at first locked down ChatGPT preventing employees from using it, later opening it up. Of course, it didn't take long before several Samsung employees leaked confidential source code[\[6\]](#). Because OpenAI uses their user's interactions to improve the model, that code is retained and could have been used to further train the model later on. Meaning with the right prompt injection, anyone could potentially pull the code out of the model. A recent example goes even further, where when any OpenAI model was prompted to repeat a word ad infinitum, it would start regurgitating training data, including all of the personally identifiable information (PII) that had snuck through the cleaning process.

#### **A note on OpenAI's policies**

OpenAI's privacy and usage policies have changed a lot over the course of writing this book. When ChatGPT was first introduced it was done as a demo specifically so OpenAI could collect user interactions and improve the model. They pretty much didn't have a privacy policy and they had disclaimers saying such. As ChatGPT grew and became an actual product this changed as clients wanted more protection. For example, OpenAI changed their policies to better serve their customers and since March 1st of 2023 they no longer use customer API data to improve their models[\[7\]](#). The wording of course indicates it is only data sent through the API. It's best if you ask your lawyers on where your company stands on using it. Regardless, the fact that terms of use have changed so much is just further proof you might want more control in this regard.

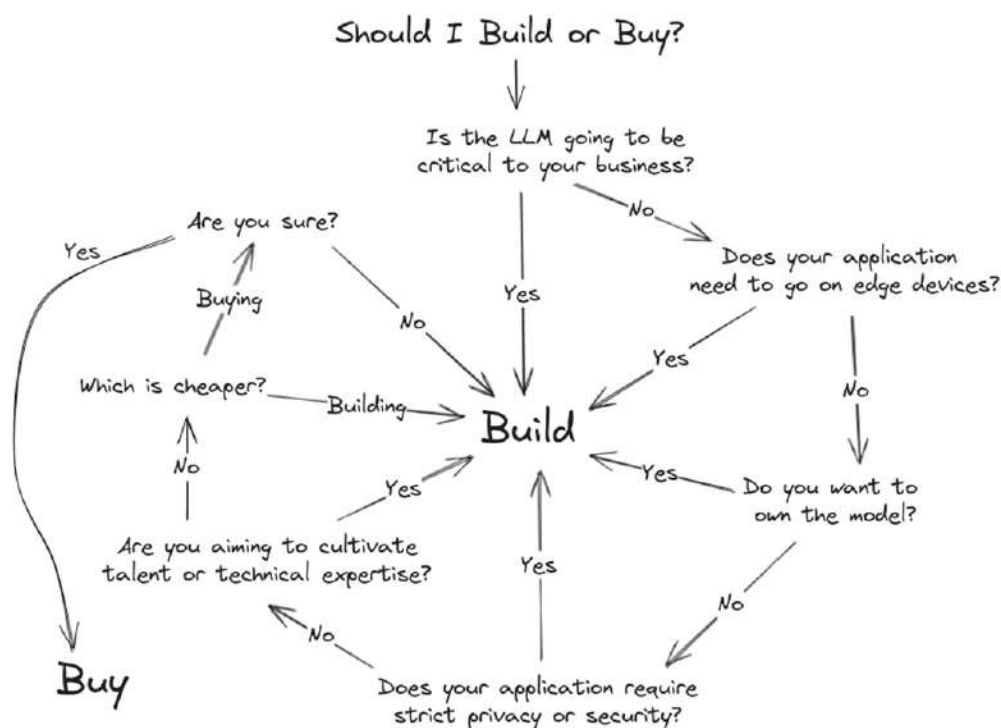
It's not just code that can easily be lost. Business plans, meeting notes, confidential emails, and even potential patent ideas are at risk. We unfortunately know of a few companies who have started sending confidential data to ChatGPT, using that model to clean and extract PII. If this strikes you as potential negligent misuse, you'd be right. This methodology directly exposes customer data, not just to OpenAI, but to any

and all 3rd-party services that they use (including AWS Mechanical Turk, Fiverr, and Freelance workers) to perform the Human Feedback part of RLHF. Don't get us wrong, it's not necessarily a security or privacy issue if you use a third party to do data processing tasks even for sensitive data, but this should be done with high levels of trust and contracts in place.

## Wrapping up

As you can see, there are lots of reasons why a company might want to own and build their own LLMs including having greater control, cutting costs, and meeting security and regulation requirements. Despite this we understand that buying is easy, and building is much more difficult so for many projects it makes sense to just buy, but before you do in figure 1.1 we share a flowchart of questions you should ask yourself first. Even though it's the more difficult path, building can be much more rewarding.

**Figure 1.1** Questions you should ask yourself before making that build vs buy decision.



One last point we think these build versus buy conversations never seem to hone in on enough is, “Por qué no los dos?” Buying gets you all the things

building is bad at: time-to-market, relatively low cost, and ease-of use. Building gets you all the things buying struggles with: privacy, control, and flexibility. Research and prototyping phases could very much benefit from buying a subscription to GPT4 or Databricks for building something quick to help raise funding or get stakeholder buy-in. Production however, often isn't an environment that lends itself to third-party solutions well.

Ultimately, whether or not you plan to build or buy we wrote this book for you. Obviously if you plan to build it, there's going to be a lot more you need to know about, so a majority of this book will be geared to these folks. In fact, we don't need to belabor the point anymore, we're going to teach you how to build in this book, but don't let that stop you from doing the right thing for your company.

### **1.2.3 A word of warning: embrace the future now**

All new technology meets resistance and has critics, but despite this, technologies keep being adopted and progress continues. In business, technology can give a company an unprecedented advantage. There's no shortage of stories of companies failing because they failed to adapt to new technologies. We can learn a lot from their failures.

Borders Books first opened its doors in 1971. After developing a comprehensive inventory management system that included advanced analytic capabilities it skyrocketed to become the second-largest Book retailer in the world, only behind Barnes & Noble. Using this new technology it disrupted the industry, allowing it to easily keep track of tens of thousands of books, opening large stores where patrons could peruse many more books than they could at smaller stores. The analytic capabilities helped it track which books were gaining popularity and gain better insights into their customers allowing it to make better business decisions. It dominated the industry for over two decades.

Borders however, failed to learn from its own history, going bankrupt in 2011, failing to adapt and being disrupted by technology this time; e-commerce. In 2001, instead of building their own platform and online store, they decided to outsource their online sales to Amazon.[\[8\]](#) A decision many



critics would say was akin to giving your competitors the key to your business. While not exactly handing over their secret sauce, it was a decision that gave up their competitive edge.

For the next seven years they turned a blind eye to the growing online sector instead focusing on expanding their physical store presence, buying out competitors and securing a coveted Starbucks deal. When Amazon released the Kindle in 2007, the book retail landscape completely changed. Barnes & Noble having run their own online store quickly pivoted and released the Nook to compete, Borders however, did nothing, or in fact, could do nothing.

By embracing e-commerce through a third party, they failed to develop the in-house expertise required to create a successful online sales strategy, leading to a substantial loss in market share. They eventually launched their own e-reader, Kobo, in late 2010, but it was too late to catch up. Their inability to fully understand and implement e-commerce technology effectively led to massive financial losses, store closures, and ultimately, the company filed for bankruptcy in 2011.

Borders Books is a cautionary tale, but there are hundreds more of similar companies who failed to adopt new technology to their own detriment. With a new technology as impactful as LLMs each company has to decide on which side of the fence they want to be on. Do they delegate implementation and deployment to large FAANG like corporations relegating to just hitting an API, or do they take charge preferring to master the technology and deploy it themselves?

The biggest lesson we hope to impart from this story is that technologies build on top of one another. Ecommerce was built on top of the internet. Failing to build their own online store meant Borders failed to build the in-house technical expertise they needed to stay in the game when the landscape shifted. We see the same things with LLMs today because the companies that are best prepared to utilize them have already gathered expertise in machine learning and data science and have some idea of what they are doing.

We don't have a crystal ball that tells us the future, but many believe that LLMs are a revolutionary new technology like the internet or electricity before it. Learning how to deploy these models, or failing to do so, may very

well be the defining moment for many companies. Not because doing so will make or break their company now, but in the future, when something even more valuable comes along that's built on top of LLMs.

Foraying into this new world of deploying LLMs may be challenging but will help your company build the technical expertise to stay on top of the game. No one really knows where this technology will lead, but learning about this technology will likely be necessary to avoid mistakes like Borders Books'.

There are many great reasons to buy your way to success, but there is at least one prevalent thought that is just absolutely wrong. It's the myth that only large corporations can work in this field because it takes millions of dollars and thousands of GPUs to train these models. Creating this impenetrable moat of cash and resources the little guy can't hope to cross. We'll be talking about this more in the next section, but any company of any size can get started and there's no better time than now to do so.

## **1.3 Debunking myths**

We have all heard from large corporations and the current leaders in LLMs how incredibly difficult it is to train an LLM from scratch and how intense it is to try to finetune them. Whether from OpenAI, BigScience, or Google they discuss large investments and the need for strong data and engineering talent. But how much of this is true and how much of it is just a corporate attempt to create a technical moat?

Most of these barriers start with the premise that you will need to train an LLM from scratch if you hope to solve your problems. Simply put, you don't! Open source models covering many dimensions of language models are constantly being released, so you more than likely don't need to start from scratch. While what they say is true that training LLMs from scratch is supremely difficult, we are still constantly learning about how to do it and are able to more and more automate the repeatable portions. In addition, since this is an active field of research frameworks and libraries are being released or updated daily and will help you start from wherever you currently are. Frameworks like oobabooga's Gradio will help you run LLMs and base models like Falcon 40B will be your starting point. All of it is covered. To

add to this, memos have circulated at large companies addressing the lack of a competitive edge that any organization currently holds over the open source community at large.

A friend once confided in me that, “I really want to get more involved in all this machine learning and data science stuff. It seems to be getting cooler every time I blink an eye. However, it feels like the only way to get involved is to go through a lengthy career change and go work for a FAANG. No, thank you. I’ve done my time at large companies, and they aren’t for me. But I hate feeling like I’m trapped on the outside.” This is the myth that inspired this book. We’re here to equip you with tools and examples to help you stop feeling trapped on the outside. We’ll help you go through the language problems that we’re trying to solve with LLMs, along with machine learning operation strategies to account for the sheer size of the models.

Oddly enough, as many believe they are trapped on the outside, many others believe they can become experts in a weekend. Just get a GPT API key and that’s it, you’re done. This has led to a lot of fervor and hype with a cool new demo popping up on social media every day. But, most of these demos never become actual products and not because people don’t want them.

To understand this, let’s discuss IBM’s Watson the world’s most advanced language model before GPT. Watson is a question and answering machine that went on to crush Jeopardy in 2011 against some of the best human contestants to ever appear on the show, Brad Rutter and Ken Jennings. Rutter the highest earning contestant ever to play the game show and Jennings a player so good he went on to win a whopping 74 times in a row. Despite facing these legends, it wasn’t even close. Watson won in a landslide. Jennings in response to the loss responded with the famous quote, “I, for one, welcome our new computer overlords.”[\[9\]](#)

Watson was the first impressive foray into language modeling and many companies were clamoring to take advantage of its capabilities. Starting in 2013, Watson started being released for commercial use. One of the biggest applications involved many attempts trying to integrate it into healthcare to solve various problems. However, none of these solutions ever really worked the way they needed to and the business never became profitable. By 2022 Watson Health was sold off.

What we find when solving language related problems is that building a prototype is easy, building a functioning product on the other hand is very, *very* difficult. There are just too many nuances to language. Many people wonder what made ChatGPT so explosive? Gaining over a million clients in just five days. Most of the answers I've heard would never satisfy an expert because ChatGPT wasn't much more impressive than GPT3 or other LLMs which had been around for several years already. I heard Sam Altman of OpenAI himself say in an interview they didn't think ChatGPT would get this much attention, they thought that would come with GPT4's release.[\[10\]](#) So why was it explosive? The magic in our opinion, is that it was the first product to truly productionize LLMs. To turn it from a demo into an actual product. It was something anyone could interact with and ask it tough questions, only to be amazed by how well it responded. A demo only has to work once, but the product has to work every time and even when millions of users are showing it to their friends saying, "Check this out!" That magic is exactly what you can hope to learn from reading this book.

We're excited about writing this book. We are excited about the possibilities of bringing this magic to you so you can bring it to the world. LLMs are at the intersection between so many fields such as linguistics, mathematics, computer science and more. While the more you know will help you, to be an expert isn't required. Expertise in any of the individual parts only raises the skill ceiling, not the floor to get in. Consider an expert in physics or music theory, they won't automatically have the skills for music production, but they will be more prepared to quickly learn. LLMs are a communication tool and communicating is a skill just about everyone needs.

Like all other skills, your proximity and willingness to get involved are the two main blockers to knowledge, not a degree or ability to notate—these only shorten your journey towards being heard and understood. If you don't have any experience in this area, it might be good to start by just developing an intuition around what an LLM is and needs by going and contributing to a project like OpenAssistant. If you're a human, that's exactly what they need. By volunteering, you can start understanding exactly what these models train on and why. If you fall anywhere from no knowledge up to being a professional machine learning engineer, we'll be imparting the knowledge necessary to shorten your conversations along with your time to

understanding considerably. If you're not interested in learning the theoretical underpinnings of the subject, we've got plenty of hands-on examples and projects to get your hands dirty.

We've all heard a story by now of LLM hallucinations, but LLMs don't need to be erratic. Companies like Lakera are working daily to improve security, while others like LangChain are making it easier to provide models with pragmatic context which makes them more consistent and less likely to deviate. Techniques such as RLHF and Chain of Thought further allow our models to align themselves with negotiations we've already accepted as people and models should understand from the get-go, such as basic addition or the current date, both of which are conceptually arbitrary. We'll help you increase your model stability from a linguistic perspective, so they'll figure out not just what are the most likely outputs, but the most useful.

Something to consider as you venture further down this path is not just the security of what goes into your model/code, but what comes out. LLMs can sometimes produce outdated, factually incorrect, or even copyrighted or licensed material, depending on what its training data contained. LLMs are unaware of any agreements people make about what is supposed to be a trade secret and what can be shared openly. That is, unless you tell it about those agreements during training or through careful prompting mechanisms during inference. Indeed, the challenges around prompt injection giving inaccurate information primarily arise due to two factors: firstly, users requesting information beyond the model's understanding; and secondly, the model developers not fully predicting how users will interact with the models or the nature of their inquiries. If you had a resource that could help you get a head start on that second problem, it would be pretty close to invaluable, wouldn't it?

Lastly, we don't want to artificially or untruthfully inflate your sense of hope with LLMs. They are resource intensive to train and run. They are hard to understand, and they are harder to get working how you want. They are new and not well-understood. The good news is that these problems are being actively worked on, and we've put in a lot of work finding implementations that are concurrent with writing and actively lessen the burden on you to know everything about the entire deep learning architecture. From

quantization to Kubernetes, we'll help you figure out everything you need to know to do this now with what you have. Maybe we'll inadvertently convince you that it's too much, and you should just purchase from a vendor. Either way, we'll help you every step of the way to help you get the results you need from this magical technology.

## 1.4 Summary

- LLMs are exciting because they work within the same framework (language) as humans
- Society has been built on language, so effective language models have limitless applications such as chatbots, programming assistants, video games, and AI assistants.
- LLMs are excellent at many tasks and can even pass high-ranking medical and law exams
- LLMs are wrecking balls not hammers, and should be avoided for simple problems, problems that require low latency, and problems with high risks.
- Reasons to buy include:
  - Quickly get up and running to conduct research and prototype use cases
  - Easy access to highly optimized production models
  - Access to vendors technical support and system
- Reasons to build include:
  - Getting a competitive edge for your business use case
  - Keeping costs low and transparent
  - Ensuring reliability of the model
  - Keeping your data safe
  - Controlling model output on sensitive or private topics
- There is no technical moat that is preventing you from competing with larger companies since open source frameworks and models provide the building blocks to pave your own path.

[1] Med-PaLM 2 has scored an 86.5% on the MedQA exam.

[2] You can see a whole list of exams passed in OpenAI's GPT-4 paper at <https://cdn.openai.com/papers/gpt-4.pdf>

[3] Google interviewed ChatGPT as a test and it passed.  
<https://www.pcmag.com/news/chatgpt-passes-google-coding-interview-for-level-3-engineer-with-183k-salary>

[4] WIRED, “It began as an AI-fueled dungeon game. Then it got much darker,” Ars Technica, May 08, 2021.  
<https://arstechnica.com/gaming/2021/05/it-began-as-an-ai-fueled-dungeon-game-then-it-got-much-darker/>

[5] M\*A\*S\*H reference for those wondering:  
<https://youtu.be/UcaWQZIPXgQ>

[6] 이코노미스트, “[단독] 우려가 현실로...삼성전자, 챗GPT 빗장 풀자마자 ‘오남용’ 속출,” 이코노미스트, Mar. 30, 2023.  
<https://economist.co.kr/article/view/ecn202303300057?s=31>

[7] See ChatGPT FAQ <https://help.openai.com/en/articles/7232945-how-can-i-use-the-chatgpt-api>

[8] A. Lowrey, “Borders bankruptcy: Done in by its own stupidity, not the Internet.,” Slate Magazine, Jul. 20, 2011.  
<https://slate.com/business/2011/07/borders-bankruptcy-done-in-by-its-own-stupidity-not-the-internet.html>

[9] J. Best, “IBM Watson: The inside story of how the Jeopardy-winning supercomputer was born, and what it wants to do next,” TechRepublic, Sep. 09, 2013. <https://www.techrepublic.com/article/ibm-watson-the-inside-story-of-how-the-jeopardy-winning-supercomputer-was-born-and-what-it-wants-to-do-next/>

[10] “A conversation with OpenAI CEO Sam Altman | Hosted by Elevate,” May 18, 2023 <https://youtu.be/uRIWgbvouEw>.

# 2 Large Language Models: A deep dive into language modeling

## **This chapter covers**

- Linguistic background for understanding meaning and interpretation
- A comparative study on language modeling techniques
- Attention and the transformer architecture
- How Large Language Models both fit into and build upon these histories

The idiom “Once upon a time,” is how we signal to each other that a story is beginning. There isn’t an idiom for productionizing LLMs, so instead this chapter delves into linguistics as it relates to the development of Large Language Models (LLMs), exploring the foundations of semiotics, linguistic features, and the progression of language modeling techniques that have shaped the field of natural language processing (NLP). We will begin by studying the basics of linguistics and its relevance to LLMs in section 2.1, highlighting key concepts such as syntax, semantics, and pragmatics, that form the basis of natural language and play a crucial role in the functioning of LLMs. We will delve into semiotics, the study of signs and symbols, and explore how its principles have informed the design and interpretation of LLMs.

We will then trace the evolution of language modeling techniques in section 2.2, providing an overview of early approaches, including N-grams, Naive Bayes classifiers, and neural network-based methods such as Multi-Layer Perceptrons (MLPs), Recurrent Neural Networks (RNNs), and Long Short-Term Memory (LSTM) networks. We will also discuss the groundbreaking shift to transformer-based models in section 2.3, which have laid the foundation for the emergence of LLMs—which are just really big transformer based models. Finally, we will introduce LLMs in 2.4 and their distinguishing features, discussing how they have built upon and surpassed earlier language modeling techniques to revolutionize the field of Natural Language Processing (NLP).



This is a book about LLMs in production. We firmly believe that if you want to turn an LLM into an actual product, understanding the technology better will improve your results and save you from making costly and time-consuming mistakes. Any engineer can figure out how to lug a big model into production and throw a ton of resources at it to make it run, but that brute force strategy completely misses the lessons people have already learned trying to do the same thing before, which is what we are trying to solve with LLMs in the first place. Having a grasp of these fundamentals will better prepare you for the tricky parts, the gotchas, and the edge cases you are going to run into when working with LLMs. By understanding the context in which LLMs emerged, we can appreciate their transformative impact on NLP and how to enable them to create a myriad of applications.

## 2.1 Language Modeling

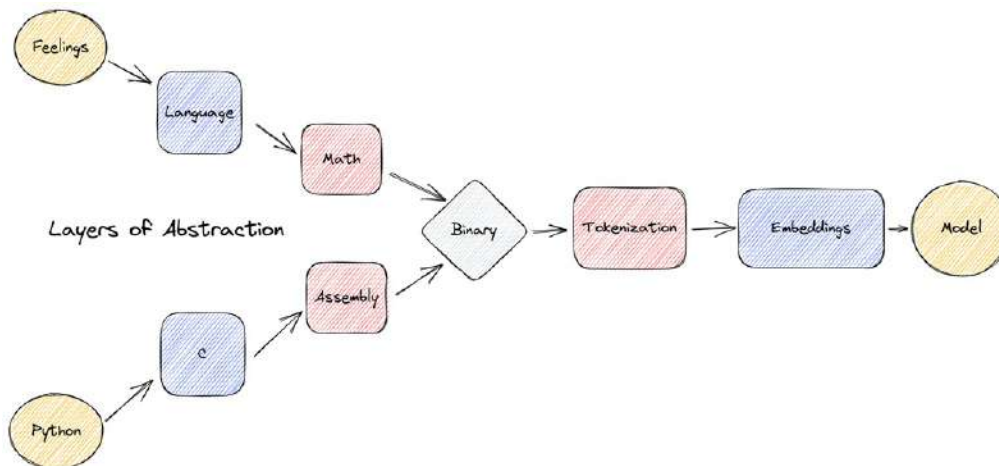
It would be a great disservice to address LLMs in any depth without first addressing language, to begin with. To that end, we will start with a brief but comprehensive overview of language modeling, focusing on the lessons that can help us with modern LLMs. Let's first discuss levels of abstraction as it will help us garner an appreciation for language modeling.

Language, as a concept, is an abstraction of the feelings and thoughts that occur to us in our heads. Feelings come first in the process of generating language, but that's not the only thing we're trying to highlight here. We're also looking at language as being unable to capture the full extent of what we are able to feel, which is why we're calling it an abstraction. It moves away from the source material and loses information. Math is an abstraction of language, focusing on logic and provability, but as any mathematician will tell you, it is a subset of language used to describe and define in an organized and logical way. From math comes another abstraction, the language of binary, a base-2 system of numerical notation consisting of either on or off.

This is not a commentary on usefulness, as binary and math are just as useful as the lower-level aspects of language, nor is it commenting on order, like we said before. With math and binary the order coincidentally lines up with the layer of abstraction. Computers can't do anything on their own and need to take commands to be useful. Binary, unfortunately, ends up taking too long

for humans to communicate important things in, so binary was also abstracted to assembly, a more human-comprehensible language for communicating with computers. This was further abstracted to the high-level assembly language, C, which has been even further abstracted to object-oriented languages like Python or Java (which one doesn't matter, we're just measuring distance from binary). The flow we just discussed is outlined in Figure 2.1.

**Figure 2.1 We compare cognitive layers of abstraction to programming layers of abstraction down to the logical binary abstraction. Python doesn't come from C, nor does it compile into C. Python is, however, another layer of abstraction distant from binary. Similarly, language follows a similar path. Each layer of abstraction creates a potential point of failure. There are also several layers of abstraction to creating a model, each of which are important in seeing the full path from our feelings to a working model.**



This is obviously a reduction, however, it's useful to understand that the feelings you have in your head are the same number of abstractions away from binary, the language the computer actually reads, as the languages most people use to program in. Some people might argue that there are more steps between Python and binary, such as compilers or using assembly to support the C language, and that's true, but there are more steps on the language side too, such as morphology, syntax, logic, dialogue, and agreement.

This can help us understand how difficult the process of getting what we want to be understood by an LLM actually is, and even help us understand language modeling techniques better. The reason we focus on binary here is simply to illustrate that there are a similar number of abstract layers to get

from an idea you have or from one of our code samples to a working model. Like the children's telephone game where participants whisper into each other's ears, each abstraction layer creates a disconnect point or barrier where mistakes can be made.

Figure 2.1 is also meant to illustrate the difficulty in not only creating reliable code and language input, but to draw attention to how important the intermediary abstraction steps like tokenization and embeddings are for the model itself. Even if you have perfectly reliable code and perfectly expressed ideas, the meaning may be fumbled by one of those processes before it ever reaches the LLM.

In this chapter we will try and help you understand what you can do to reduce the risks of these failure points, whether that be on the language, coding, or modeling side. Unfortunately, it's a bit tricky to strike a balance between giving you too much linguistics that doesn't immediately matter for the task at hand versus giving you too much technical knowledge that, while useful, doesn't help you develop an intuition for language modeling as a practice. With this in mind, you should know that linguistics can be traced thousands of years back in our history and there's lots to learn from it. We've included a brief overview for interested readers of how language modeling has progressed over time in Appendix A and encourage you to take a look.

Let's start with our focus on the building blocks that constitute language itself. We expect our readers to have at least attempted language modeling before and to maybe have heard of libraries like PyTorch and Tensorflow, but we do not expect most of our readers to have considered the language side of things before. By understanding the essential features that make up language, we can better appreciate the complexities involved in creating effective language models and how these features interact with one another to form the intricate web of communication that connects us all. In the following section, we will examine the various components of language, such as phonetics, pragmatics, morphology, syntax, and semantics, as well as the role they play in shaping our understanding and usage of languages around the world. Let's take a moment to explore how we currently understand language along with the challenges we face that LLMs are meant to solve.

## 2.1.1 Linguistic Features

Our current understanding of language is that language is made up of at least 5 parts: Phonetics, Syntax, Semantics, Pragmatics, and Morphology. Each of these portions contributes significantly to the overall experience and meaning being ingested by the listener in any conversation. Not all of our communication uses all of these forms, for example, the book you're currently reading is devoid of phonetics, which is one of the reasons why so many people think text messages are unsuited for a more serious or complex conversation. Let's work through what each of these is to figure out how to present them to a language model for a full range of communicative power.

### Phonetics

Probably the easiest for a language model to ingest, phonetics involves the actual sound of the language. This is where accent manifests and deals with the production and perception of speech sounds, with phonology focusing on the way sounds are organized within a particular language system. Similarly to computer vision, while a sound isn't something necessarily easy to deal with as a whole, there's no ambiguity for how to parse, vectorize, or tokenize the actual sound waves. They have a numerical value attached to each part, the crest, the trough, and the slope during each frequency cycle. It is vastly easier than text to be tokenized and processed by a computer while being no less complex. Sound inherently contains more encoded meaning than the text as well, for example, imagine someone saying the words "yeah, right," to you. Could be sarcastic, could be congratulatory, depending on the tone and English isn't even tonal! Phonetics, unfortunately, doesn't have Terabyte-sized datasets commonly associated with it, and performing data acquisition and cleaning on phonetic data, especially on the scale needed to train an LLM is difficult at best. In an alternate world where audio data was more prevalent than text data, and took up a smaller memory footprint, phonetic-based or phonetic-aware LLMs would be much more sophisticated, and creating that world is a solid goal to work towards.

Anticipating this problem, a system created in 1888 called the International Phonetic Alphabet (IPA) has been revised in both the 20th and 21st centuries to be more concise, more consistent, and more clear, which could be a way to

insert phonetic awareness into text data. IPA functions as an internationally standardized version of every language's sound profile. A sound profile is the set of sounds that a language uses, for example in English, we never have the /ʃ/ (she, shirt, sh) next to the /v/ sound. IPA is used to write sounds, rather than writing an alphabet or logograms, as most languages do. For example, you could simply describe how to pronounce the word "cat" using these symbols: /k/, /æ/, and /t/. That's of course a *very* simplified version of it, but for models it doesn't have to be. You can describe tone and aspiration as well. This could be a happy medium between text and speech, capturing some phonetic information. Think of the phrase "what's up?" Your pronunciation and tone can drastically change how you understand that phrase, sometimes sounding like a friendly "wazuuuuup," and other an almost threatening, "'sup," which IPA would fully capture. IPA isn't a perfect solution though, for example, it doesn't solve the problem of replicating tone very well.

Phonetics is listed first here because it's the place that LLMs have been applied to the least out of all the features and therefore have the largest space for improvement. Even modern TTS and Voice cloning models for the most part end up converting the sound to a spectrogram and analyzing that image rather than incorporating any type of phonetic language modeling. This is something to look for as far as research goes in the coming months and years.

## **Syntax**

This is the place where current LLMs are highest-performing, both in parsing syntax from the user and generating its own. Syntax is generally what we think of as grammar and word order, and is the study of how words can combine to form phrases, clauses, and sentences. Syntax is also the first place that language learning programs start to help people acquire new languages, especially based on where you're coming from natively. For example, it is important for a native English speaker learning Turkish to know that the syntax is completely different, and you can often build entire sentences in Turkish that are just one long compound word, whereas in English, we never put our subject and verb together into one word.

Syntax is largely separate from meaning in language, as the famous sentence

from Noam Chomsky the so-called father of syntax demonstrates: “Colorless green ideas sleep furiously.” Everything about that sentence is both grammatically correct and semantically understandable. The problem isn’t that it doesn’t make sense, it’s that it does, and the encoded meanings of those words conflict. This is a reduction, however, you can think of all the times LLMs give nonsense answers as this phenomenon manifesting. Unfortunately for us, the syntax is also where ambiguity is the most commonly found. Consider the sentence, “I saw an old man and woman.” Now answer the question: is the woman also old? This is syntactic ambiguity, where we aren’t sure whether the modifier “old” applies to all people in the following phrase or just the one it immediately precedes. This is less consequential than the fact that semantic and pragmatic ambiguity also show up in syntax. Consider this sentence now, “I saw a man on a hill with a telescope,” and answer the question: Where is the speaker, and what are they doing? Is the speaker on the hill cutting a man in half using a telescope? Likely, you didn’t even consider this option when you read the sentence, because when we interpret syntax, all of our interpretations are at least semantically and pragmatically informed. We know from lived experience that that interpretation isn’t at all likely, so we throw it out immediately, usually without even taking time to process that we’re eliminating it from the pool of probable meanings. Single-modality LLMs will always have this problem, and multi-modal LLMs can (so far) only asymptote towards the solution.

It shouldn’t take any logical leap for why LLMs need to be syntax-aware in order to be high-performing. LLMs that don’t get word order correct or generate nonsense aren’t usually described as “good.” LLMs being syntax-dependent is something that has prompted even Chomsky to call LLMs “stochastic parrots.” In the authors’ opinions, GPT2 in 2018 was when language modeling solved syntax as a completely meaning-independent demonstration, and we’ve been happy to see the more recent attempts to combine the syntax that GPT2 output so well with encoded and entailed meaning, which we’ll get into now.

## **Semantics**

Semantics are the literal encoded meaning of words in utterances, and it

changes at breakneck speed in waves. People automatically optimize semantic meaning, only using words that they consider meaningful in the current language epoch. If you've ever created or used an embedding with language models (word2vec, ELMO, BERT [the E is for Embedding], MUSE, etc.]) you've used a semantic approximation. Words often go through semantic shifts, and while we won't cover all of this topic nor go in-depth, here are some common ones you may already be familiar with: narrowing, a broader meaning to a more specific one, broadening, the inverse of narrowing going from a specific meaning to a broad one, and reinterpretations, going through whole or partial transformations. These shifts do not have some grand logical underpinning. They don't even have to correlate with reality, nor do speakers of a language hardly ever consciously think about the changes as they're happening. That doesn't stop the change from occurring, and in the context of language modeling it doesn't stop us from having to keep up with that change.

Some examples: narrowing includes "deer" which in Old and Middle English just meant any wild animal, even a bear or a cougar, and now means only one kind of forest animal. For broadening we have "dog" which used to refer to only one canine breed from England, and now can be used to refer to any domesticated canine. One fun tangent about dog-broadening is in the *FromSoft* game *Elden Ring*, where because of a limited message system between players, they will use "dog" to refer to anything from a turtle to a giant spider and literally everything in between. For reinterpretation, we can consider "pretty" which used to mean clever or well-crafted, not visually attractive. Another good example is "bikini" which went from referring to a particular atoll, to referring to clothing you might have worn when visiting that atoll to people acting as if the "bi-" was referring to the two-piece structure of the clothing, thus inventing the tankini and monokini. Based on expert research and decades of study, we can think of language as being constantly compared and reevaluated by native language speakers out of which common patterns emerge. The spread of those patterns is closely studied in sociolinguistics and is largely out-of-scope for the current purpose, but can quickly come into scope as localization (l10n) or internationalization (i18n) for LLMs arises as a project requirement. Sociolinguistic phenomena such as prestige can help in designing systems that work well for everyone.

In the context of LLMs, so-called semantic embeddings are vectorized versions of text that attempt to mimic semantic meaning. The most popular way of doing this currently is by tokenizing or assigning an arbitrary number in a dictionary to each subword in an utterance (think prefixes, suffixes, and morphemes generally), applying a continuous language model to increase the dimensionality of each token within the vector so that there's a larger vector representing each index of the tokenized vector, then applying a positional encoding to each of those vectors to capture word order. Each subword ends up being compared to other words in the larger dictionary based on how it's used. We'll show an example of this later. Something to consider when thinking about word embeddings is that they struggle to capture deep encoded meaning of those tokens, and simply adding more dimensions to the embeddings hasn't shown marked improvement. One evidence that embeddings are working in a similar way to humans is that you can apply a distance function to related words and see that they are closer together than unrelated words. This is another area to expect groundbreaking research in the coming years and months for how to capture and represent meaning more completely.

## **Pragmatics**

Sometimes omitted from linguistics, due to its referent being all the non-linguistic context affecting a listener's interpretation and the speaker's decision to express things in a certain way. Pragmatics refers in a large part to dogmas followed in cultures, regions, socio-economic classes, and shared lived experiences played off of to take shortcuts in conversations using entailment.

If I were to say, "A popular celebrity was just taken into the ICU," your pragmatic interpretation based on lived experience might be to assume that a well-beloved person has been badly injured and is now undergoing medical treatment in a well-equipped hospital. You may wonder about which celebrity it is, whether they will have to pay for the medical bills, or if the injury was self-inflicted, also based on your lived experience. None of these things can be inferred directly from the text and its encoded meaning by itself. You would need to know that ICU stands for a larger set of words, and what those words are. You would need to know what a hospital is and why



someone would need to be taken there instead of going there themselves. If any of these feel obvious, good. You live in a society and your pragmatic knowledge of that society overlaps well with the example provided. If I share an example from a less-populated society, “Janka got her grand-night lashings yesterday, she’s gonna get Peter tomorrow” you might be left scratching your head. If you are, realize this probably looks like how a lot of text data ends up looking to an LLM (anthropomorphization acknowledged). For those wondering, this sentence comes from Slovak Easter traditions. There’s a lot of meaning here that will just be missed and go unexplained if you are unaccustomed to these particular traditions as they stand in that culture. I personally have had the pleasure of trying to explain the Easter Bunny and its obsession with eggs to foreign colleagues and enjoyed the satisfaction of looking off my rocker.

In the context of LLMs, we can effectively group all out-of-text context into pragmatics. This means LLMs start without any knowledge of the outside world, and do not gain it during training. They only gain a knowledge of how humans respond to particular pragmatic stimuli. LLMs do not understand social class or race or gender or presidential candidates or anything else that might spark some type of emotion in you based on your life experience. Pragmatics isn’t something that we expect will be able to be directly incorporated into a model at any point, because models cannot live in society, but we have already seen the benefits of incorporating it indirectly through data engineering and curation, prompting mechanisms like RAG, and supervised fine-tuning on instruction datasets. This is a point where we expect great improvements in the future, but would like to emphasize that it’s an asymptoting solution because language is ultimately still an abstraction.

Pragmatic structure gets added whether you mean to add it or not as soon as you acquire the data you are going to train on. You can think of this type of pragmatic structure as bias, not inherently good or bad, but impossible to get rid of. Later down the line you get to pick what types of bias you’d like your data to keep by normalizing and curating, augmenting particular underrepresented points and cutting overrepresented or noisy examples. Instruction datasets show us how you can harness pragmatic structure in your training data to create incredibly useful bias, like biasing your model to answer a question when asked, instead of attempting to categorize the

sentiment of the question.

Pragmatics and context all revolve around entailment. An entailment, is a pragmatic marker within your data, as opposed to the literal text that your dataset contains. For example, let's say you have a model attempting to take an input like "write me a speech about frogs eating soggy socks that doesn't rhyme and where the first letters of each line spell amphibian," and actually follow that instruction. You can immediately tell that this input is asking for a lot. The balance for you as a data engineer would be to make sure that everything the input is asking for is explicitly accounted for in your data. You need to have examples of speeches, examples of what frogs and socks are and how they behave, and examples of acrostic poems. If you don't, the model might be able to understand just from whatever entailments exist in your dataset, but it's pretty up in the air. If you go the extra mile and keep track of entailed vs explicit information and tasks in your dataset, along with data distributions, you'll have examples to answer, "What is the garbage-in resulting in our garbage-out?"

LLMs struggle to pick up on pragmatics, even more so than people, but they do pick up on the things that your average standard deviation of people would. They can even replicate responses from people outside that standard deviation, but pretty inconsistently without the exact right stimulus. That means that it's difficult for a model to give you an expert answer on a problem the average person doesn't know about without providing the correct bias and entailment during training and in the prompt, for example, just including "masterpiece" at the beginning of an image generation prompt will elicit different and usually higher quality generations, but only if that distinction was present in the training set and only if you're asking for an image where "masterpiece" is a compliment. Instruction-based datasets attempt to manufacture those stimuli during training by asking questions and giving instructions that entail representative responses. It is impossible to account for every possible situation in training, and you may inadvertently create new types of responses from your end users by trying to account for everything. After training, you can coax particular information from your model through prompting, which has a skill ceiling based on what your data originally entailed.

## Morphology

Morphology is the study of word structures and how they are formed from smaller units called morphemes. Morphemes are the smallest units of meaning, like the "re-" in "redo" or "relearn." However, not all parts of words are morphemes, such as "ra-" in "ration" or "na-" in "nation," and some can be unexpected like "helico-" as in "helicoid," and "-pter" as in "pterodactyl."

Understanding how words are constructed helps create better language models and parsing algorithms, which are essential for tasks like tokenization. Tokens are somewhere between words and morphemes, they are statistically the most-likely candidates for units of meaning, but don't necessarily correspond to existing morphemes. People do not consciously decide what their units of meaning are going to be, and as such they are often illogical. The effectiveness of a language model can depend on how well it can understand and process these tokens. For instance, in tokenization, a model needs to store a set of dictionaries to convert between words and their corresponding indices. One of these tokens is usually an "/<UNK/>" token, which represents any word that the model does not recognize. If this token is used too frequently, it can hinder the model's performance, either because the model's vocabulary is too small or because the tokenizer is not using the right algorithm for the task.

Consider a scenario where you want to build a code completion model, but you're using a tokenizer that only recognizes words separated by whitespace, like the nltk "punkt" tokenizer. When it encounters the string "def add\_two\_numbers\_together(x, y):," it will pass "[def, /<UNK/>, y]" to the model. This causes the model to lose valuable information, not only because it doesn't recognize the punctuation, but also because the important part of the function's purpose is replaced with an unknown token due to the tokenizer's morphological algorithm. To improve the model's performance, a better understanding of word structure and the appropriate parsing algorithms is needed.

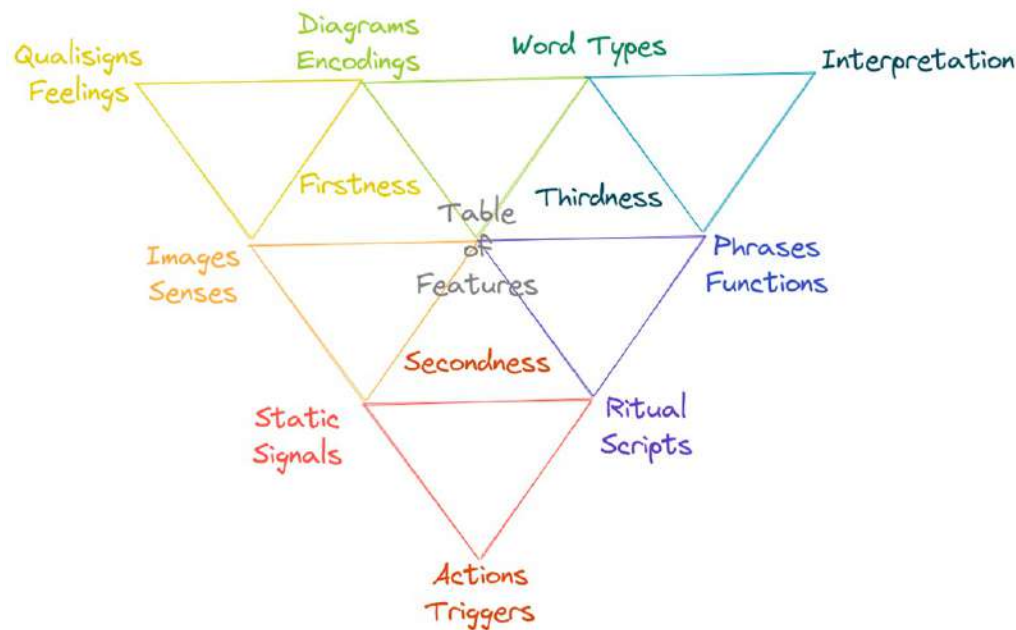
### 2.1.2 Semiotics

After exploring the fundamental features of language and examining their

significance in the context of large language models, it is important to consider the broader perspective of meaning-making and interpretation in human communication. Semiotics, the study of signs and symbols, offers a valuable lens through which we can better understand how people interpret and process language. In the following section, we will delve into the realm of semiotics, examining the relationship between signs, signifiers, and abstractions, as well as how these elements are utilized by LLMs to generate meaningful output. This discussion will provide a deeper understanding of the intricate processes through which LLMs manage to mimic human-like understanding of language, while also shedding light on the challenges and limitations they face in this endeavor. It should be noted that the authors do not believe that mimicking human behavior is necessarily the right answer for LLM improvement, only that mimicry is how the field has evaluated itself so far.

In our introduction to semiotics let's consider Figure 2.2 an adapted Peircean semiotic triangle. These are used to organize base ideas into sequences of firstness, secondness, and thirdness, with firstness being at the top left, secondness at the bottom, and thirdness being at the top right. If you've ever seen a semiotic triangle before, you may be surprised at the number of corners and orientation. To explain, we've turned them upside down to make it slightly easier to read, and because the system is recursive, we're showing how the system can model the entire process and each piece individually simultaneously. While the whole concept of these ideas is very cool, it's outside of the scope of this book to really delve into the philosophy. Instead, we can focus on the cardinal parts of those words (first, second, third) to show the sequence things are processed in.

**Figure 2.2 This is a recursive Peircean semiotic triangle. It's a system of organizing the process of extracting meaning from anything, in our case from language. Each point on the triangle illustrates one of the minimal parts needed to synthesize meaning within whatever the system is being used to describe, so each of these points are minimal units in meaning for language. Firstness, Secondness, and Thirdness are not points on the triangle, more markers for the people versed in Semiotics to be able to orient themselves in this diagram.**



We can also look at each intersection of each of the triangles to gain an idea of why things are presented in the order they are. Feelings can be attached to images and encodings way before they can be attached to words and tables. Ritual and common scripts give a space for interpreted action that's just second nature and doesn't have to be thought about, similar to how most phrases just come together from words without the native speaker needing to perform metacognition about each word individually. All of these eventually lead to an interpretation or a document (a collection of utterances), and in our case, that interpretation should be reached by the LLM. This is why, for example, prompt engineering can boost model efficacy. Foundation LLMs that have trained on millions of examples of ritual scripts can replicate the type of script significantly better when you explicitly tell the model in the prompt which script needs to be followed. Try asking the model to give a step-by-step explanation, maybe prepend your generation with "Let's think about this step-by-step", you will see the model will generate step-by-step scripts based on previous scripts it's seen play out.

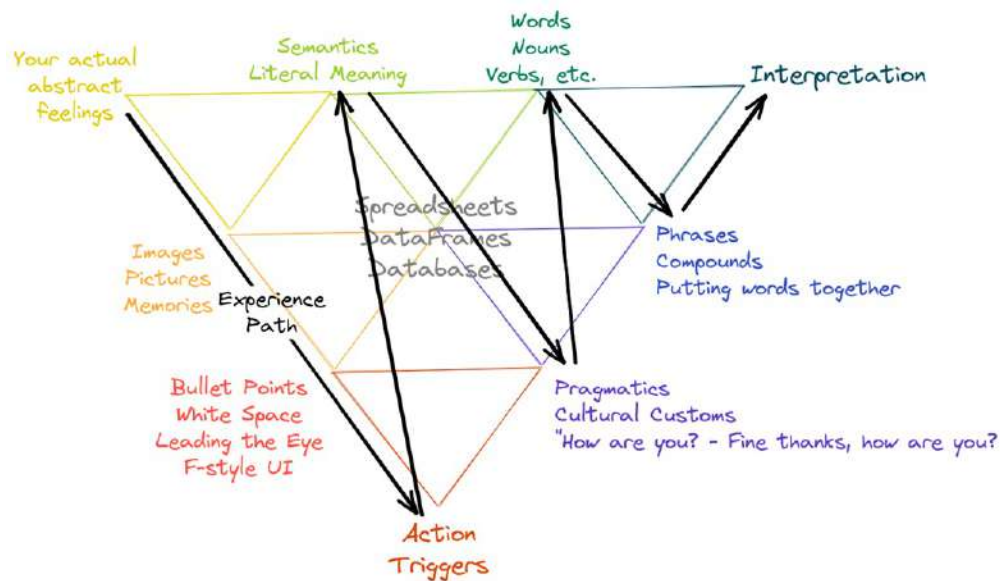
For those interested, there are specific ways of reading these figures and a whole field of semiotics to consider, however, it's not guaranteed that you'll be able to create the best LLMs by understanding the whole thing. Instead of diving really deeply into this, we'll consider the bare minimum that can help you build the best models, UX, and UI for everyone to interact with. For

example, one aspect of the process of creating meaning is recursiveness. When someone is talking to you, and they say something that doesn't make sense (is "meaningless" to you), what do you do? Generally, people will ask one or more clarifying questions to figure out what the meaning is, and the process is started over and over until the meaning is put across. The most state-of-the-art models that are currently on the market do not do this, but they can be made to do it through very purposeful prompting. Many people wouldn't even know to do that without having it pointed out to them. In other words, this is a brief introduction to semiotics. You don't need to be able to give in-depth and accurate coordinate-specific explanations to experts in the semiotic field by the end of this section. The point I'm really trying to push is that this is an organizational system showcasing the minimum number of things you actually need to create a full picture of meaning for another person to interpret. We are not giving the same amount of the same kinds of information to our models during training, but if we did, it would result in a marked improvement in model behavior.

These figures are meant to represent a minimal organizational model, where each of these pieces is essential. Let's consider Figure 2.3 which walks through an example of using a semiotic triangle. Consider Images, Pictures, and Memories and think about what it would be like to try and absorb the knowledge in this book without your eyes to process images, and without orthography (a writing system) to abstract the knowledge. Looking at Bullet Points, etc., how could you read this book without sections, whitespace between letters, and bullet points to show you the order and structure to process information with? Look at Semantics and literal encoded meaning and imagine the book without diagrams and words that didn't have dictionary definitions. Looking at the spreadsheets in the middle, that could be a book without any tables or comparative informational organizers, including these figures. What would it be like trying to read this book without a culture or society that has habits and dogma to use as a lens for our interpretations? All of these points form our ability to interpret information, along with the lens that we end up passing our information through to recognize patterns.

**Figure 2.3 Starting at the top left corner, follow the arrows to see the general order that we use to build our interpretations and extract meaning from things we interact with. Here, we've replaced the descriptive words with some examples of each point. Try to imagine interpreting this diagram without any words, without examples, without the arrows, or even without the pragmatic context**

of knowing what a figure in a book like this is supposed to be for.



So the important question then is: how many of these things do you see LLMs having access to to return meaningful interpretations? Does an LLM have access to feelings or societal rituals? Currently, they do not but think about this as we go through traditional and newer techniques for NLP inference and think about what different models have access to.

### 2.1.3 Multilingual NLP

The last challenge that we need to touch on before we evaluate previous NLP techniques and current-generation LLMs is the foundation of linguistics and the reason LLMs even exist. People have wanted to understand or exploit each other since the first civilizations made contact. These cases have resulted in the need for translators, and this need has only exponentially increased as the global economy has grown and flourished.

It's pretty simple math for business as well. Did you know that there are almost as many native speakers of Bengali as there are native speakers of English? If this is the first time you've heard of the Bengali language, this should hopefully color your perception that there is a valuable market for multilingual models. There are billions of people in the world, but only about a third of one billion of them speak English natively. If your model is anglocentric, like most are, you are missing out on 95% of the people in the

world as customers and users. Spanish and Mandarin Chinese are easy wins in this area, but more people don't even go that far.

There are many more politically-charged examples of calling things the same or different languages that are out of the scope of this book. These are most often because of external factors like government involvement. Keeping these two points in mind—that a monolingual system focusing on English doesn't have the coverage or profit potential as many businesses act like and that the boundaries between languages and dialects are unreliable at best and systematically harmful at worst—it should highlight the dangerous swamp of opinions. Many businesses and research scientists don't even pretend to want to touch this swamp with a 50-foot pole when designing a product or system.

Unfortunately, no easy solutions exist at this time. However, the consideration of these factors can help you as a scientist or engineer (and hopefully an ethical person) to design LLMs that, at the very least, don't exacerbate and negatively contribute to the problems that already exist. The first step in this process is deciding on a directional goal from the beginning of the project, either towards localization (L10n) or internationalization (I18n). Localization is an approach exemplified by Mozilla, which has a different version of their browser available through crowdsourced L10n in over 90 languages with no indications of stopping that effort.

Internationalization is similar, but in the opposite direction, for example, Ikea tries to put as few words as possible in their instructional booklets, opting instead for internationally recognized symbols and pictures to help customers navigate the DIY projects. Deciding at the beginning of the project cuts down the effort required to expand to either solution exponentially. Large enough to switch the perception of translation and formatting from a cost to an investment. In the context of LLMs and their rapid expansion across the public consciousness, it becomes even more important to make that consideration early. Hitting the market with a world-changing technology that automatically disallows most of the world from interacting with it devalues those voices. Having to wait, jeopardizes their economic prospects.

Before continuing, let's take a moment to reflect on what we've discussed so far. We've hit important points in linguistics illustrating concepts for us to consider like understanding that the structure of language is separate from its



meaning. We have demonstrated quite a journey that each of us takes, both personally and as a society, towards having the metacognition to understand and represent language in a coherent way for computers to work with. This understanding will only improve as we deepen our knowledge of cognitive fields, and as we solve for the linguistic features that we encounter. Going along with Figure 2.1, we will now demonstrate the computational path for language modeling that we have followed, and explore how it has and hasn't solved for any of those linguistic features or strived to create meaning. Let's move into evaluating the various techniques for representing a language algorithmically.

## 2.2 Language Modeling Techniques

Having delved into the fundamental features of language, the principles of semiotics, and the ways in which large language models interpret and process linguistic information, we now transition into a more practical realm. In the following section, we will explore the various natural language processing techniques that have been developed and employed to create these powerful language models. By examining the strengths and weaknesses of each approach, we will gain valuable insights into the effectiveness of these techniques in capturing the essence of human language and communication. This knowledge will not only help us appreciate the advancements made in the field of NLP but also enable us to better understand the current limitations of these models and the challenges that lie ahead for future research and development.

Let's take a second to just go over some data processing that will be universal to all language modeling. First, we'll need to decide how we want to break up the words and symbols that we'll be passing into our model, effectively deciding what a token will be in our model. Then, we'll need a way to convert those tokens to numerical values and back again. Then, we'll need to pick how our model will actually process the tokenized inputs. Each of the following techniques will build upon the previous techniques in at least one of these ways.

The first of these techniques is called a Bag of Words (BoW) model, and it consists of simply counting words as they appear in text. It can be

accomplished very easily with a dictionary that scans through text, creating a new vocabulary entry for each new word as a key and an incrementing value starting at 1.

```
sentence = "What is a bag of words and what does it do for me whe  
clean_text = sentence.lower().split(" ")  
bow = {word:clean.count(word) for word in clean_text}  
print(bow)  
#{'what': 2, 'is': 1, 'a': 1, 'bag': 1, 'of': 1, 'words': 1, 'and
```

Considering its simplicity, even this model based entirely on frequency can be quite powerful when trying to gain insight into a speaker's intentions or at least their idiosyncrasies. For example, you could run a simple BoW model on inaugural speeches of US presidents, searching for the words freedom, economy, and enemy to gain a pretty good insight about which presidents assumed office under peacetime, during wartime, and during times of monetary strife, just based on how many times each word was mentioned. The BoW model's weaknesses are many, however, as the model provides no images, semantics, pragmatics, phrases, or feelings. Just above, there are two instances of "words," but because our tokenization strategy is just whitespace it didn't increment the key in the model. It doesn't have any mechanisms to evaluate context or phonetics, and because it divides words by default on white space (you can obviously tokenize however you want, but try tokenizing on subwords and see what happens with this model—spoiler it is bad), it doesn't account for morphology either. Altogether, it should be considered a weak model for representing language, but a strong baseline for evaluating other models against. In order to solve the problem of Bag of Words models not capturing any sequence data, N-Gram models were conceived.

### **2.2.1 N-Gram and Corpus-based techniques**

N-Gram models represent a marked but efficient improvement to BoW, where you are able to give the model a sort of context, represented by N. They are relatively simple statistical models that allow you to generate words based on the N-1 context space. Looking at Listing 2.1, I'm using trigrams which means N=3. I clean the text and give it minimal padding/formatting to help the model, then we train using everygrams, which is meant to prioritize

flexibility over efficiency so that you could train a pentagram or a septagram (N=5, N=7) model if you wanted instead. At the end of the listing where I'm generating, I can give the model up to 2 tokens to help it figure out how to generate further. N-Gram models were not created, and have never even claimed to attempt, complete modeling systems of linguistic knowledge, but they are widely useful in practical applications. They ignore all linguistic features, including syntax, and only attempt to draw probabilistic connections between words appearing in an N-length phrase.

#### **Listing 2.1 Generative N-Grams Language Model Implementation**

```
from nltk.corpus.reader import PlaintextCorpusReader
from nltk.util import everygrams
from nltk.lm.preprocessing import (
    pad_both_ends,
    flatten,
    padded_everygram_pipeline,
)
from nltk.lm import MLE

# Create a corpus from any number of plain .txt files
my_corpus = PlaintextCorpusReader("./", ".*\\.txt")

for sent in my_corpus.sents(fileids="hamlet.txt"):
    print(sent)

# Pad each side of every line in the corpus with <s> and </s> to
padded_trigrams = list(
    pad_both_ends(my_corpus.sents(fileids="hamlet.txt")[1104], n=
)
list(everygrams(padded_trigrams, max_len=3))

list(
    flatten(
        pad_both_ends(sent, n=2)
        for sent in my_corpus.sents(fileids="hamlet.txt")
    )
)

# Allow everygrams to create a training set and a vocab object fr
train, vocab = padded_everygram_pipeline(
    3, my_corpus.sents(fileids="hamlet.txt")
)

# Instantiate and train the model we'll use for N-Grams, a Maximu
```

```
# This model will take the everygrams vocabulary, including the <
lm = MLE(3)
len(lm.vocab)

lm.fit(train, vocab)
print(lm.vocab)
len(lm.vocab)

# And finally, language can be generated with this model and cond
lm.generate(6, ["to", "be"])
```

The above code is all that you need to create a generative N-gram model. For those of you interested in being able to evaluate that model further, we've included the below code to grab probabilities and log scores, or analyze the entropy and perplexity of a particular phrase. Because this is all frequency-based, even though it's mathematically significant, it still does a pretty bad job of describing how perplexing or frequent real-world language actually is.

```
# Any set of tokens up to length=n can be counted easily to deter
print(lm.counts)
lm.counts[["to"]]["be"]

# Any token can be given a probability of occurrence, and can be
print(lm.score("be"))
print(lm.score("be", ["to"]))
print(lm.score("be", ["not", "to"]))

# This can be done as a log score as well to avoid very big and v
print(lm.logscore("be"))
print(lm.logscore("be", ["to"]))
print(lm.logscore("be", ["not", "to"]))

# Sets of tokens can be tested for entropy and perplexity as well
test = [("to", "be"), ("or", "not"), ("to", "be")]
print(lm.entropy(test))
print(lm.perplexity(test))
```

While this code example illustrates creating a trigram language model, unfortunately, not all phrases needing to be captured are only 3 tokens long. For example, from Hamlet, “To be or not to be,” consists of one phrase with 2 words and one phrase with 4 words. One point of note, even though N-Grams are typically very small language models, it is possible to make an N-Gram LLM by just making N=1000000000 or higher but don't expect to get even one ounce of use out of it. Just because we made it big doesn't make it

better or mean it'll have any practical application. 99.9% of all text and 100% of all meaningful text contains fewer than one billion tokens appearing more than once and that computational power can be much better spent elsewhere.

Looking back at Figure 2.2, N-Grams really only use static signals (whitespace, orthography) and words to try to extract any meaning. It tries to measure phrases manually, assuming that all of the phrases will be the same length. That said, N-Grams can be used to create powerful baselines for text analysis, and if the pragmatic context of the utterance is already known by the analyst, they can be used to give quick and accurate insight into real-world scenarios. This type of phrasal modeling also fails to capture any semantic encodings that individual words could have. In order to solve these problems, Bayesian statistics were applied to language modeling.

## 2.2.2 Bayesian Techniques

Bayes' theorem is one of the most mathematically sound and simple theories present in describing the occurrence of your output within your input space. Essentially, it calculates the probability of an event occurring based on prior knowledge. The theorem posits that the probability of a hypothesis being true given evidence, for example that a sentence has a positive sentiment, is equal to the probability of the evidence occurring given the hypothesis is true multiplied by the probability of the hypothesis occurring, all divided by the probability of the evidence being true. Or, expressed mathematically:

$$P(\text{hypothesis} \mid \text{evidence}) = (P(\text{evidence} \mid \text{hypothesis}) * P(\text{hypothesis})) / P(\text{evidence})$$

or

$$P(A|B) * P(B) = P(B|A) * P(A)$$

Because this isn't a math book, we'll dive into Bayes' theorem the exact same depth we dove into any of the linguistics concepts above and trust the interested reader to search for more.

Unfortunately, even though the theorem represents the data in a mathematically sound way, it doesn't account for any stochasticity or

multiple meanings of words. One word you can always throw at a Bayesian model to confuse it is the word, “it.” Any demonstrative pronoun ends up getting assigned values in the same LogPrior and LogLikelihood way as all of the other words and gets a static value, which is antithetical to the usage of those words. For example, if you’re trying to perform sentiment analysis on an utterance, it would be better for you to assign all pronouns a null value than to even let them go through the Bayesian training. It should be noted also that Bayesian techniques don’t end up creating generative language models the way the rest of these will. Because of the nature of Bayes’ theorem validating a hypothesis, these models work for classification, and can bring powerful augmentation to a generative language model.

In Listing 2.2 we show how to create a Naive Bayes classification language model. Instead of using a package like sklearn or something that would make writing the code a little easier, we opted to write out what we were doing, so it’s a bit longer, but should be more informative about how it works. We are using the least-complex version of a Naive Bayes model. We haven’t made it multinomial or added anything fancy and this could obviously work better if you opted to upgrade it for any problem you want. And we highly recommend you do.

#### **Listing 2.2 Categorical Naive Bayes Language Model Implementation**

```
from utils import process_utt, lookup
from nltk.corpus.reader import PlaintextCorpusReader
import numpy as np

my_corpus = PlaintextCorpusReader("./", ".*\\.txt")

sents = my_corpus.sents(fileids="hamlet.txt")

def count_utts(result, utts, ys):
    """
    Input:
        result: a dictionary that is used to map each pair to its
        utts: a list of utts
        ys: a list of the sentiment of each utt (either 0 or 1)
    Output:
        result: a dictionary mapping each pair to its frequency
    """
```

```

for y, utt in zip(ys, utts):
    for word in process_utt(utt):
        # define the key, which is the word and label tuple
        pair = (word, y)

        # if the key exists in the dictionary, increment the
        if pair in result:
            result[pair] += 1

        # if the key is new, add it to the dict and set the c
        else:
            result[pair] = 1

return result

result = {}
utts = [" ".join(sent) for sent in sents]
ys = [sent.count("be") > 0 for sent in sents]
count_utts(result, utts, ys)

freqs = count_utts({}, utts, ys)
lookup(freqs, "be", True)
for k, v in freqs.items():
    if "be" in k:
        print(f"{k}:{v}")

def train_naive_bayes(freqs, train_x, train_y):
    """
    Input:
        freqs: dictionary from (word, label) to how often the wor
        train_x: a list of utts
        train_y: a list of labels correponding to the utts (0,1)
    Output:
        logprior: the log prior.
        loglikelihood: the log likelihood of you Naive bayes equa
    """
    loglikelihood = {}
    logprior = 0

    # calculate V, the number of unique words in the vocabulary
    vocab = set([pair[0] for pair in freqs.keys()])
    V = len(vocab)

    # calculate N_pos and N_neg

```

```

N_pos = N_neg = 0
for pair in freqs.keys():
    # if the label is positive (greater than zero)
    if pair[1] > 0:
        # Increment the number of positive words (word, label)
        N_pos += lookup(freqs, pair[0], True)

    # else, the label is negative
    else:
        # increment the number of negative words (word,label)
        N_neg += lookup(freqs, pair[0], False)

# Calculate D, the number of documents
D = len(train_y)

# Calculate the number of positive documents
D_pos = sum(train_y)

# Calculate the number of negative documents
D_neg = D - D_pos

# Calculate logprior
logprior = np.log(D_pos) - np.log(D_neg)

# For each word in the vocabulary...
for word in vocab:
    # get the positive and negative frequency of the word
    freq_pos = lookup(freqs, word, 1)
    freq_neg = lookup(freqs, word, 0)

    # calculate the probability that each word is positive, a
    p_w_pos = (freq_pos + 1) / (N_pos + V)
    p_w_neg = (freq_neg + 1) / (N_neg + V)

    # calculate the log likelihood of the word
    loglikelihood[word] = np.log(p_w_pos / p_w_neg)

return logprior, loglikelihood

def naive_bayes_predict(utt, logprior, loglikelihood):
    """
    Input:
        utt: a string
        logprior: a number
        loglikelihood: a dictionary of words mapping to numbers
    Output:

```



```

    """ p: the sum of all the logliklihoods + logprior
    """
    # process the utt to get a list of words
    word_l = process_utt(utt)

    # initialize probability to zero
    p = 0

    # add the logprior
    p += logprior

    for word in word_l:
        # check if the word exists in the loglikelihood dictionary
        if word in loglikelihood:
            # add the log likelihood of that word to the probability
            p += loglikelihood[word]

    return p

def test_naive_bayes(test_x, test_y, logprior, loglikelihood):
    """
    Input:
        test_x: A list of utts
        test_y: the corresponding labels for the list of utts
        logprior: the logprior
        loglikelihood: a dictionary with the loglikelihoods for each word
    Output:
        accuracy: (# of utts classified correctly)/(total # of utts)
    """
    accuracy = 0 # return this properly

    y_hats = []
    for utt in test_x:
        # if the prediction is > 0
        if naive_bayes_predict(utt, logprior, loglikelihood) > 0:
            # the predicted class is 1
            y_hat_i = 1
        else:
            # otherwise the predicted class is 0
            y_hat_i = 0

        # append the predicted class to the list y_hats
        y_hats.append(y_hat_i)

    # error = avg of the abs vals of the diffs between y_hats and test_y
    error = sum(

```

```

    [abs(y_hat - test) for y_hat, test in zip(y_hats, test_y)
    ) / len(y_hats)

# Accuracy is 1 minus the error
accuracy = 1 - error

return accuracy

if __name__ == "__main__":
    logprior, loglikelihood = train_naive_bayes(freqs, utts, ys)
    print(logprior)
    print(len(loglikelihood))

    my_utt = "To be or not to be, that is the question."
    p = naive_bayes_predict(my_utt, logprior, loglikelihood)
    print("The expected output is", p)

    print(
        "Naive Bayes accuracy = %0.4f"
        % (test_naive_bayes(utts, ys, logprior, loglikelihood))
    )

```

This theorem doesn't create the same type of language model, but one with a list of probabilities associated with one hypothesis. As such, Bayesian language models can't be used effectively to generate language, but it can be very powerfully implemented for classification tasks. In my opinion though, Bayesian models are often overhyped for even this task. One of the crowning achievements of my career was replacing and removing a Bayesian model from production.

In Bayesian models one of the big issues is essentially that all sequences are completely unconnected, like BoW models, moving us to the opposite end of sequence modeling from N-Grams. Similarly to a pendulum, language modeling swings back towards sequence modeling and language generation with Markov Chains.

### 2.2.3 Markov Chains

Often called Hidden Markov Models (HMMs), Markov Chains essentially add state to the N-Gram models mentioned before, storing probabilities using hidden states. They are often used to help parse text data for even larger

models, doing things like Part-of-Speech tagging (PoS Tagging, marking words with their part of speech) and Named Entity Recognition (NER, marking identifying words with their referent and usually type, e.g. LA - Los Angeles - City) on textual data. Markov models, building upon the previous Bayesian models, rely completely on stochasticity (predictable randomness) in the tokens encountered. The idea is similarly mathematically sound, however, that the probability of anything happening *next* depends completely upon the state of *now*. So instead of modeling words based solely on their historical occurrence and drawing a probability from that, we model their future and past collocation based on what is currently occurring. So the probability of “happy” occurring goes down to almost zero if “happy” was just output, but goes up significantly if “am” has just occurred. Markov chains are so intuitive that they were incorporated into later iterations of Bayesian statistics, and are still used in production systems today.

In Listing 2.3 we train a Markov chain generative language model. This is the first model where we’ve used a specific tokenizer, which in this case will just tokenize based on the white space between words. This is also only the second time we’ve referred to a collection of utterances meaning to be viewed together as a document. As you play around with this one, pay close attention and make some comparisons yourself for how well the HMM generates compared to even a large N-gram model.

#### **Listing 2.3 Generative Hidden Markov Language Model Implementation**

```
import re
import random
from nltk.tokenize import word_tokenize
from collections import defaultdict, deque

class MarkovChain:
    def __init__(self):
        self.lookup_dict = defaultdict(list)
        self._seeded = False
        self.__seed_me()

    def __seed_me(self, rand_seed=None):
        if self._seeded is not True:
            try:
                if rand_seed is not None:
```

```

        random.seed(rand_seed)
    else:
        random.seed()
    self._seeded = True
except NotImplementedError:
    self._seeded = False

def add_document(self, str):
    preprocessed_list = self._preprocess(str)
    pairs = self.__generate_tuple_keys(preprocessed_list)
    for pair in pairs:
        self.lookup_dict[pair[0]].append(pair[1])

def _preprocess(self, str):
    cleaned = re.sub(r"\W+", " ", str).lower()
    tokenized = word_tokenize(cleaned)
    return tokenized

def __generate_tuple_keys(self, data):
    if len(data) < 1:
        return

    for i in range(len(data) - 1):
        yield [data[i], data[i + 1]]

def generate_text(self, max_length=50):
    context = deque()
    output = []
    if len(self.lookup_dict) > 0:
        self.__seed_me(rand_seed=len(self.lookup_dict))
        chain_head = [list(self.lookup_dict)[0]]
        context.extend(chain_head)

        while len(output) < (max_length - 1):
            next_choices = self.lookup_dict[context[-1]]
            if len(next_choices) > 0:
                next_word = random.choice(next_choices)
                context.append(next_word)
                output.append(context.popleft())
            else:
                break
        output.extend(list(context))
    return " ".join(output)

if __name__ == "__main__":
    with open("hamlet.txt", "r", encoding="utf-8") as f:

```

```
text = f.read()
HMM = MarkovChain()
HMM.add_document(text)

print(HMM.generate_text(max_length=25))
```

This code shows a basic implementation of a Markov model for generation, and we'd encourage the reader to experiment with it, give it text from songs from your favorite musicians or books from your favorite authors and see whether what comes out actually sounds like them. HMMs are incredible fast and often used in predictive text or predictive search applications. Markov models represent the first comprehensive attempt to actually model language from a descriptive linguistic perspective, as opposed to a prescriptive one, which is interesting, because Markov did not originally intend for linguistic modeling, only to win an argument about continuous independent states. Later, Markov used Markov Chains to model vowel distribution in a Pushkin novel, so he was at least aware of the possible applications.

The difference between descriptive and prescriptive linguistics is that one focuses on how things *ought* to be, while the other focuses on how things *are*. From a language modeling perspective, it has proven vastly more effective to describe what language is doing from a corpus or Markov perspective, rather than to attempt to prescribe how language ought to behave. Unfortunately, a current state by itself cannot be used to give context beyond the now, so historical or societal context is unable to be represented effectively in a Markov model. Semantic encoding of words also becomes a problem, as is represented in the code example, Markov chains will output syntactically correct chains of words that semantically are nonsense, similar to “colorless green ideas sleep furiously.” To attempt to solve this problem, “continuous” models were developed to allow for a “semantic embedding” representation of tokens.

## 2.2.4 Continuous Language Modeling

A Continuous Bag of Words (CBoW), much like its namesake, the Bag of Words, is a frequency-based approach to analyzing language, meaning that it models words based on how often they occur. The next word in an utterance has never been determined based on probability or frequency. Due to this, the

example given will be for how to create word embeddings to be ingested or compared by other models using a CBoW. We'll use a neural network for this to give you a good methodology.

This is the first language modeling technique we'll see that essentially slides a context window over a given utterance (the context window is an N-gram model) and attempts to guess what word is in the middle based upon surrounding words in the window. For example, let's say your window has a length of 5, and your sentence is, "Learning about linguistics makes me happy," you would give the CBoW ['learning', 'about', 'makes', 'me'] and try to get the model to guess "linguistics," based upon how many times the model has seen that word occur in similar places previously. This should show you why generation is difficult for models trained like this, because if you give the model ['makes', 'me', '</s>'] as input, first of all it only has 3 pieces of information to try to figure out instead of 4, and it also will be biased towards only guessing words it has seen at the end of sentences before, as opposed to getting ready to start new clauses. It's not all bad though, one feature that makes continuous models stand out for embeddings is that it doesn't just have to look at words before the target word, they can also use words that come after the target to gain some semblance of context.

In Listing 2.4 we create our first continuous model. In our case, to keep things as simple as possible, we use a bag of words for the language processing and a one-layer neural network with two parameters for the embedding estimation, although both of those could be substituted out for any other models. For example, you could substitute N-grams for the BoW and a Naive Bayes for the neural network, and get a Continuous Naive N-gram model. The point is that the actual models used in this technique are a bit arbitrary, it's more the Continuous technique that's important. To illustrate this further, we don't use any packages other than numpy to do the math for the neural network, even though it's our first one appearing in this section.

Pay special attention to the steps below, initializing the model weights, the ReLU activation function, the final softmax layer, forward and backpropagation, and then how it all fits together in the `gradient_descent` function. These are pieces in the puzzle that you will see crop up again and

again, regardless of programming language or framework. You will need to initialize models, pick activation functions, pick final layers, and define forward and backward propagation in Tensorflow, Pytorch and HuggingFace, and if you ever start creating your own models as opposed to using someone else's.

#### **Listing 2.4 Generative Continuous Bag of Words Language Model Implementation**

```
import nltk
import numpy as np
from utils import get_batches, compute_pca, get_dict
import re
from matplotlib import pyplot

# Create our corpus for training
with open("hamlet.txt", "r", encoding="utf-8") as f:
    data = f.read()

# Slightly clean the data by removing punctuation, tokenizing by
data = re.sub(r"[,!?;-]", ".", data)
data = nltk.word_tokenize(data)
data = [ch.lower() for ch in data if ch.isalpha() or ch == "."]
print("Number of tokens:", len(data), "\n", data[500:515])

# Get our Bag of Words, along with a distribution
fdist = nltk.FreqDist(word for word in data)
print("Size of vocabulary:", len(fdist))
print("Most Frequent Tokens:", fdist.most_common(20))

# Create 2 dictionaries to speed up time-to-convert and keep track
word2Ind, Ind2word = get_dict(data)
V = len(word2Ind)
print("Size of vocabulary:", V)

print("Index of the word 'king':", word2Ind["king"])
print("Word which has index 2743:", Ind2word[2743])

# Here we create our Neural network with 1 layer and 2 parameters
def initialize_model(N, V, random_seed=1):
    """
    Inputs:
        N: dimension of hidden vector
        V: dimension of vocabulary
        random_seed: seed for consistent results in tests
    Outputs:
```

```

        W1, W2, b1, b2: initialized weights and biases
    """
    np.random.seed(random_seed)

    W1 = np.random.rand(N, V)
    W2 = np.random.rand(V, N)
    b1 = np.random.rand(N, 1)
    b2 = np.random.rand(V, 1)

    return W1, W2, b1, b2

# Create our final classification layer, which makes all possible
def softmax(z):
    """
    Inputs:
        z: output scores from the hidden layer
    Outputs:
        yhat: prediction (estimate of y)
    """
    yhat = np.exp(z) / np.sum(np.exp(z), axis=0)
    return yhat

# Define the behavior for moving forward through our model, along
def forward_prop(x, W1, W2, b1, b2):
    """
    Inputs:
        x: average one-hot vector for the context
        W1, W2, b1, b2: weights and biases to be learned
    Outputs:
        z: output score vector
    """
    h = W1 @ x + b1
    h = np.maximum(0, h)
    z = W2 @ h + b2
    return z, h

# Define how we determine the distance between ground truth and m
def compute_cost(y, yhat, batch_size):
    logprobs = np.multiply(np.log(yhat), y) + np.multiply(
        np.log(1 - yhat), 1 - y
    )
    cost = -1 / batch_size * np.sum(logprobs)
    cost = np.squeeze(cost)
    return cost

# Define how we move backward through the model and collect gradi
def back_prop(x, yhat, y, h, W1, W2, b1, b2, batch_size):

```



```

"""
Inputs:
    x: average one hot vector for the context
    yhat: prediction (estimate of y)
    y: target vector
    h: hidden vector (see eq. 1)
    W1, W2, b1, b2: weights and biases
    batch_size: batch size
Outputs:
    grad_W1, grad_W2, grad_b1, grad_b2: gradients of weights
"""
l1 = np.dot(W2.T, yhat - y)
l1 = np.maximum(0, l1)
grad_W1 = np.dot(l1, x.T) / batch_size
grad_W2 = np.dot(yhat - y, h.T) / batch_size
grad_b1 = np.sum(l1, axis=1, keepdims=True) / batch_size
grad_b2 = np.sum(yhat - y, axis=1, keepdims=True) / batch_size

return grad_W1, grad_W2, grad_b1, grad_b2

# Put it all together and train
def gradient_descent(data, word2Ind, N, V, num_iters, alpha=0.03)
    """
    This is the gradient_descent function

    Inputs:
        data: text
        word2Ind: words to Indices
        N: dimension of hidden vector
        V: dimension of vocabulary
        num_iters: number of iterations
    Outputs:
        W1, W2, b1, b2: updated matrices and biases

    """
    W1, W2, b1, b2 = initialize_model(N, V, random_seed=8855)
    batch_size = 128
    iters = 0
    C = 2
    for x, y in get_batches(data, word2Ind, V, C, batch_size):
        z, h = forward_prop(x, W1, W2, b1, b2)
        yhat = softmax(z)
        cost = compute_cost(y, yhat, batch_size)
        if (iters + 1) % 10 == 0:
            print(f"iters: {iters+1} cost: {cost:.6f}")
        grad_W1, grad_W2, grad_b1, grad_b2 = back_prop(
            x, yhat, y, h, W1, W2, b1, b2, batch_size

```

```

    )
    w1 = w1 - alpha * grad_w1
    w2 = w2 - alpha * grad_w2
    b1 = b1 - alpha * grad_b1
    b2 = b2 - alpha * grad_b2
    iters += 1
    if iters == num_iters:
        break
    if iters % 100 == 0:
        alpha *= 0.66

    return w1, w2, b1, b2

# Train the model
C = 2
N = 50
word2Ind, Ind2word = get_dict(data)
V = len(word2Ind)
num_iters = 150
print("Call gradient_descent")
w1, w2, b1, b2 = gradient_descent(data, word2Ind, N, V, num_iters)
Call gradient descent
Iters: 10 loss: 0.525015
Iters: 20 loss: 0.092373
Iters: 30 loss: 0.050474
Iters: 40 loss: 0.034724
Iters: 50 loss: 0.026468
Iters: 60 loss: 0.021385
Iters: 70 loss: 0.017941
Iters: 80 loss: 0.015453
Iters: 90 loss: 0.012099
Iters: 100 loss: 0.012099
Iters: 110 loss: 0.011253
Iters: 120 loss: 0.010551
Iters: 130 loss: 0.009932
Iters: 140 loss: 0.009382
Iters: 150 loss: 0.008889

```

The CBoW example is our first code example to showcase a full and effective training loop in machine learning. Within all of that, we asked the reader to pay special attention to the steps in a training loop, especially the activation function, ReLU. As we expect the reader to be at least familiar with various ML paradigms, including different activations, we won't explain the ReLU here, rather why you should use it and why you shouldn't. ReLUs, while solving the vanishing gradient problem, don't solve the exploding

gradient problem, and they sharply destroy all negative comparisons within the model. Better situational variants include the ELU, which allows negative numbers normalizing to alpha, or the GEGLU/SWIGLU, which works well in increasingly perplex scenarios, like language. However, people often use ReLUs, not because they are the best in a situation, but because they are easy-to-understand, easy-to-code, and intuitive, even more so than the activations they were created to replace like the sigmoid or tanh.

A lot of this ends up being abstracted with packages and the like, but knowing what's going on under the hood will be very helpful for you as someone putting LLMs in production. You should be able to predict with some certainty how different models will behave in various situations. The next section will dive into one of those abstractions, in this case being the abstraction that is created by the continuous modeling technique.

## **2.2.5 Embeddings**

Harkening back to our features of language, it should be easy to connect why continuous-style language modeling was such a breakthrough.

Embeddings take the tokenized vectors we've created that don't contain any meaning, and attempt to insert that meaning based on observations that can be made about the text, such as word order and subwords appearing in similar contexts. Despite the primary mode of meaning being collocation (co-located, words that appear next to each other), they prove useful and even show some similarities to human-encoded word meaning.

The quintessential example from Word2Vec, one of the first pre-trained vector embeddings, was taking the vector for “king” subtracting the vector for “man” adding the vector for “woman” and finding the nearest neighbor to the sum was the vector for the word “queen”. This makes sense to us as it mimics human semantics. One of the major differences is one that's already been mentioned a couple of times, pragmatics. Humans use pragmatic context to inform semantic meaning, understanding that just because you said, “I need food,” doesn't mean you are actually in physical danger without it. Embeddings are devoid of any influence outside of pure usage, which feels like it could be how humans learn as well, and there are good arguments on all sides here. The one thing holding is that if we can somehow give models

more representative data, that may open the door to more effective embeddings, but it's a chicken and egg problem because more effective embeddings give better model performance.

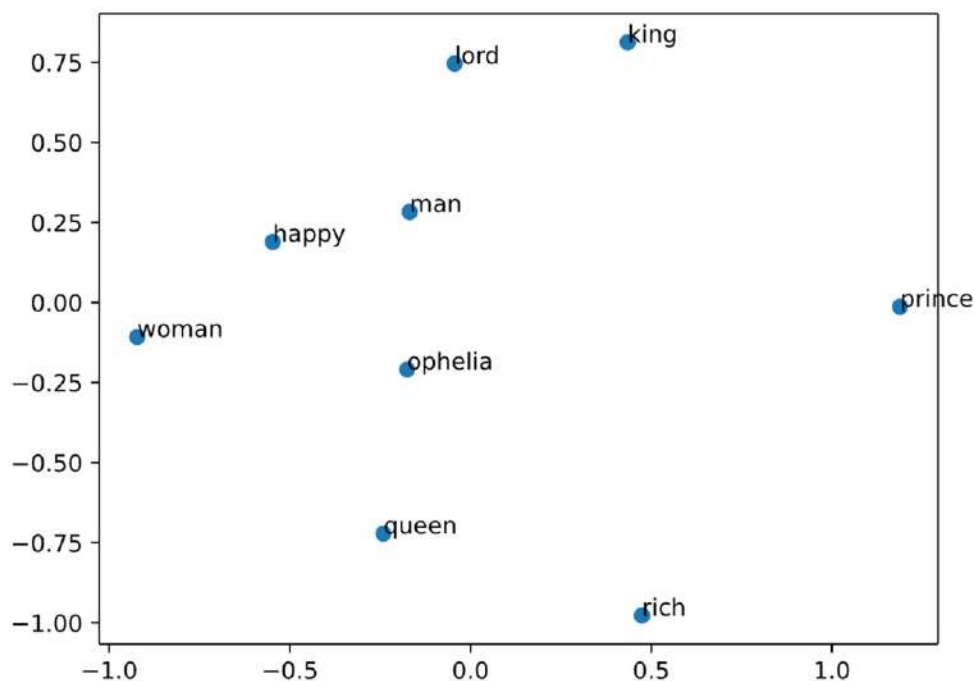
In Listing 2.5, we'll dive into how to visualize embeddings using pyplot. We will be going more in-depth into embeddings in later chapters. This is helpful for model explainability and also for validation during your pre-training step. If you see that your semantically similar embeddings are relatively close to each other on the graph, then you're likely going in the right direction.

#### **Listing 2.5 Embedding Visualization**

```
# After listing 2.4 is done and gradient descent has been execute
words = [
    "King",
    "Queen",
    "Lord",
    "Man",
    "Woman",
    "Prince",
    "Ophelia",
    "Rich",
    "Happy",
]
embs = (w1.T + w2) / 2.0
idx = [word2Ind[word] for word in words]
X = embs[idx, :]
print(X.shape, idx)

result = compute_pca(X, 2)
pyplot.scatter(result[:, 0], result[:, 1])
for i, word in enumerate(words):
    pyplot.annotate(word, xy=(result[i, 0], result[i, 1]))
pyplot.show()
```

**Figure 2.4 A visualization technique for word embeddings. Visualizing embeddings can be important for model explainability.**



As we can see in figure 2.4, this is a successful, but a very sparse embedding representation that we trained from our CBoW model. Getting those semantic representations (embeddings) to be denser is the main place that we can see improvement in this field, although many successful experiments have been run where denser semantic meaning has been supplanted with greater pragmatic context through instruct and different thought chaining techniques. We will address Chain of Thought (CoT) and other techniques later. For now, let's pivot to discussing why our continuous embedding technique can even be successful, given frequency-based models are characteristically difficult to correlate with reality. All of this starts with the Multilayer Perceptron, more than half a century ago.

## 2.2.6 Multilayer Perceptrons

MLPs are the embodiment of the sentiment, “Machines are really good at doing one thing, so I wish we could just use a bunch of machines that are really good at the one thing to make one that's good at a lot of things.” Every weight and bias in the neural network of the MLP is good at doing one thing, which could be detecting one or more features. So, we bind a whole bunch of them together to detect larger, more complex features. MLPs serve as the

primary building block in most neural network architectures. The key distinctions between architectures, such as convolutional neural networks and recurrent neural networks, mainly arise from data loading methods and the handling of tokenized and embedded data as it flows through the layers of the model, rather than the functionality of individual layers, particularly the fully-connected layers.

In Listing 2.6 we provide a more dynamic class of neural network that can have as many layers and parameters as deemed necessary for your task. We give a more-defined and explicit class using pytorch to give you the tools to implement the MLP for use in whatever you'd like, both from scratch and in a popular framework.

**Listing 2.6 Multilayer Perceptron Pytorch Class Implementation**

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class MultiLayerPerceptron(nn.Module):
    def __init__(
        self,
        input_size,
        hidden_size=2,
        output_size=3,
        num_hidden_layers=1,
        hidden_activation=nn.Sigmoid,
    ):
        """Initialize weights.
        Args:
            input_size (int): size of the input
            hidden_size (int): size of the hidden layers
            output_size (int): size of the output
            num_hidden_layers (int): number of hidden layers
            hidden_activation (torch.nn.*): the activation class
        """
        super(MultiLayerPerceptron, self).__init__()
        self.module_list = nn.ModuleList()
        interim_input_size = input_size
        interim_output_size = hidden_size
        torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

        for _ in range(num_hidden_layers):
```

```

        self.module_list.append(
            nn.Linear(interim_input_size, interim_output_size
        )
        self.module_list.append(hidden_activation())
        interim_input_size = interim_output_size

    self.fc_final = nn.Linear(interim_input_size, output_size

    self.last_forward_cache = []

def forward(self, x, apply_softmax=False):
    """The forward pass of the MLP

    Args:
        x_in (torch.Tensor): an input data tensor.
            x_in.shape should be (batch, input_dim)
        apply_softmax (bool): a flag for the softmax activation
            should be false if used with the Cross Entropy loss
    Returns:
        the resulting tensor. tensor.shape should be (batch,
    """
    for module in self.module_list:
        x = module(x)

    output = self.fc_final(x)

    if apply_softmax:
        output = F.softmax(output, dim=1)

    return output

```

From the code we can see, as opposed to the CBoW implementation which had a static two layers, this MLP is not static in size until it has been instantiated. If you wanted to give this model one million layers, you would just have to put `num_hidden_layers=1000000` when you instantiate the class, although just because you *can* give a model that many parameters it won't make it immediately better. LLMs are more than just a lot of layers. Like RNNs and CNNs, the magic of LLMs is in how data goes in and moves through the model. To illustrate, let's look at the RNN and one of its variations.

## 2.2.7 RNNs and LSTMs

Recurrent Neural Networks (RNNs) are a class of neural networks designed

to analyze sequences, based on the weaknesses in previous language modeling techniques. A sequence can be thought of as an ordered array, where the sum of the whole array changes value if any of the parts are moved around. The logic goes that if language is presented in a sequence, then maybe it should be processed in a sequence, as opposed to one token at a time. RNNs accomplish this by using logic we've seen before, both in MLPs and in Markov Chains, where an internal state or memory is referred to when new inputs are processed, and creating cycles when connections between nodes are detected as being useful.

In fully recurrent networks, like the one in Listing 2.7, all nodes start out initially connected to all subsequent nodes, but those connections can be set to zero to simulate them being broken if they are not useful. This solves one of the biggest problems that earlier models suffered from, static input size, and enables an RNN and its variants to process variable length inputs. Unfortunately, longer sequences create a new problem. Because each neuron in the network has connections to subsequent neurons, longer sequences create smaller changes to the overall sum, making the gradients smaller and eventually vanishing, even with important words. This is called a vanishing gradient, but other problems exist too, such as exploding and diminishing gradients.

For example, let's consider these sentences with the task sentiment analysis, "I loved the movie last night," and, "The movie I went to see last night was the very best I had ever expected to see." These sentences can be considered semantically similar, even if they aren't exactly the same. When moving through an RNN, each word in the first sentence is worth more, and the consequence is that the first sentence has a higher positive rating than the second sentence, just because of the first sentence being shorter. The inverse is true also, exploding gradients are also a consequence of this sequence processing, which makes training deep RNNs difficult.

To solve this problem, long short-term memories (LSTMs), which are a type of RNN, use memory cells and gating mechanisms to keep being able to process sequences of variable length, but without the problems of longer and shorter sequences being comprehended differently. Anticipating multilingual scenarios and understanding that people don't think about language in only



one direction, LSTMs can also process sequences bidirectionally by concatenating the outputs of two RNNs, one reading the sequence from left to right, and the other from right to left. This bidirectionality improves results, allowing for information to be seen and remembered even after thousands of tokens have passed.

In Listing 2.7 we give classes for both an RNN and an LSTM. In the code in the repo associated with this book, you can see the results of training both the RNN and LSTM, where the takeaway is that the LSTM gets better accuracy on both training and validation sets in half as many epochs (25 vs 50 with RNN). One of the innovations to note is the packed embeddings that utilize padding, extending all variable-length sequences to the maximum length in order to allow processing any length input, as long as it is shorter than the maximum.

**Listing 2.7 Recurrent Neural Network and Long Short-Term Memory Pytorch Class Implementations**

```
import torch
from gensim.models import Word2Vec
from sklearn.model_selection import train_test_split

# Create our corpus for training
with open("./chapters/chapter_2/hamlet.txt", "r", encoding="utf-8") as f:
    data = f.readlines()

# Embeddings are needed to give semantic value to the inputs of a
# embedding_weights = torch.Tensor(word_vectors.vectors)

EMBEDDING_DIM = 100
model = Word2Vec(data, vector_size=EMBEDDING_DIM, window=3, min_count=1)
word_vectors = model.wv
print(f"Vocabulary Length: {len(model.wv)}")
del model

padding_value = len(word_vectors.index_to_key)
embedding_weights = torch.Tensor(word_vectors.vectors)

class RNN(torch.nn.Module):
    def __init__(
        self,
        input_dim,
```

```

        embedding_dim,
        hidden_dim,
        output_dim,
        embedding_weights,
    ):
        super().__init__()
        self.embedding = torch.nn.Embedding.from_pretrained(
            embedding_weights
        )
        self.rnn = torch.nn.RNN(embedding_dim, hidden_dim)
        self.fc = torch.nn.Linear(hidden_dim, output_dim)

    def forward(self, x, text_lengths):
        embedded = self.embedding(x)
        packed_embedded = torch.nn.utils.rnn.pack_padded_sequence(
            embedded, text_lengths
        )
        packed_output, hidden = self.rnn(packed_embedded)
        output, output_lengths = torch.nn.utils.rnn.pad_packed_sequence(
            packed_output
        )
        return self.fc(hidden.squeeze(0))

```

```

INPUT_DIM = 4764
EMBEDDING_DIM = 100
HIDDEN_DIM = 256
OUTPUT_DIM = 1

```

```

model = RNN(
    INPUT_DIM, EMBEDDING_DIM, HIDDEN_DIM, OUTPUT_DIM, embedding_w
)

```

```

optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
criterion = torch.nn.BCEWithLogitsLoss()
device = torch.device("cuda" if torch.cuda.is_available() else "c

```

```

class LSTM(torch.nn.Module):
    def __init__(
        self,
        input_dim,
        embedding_dim,
        hidden_dim,
        output_dim,
        n_layers,
        bidirectional,
    ):

```

```

        dropout,
        embedding_weights,
    ):
        super().__init__()
        self.embedding = torch.nn.Embedding.from_pretrained(
            embedding_weights
        )
        self.rnn = torch.nn.LSTM(
            embedding_dim,
            hidden_dim,
            num_layers=n_layers,
            bidirectional=bidirectional,
            dropout=dropout,
        )
        self.fc = torch.nn.Linear(hidden_dim * 2, output_dim)
        self.dropout = torch.nn.Dropout(dropout)

    def forward(self, x, text_lengths):
        embedded = self.embedding(x)
        packed_embedded = torch.nn.utils.rnn.pack_padded_sequence(
            embedded, text_lengths
        )
        packed_output, (hidden, cell) = self.rnn(packed_embedded)
        hidden = self.dropout(
            torch.cat((hidden[-2, :, :], hidden[-1, :, :]), dim=1)
        )
        return self.fc(hidden.squeeze(0))

```

```

INPUT_DIM = padding_value
EMBEDDING_DIM = 100
HIDDEN_DIM = 256
OUTPUT_DIM = 1
N_LAYERS = 2
BIDIRECTIONAL = True
DROPOUT = 0.5

```

```

model = LSTM(
    INPUT_DIM,
    EMBEDDING_DIM,
    HIDDEN_DIM,
    OUTPUT_DIM,
    N_LAYERS,
    BIDIRECTIONAL,
    DROPOUT,
    embedding_weights,
)

```

```

optimizer = torch.optim.Adam(model.parameters())
criterion = torch.nn.BCEWithLogitsLoss()
device = torch.device("cuda" if torch.cuda.is_available() else "c

def binary_accuracy(preds, y):
    rounded_preds = torch.round(torch.sigmoid(preds))
    correct = (rounded_preds == y).float()
    acc = correct.sum()/len(correct)
    return acc

def train(model, iterator, optimizer, criterion):
    epoch_loss = 0
    epoch_acc = 0
    model.train()
    for batch in iterator:
        optimizer.zero_grad()
        predictions = model(batch["text"], batch["length"]).squeeze_
        loss = criterion(predictions, batch["label"])
        acc = binary_accuracy(predictions, batch["label"])
        loss.backward()
        optimizer.step()
        epoch_loss += loss.item()
        epoch_acc += acc.item()

    return epoch_loss / len(iterator), epoch_acc / len(iterator)

def evaluate(model, iterator, criterion):
    epoch_loss = 0
    epoch_acc = 0
    model.eval()
    with torch.no_grad():
        for batch in iterator:
            predictions = model(batch["text"], batch["length"]).s
            loss = criterion(predictions, batch["label"])
            acc = binary_accuracy(predictions, batch["label"])

            epoch_loss += loss.item()
            epoch_acc += acc.item()

    return epoch_loss / len(iterator), epoch_acc / len(iterator)

batch_size = 128 # Usually should be a power of 2 because it's th

def iterator(X, y):
    size = len(X)

```

```

permutation = np.random.permutation(size)
iterate = []
for i in range(0, size, batch_size):
    indices = permutation[i:i+batch_size]
    batch = {}
    batch['text'] = [X[i] for i in indices]
    batch['label'] = [y[i] for i in indices]

    batch['text'], batch['label'] = zip(*sorted(zip(batch['text']
    batch['length'] = [len(utt) for utt in batch['text']]
    batch['length'] = torch.IntTensor(batch['length'])
    batch['text'] = torch.nn.utils.rnn.pad_sequence(batch['text']
    batch['label'] = torch.Tensor(batch['label']))

    batch['label'] = batch['label'].to(device)
    batch['length'] = batch['length'].to(device)
    batch['text'] = batch['text'].to(device)

    iterate.append(batch)

return iterate

index_utt = word_vectors.key_to_index

#You've got to determine some labels for whatever you're training
X_train, X_test, y_train, y_test = train_test_split(index_utt, la
X_train, X_val, y_train, y_val = train_test_split(X_train, y_trai

train_iterator = iterator(X_train, y_train)
validate_iterator = iterator(X_val, y_val)
test_iterator = iterator(X_test, y_test)

print(len(train_iterator), len(validate_iterator), len(test_itera

N_EPOCHS = 25

for epoch in range(N_EPOCHS):
    train_loss, train_acc = train(
        model, train_iterator, optimizer, criterion
    )
    valid_loss, valid_acc = evaluate(model, validate_iterator, cr

    print(
        f"| Epoch: {epoch+1:02} | Train Loss: {train_loss: .3f} |
    )
#Training on our dataset

```

|           |                   |                   |           |
|-----------|-------------------|-------------------|-----------|
| Epoch: 01 | Train Loss: 0.560 | Train Acc: 70.63% | Validatio |
| Epoch: 05 | Train Loss: 0.391 | Train Acc: 82.81% | Validatio |
| Epoch: 10 | Train Loss: 0.270 | Train Acc: 89.11% | Validatio |
| Epoch: 15 | Train Loss: 0.186 | Train Acc: 92.95% | Validatio |
| Epoch: 20 | Train Loss: 0.121 | Train Acc: 95.93% | Validatio |
| Epoch: 25 | Train Loss: 0.100 | Train Acc: 96.28% | Validatio |

Looking above at our classes and instantiations, you should see that the LSTM is not vastly different from the RNN. The only differences in the `init` input variables are `n_layers` (for convenience, you can also specify it with RNNs), `bidirectional`, and `dropout`. Bidirectional allows LSTMs to look ahead in sequences to help with meaning and context, but also drastically helps with multilingual scenarios, as left-to-right languages like English are not the only format for orthography. Dropout, another huge innovation, changes the paradigm of overfitting from being only data-dependent, and helps the model not overfit by turning off random nodes layer-by-layer during training to force all nodes not to correlate to each other. The only differences in the out-of-model parameters is that the best optimizer for an RNN is SGD, like our CBoW, and the LSTM uses Adam (could use any, including AdamW). Below, we define our training loop and train the LSTM. Compare this training loop to the one defined in Listing 2.4 in the `gradient_descent` function.

One of the amazing things demonstrated in the code here is how much quicker the LSTM can learn, compared to previous model iterations, thanks to both bidirectionality and dropout. The previous models, though training faster, take hundreds of epochs to get the same performance as an LSTM in just 25 epochs. The performance on the validation set, as its name implies, adds validity to the architecture, performing inference during training on examples it has not trained on and keeping accuracy fairly close to the training set.

The problems with these models are not as pronounced, manifesting primarily as being incredibly resource-heavy, especially when being applied to longer, more detail-oriented problems, like healthcare and law. Despite the incredible advantages of Dropout and Bidirectional processing, they both at least double the amount of processing power required to train, so while inference ends up being only 2-3x as expensive as an MLP of the same size, training becomes 10-12x as expensive. They solved exploding gradients

nicely, but exploded the compute required to train instead. To combat this a shortcut was devised and implemented which allowed any model, including an LSTM, to figure out which parts of a sequence were the most influential and which parts could be safely ignored, known as attention.

## 2.2.8 Attention

Attention is a mathematical shortcut for solving larger context windows faster by telling the model through an emergent mathematical formula which parts of an input to consider and how much. This is all based upon an upgraded version of a dictionary, where instead of just Key and Value pairs, a contextual Query is added. We will go more into Attention in later chapters. For now, know that the below code is the big differentiator between older NLP techniques and more modern ones.

Attention solves the slowness of training LSTMs, but keeps the high performance on a low number of epochs. There are multiple types of attention as well. The dot product attention method captures the relationships between each word (or embedding) in your query and every word in your key. When queries and keys are part of the same sentences, this is known as bi-directional self-attention. However, in certain cases, it is more suitable to only focus on words that precede the current one. This type of attention, especially when queries and keys come from the same sentences, is referred to as causal attention. Language modeling further improves by masking parts of a sequence and forcing the model to guess what should be behind the mask. Both Dot Product Attention and masked attention are demonstrated with functions below.

### Listing 2.8 Multi-Head Attention Implementation

```
import numpy as np
from scipy.special import softmax

# Step 1: Input: 3 inputs, d_model=4
x = np.array([[1.0, 0.0, 1.0, 0.0],
              [0.0, 2.0, 0.0, 2.0],
              [1.0, 1.0, 1.0, 1.0]])

# Step 2: weights 3 dimensions x d_model=4
```

```

w_query = np.array([1,0,1],
                    [1,0,0],
                    [0,0,1],
                    [0,1,1]))
w_key = np.array([[0,0,1],
                  [1,1,0],
                  [0,1,0],
                  [1,1,0]])
w_value = np.array([[0,2,0],
                    [0,3,0],
                    [1,0,3],
                    [1,1,0]])

# Step 3: Matrix Multiplication to obtain Q,K,V
## Query: x * w_query
Q = np.matmul(x,w_query)
## Key: x * w_key
K = np.matmul(x,w_key)
## Value: x * w_value
V = np.matmul(x,w_value)

# Step 4: Scaled Attention Scores
## Square root of the dimensions
k_d = 1
attention_scores = (Q @ K.transpose())/k_d

# Step 5: Scaled softmax attention scores for each vector
attention_scores[0] = softmax(attention_scores[0])
attention_scores[1] = softmax(attention_scores[1])
attention_scores[2] = softmax(attention_scores[2])

# Step 6: attention value obtained by score1/k_d * V
attention1 = attention_scores[0].reshape(-1,1)
attention1 = attention_scores[0][0]*V[0]
attention2 = attention_scores[0][1]*V[1]
attention3 = attention_scores[0][2]*V[2]

# Step 7: summed the results to create the first line of the outp
attention_input1 = attention1 + attention2 + attention3

# Step 8: Step 1 to 7 for inputs 1 to 3
## Because this is just a demo, we'll do a random matrix of the r
attention_head1 = np.random.random((3,64))

# Step 9: We train all 8 heads of the attention sub-layer using s
## Again, it's a demo
z0h1 = np.random.random((3,64))

```



```

z1h2 = np.random.random((3, 64))
z2h3 = np.random.random((3, 64))
z3h4 = np.random.random((3, 64))
z4h5 = np.random.random((3, 64))
z5h6 = np.random.random((3, 64))
z6h7 = np.random.random((3, 64))
z7h8 = np.random.random((3, 64))

# Step 10: Concatenate heads 1 through 8 to get the original 8x64
Output_attention = np.hstack((z0h1, z1h2, z2h3, z3h4, z4h5, z5h6, z6h7,

# Here's a function that performs all of these steps:
def dot_product_attention(query, key, value, mask, scale=True):
    assert query.shape[-1] == key.shape[-1] == value.shape[-1], "
    if scale:
        depth = query.shape[-1]
    else:
        depth = 1
    dots = np.matmul(query, np.swapaxes(key, -1, -2)) / np.sqrt(d
    if mask is not None:
        dots = np.where(mask, dots, np.full_like(dots, -1e9))
    logsumexp = scipy.special.logsumexp(dots, axis=-1, keepdims=T
    dots = np.exp(dots - logsumexp)
    attention = np.matmul(dots, value)
    return attention

# Here's a function that performs the previous steps but adds cau
def masked_dot_product_self_attention(q, k, v, scale=True):
    mask_size = q.shape[-2]
    mask = np.tril(np.ones((1, mask_size, mask_size), dtype=np.bo
    return DotProductAttention(q, k, v, mask, scale=scale)

```

Above, in the full implementation of Attention you may have noticed some terminology you're familiar with, namely Key and Value, but you may not have been introduced to Query before. Key and Value pairs are familiar because of dictionaries and lookup tables, where we map a set of keys to an array of values. Query should feel intuitive as a sort of search for retrieval. The Query is compared to the Keys, from which a Value is retrieved in a normal operation.

In Attention, the Query and Keys undergo dot product similarity comparison to obtain an attention score, which is later multiplied by the Value in order to get an ultimate score for how much Attention the model should pay to that portion of the sequence. This can get more complex, depending upon your

model's architecture because both encoder and decoder sequence lengths have to be accounted for, but suffice it to say for now that the most efficient way to model in this space is to project all input sources into a common space and compare using dot product for efficiency.

This code explanation was a bit more math-heavy than the previous examples, but it is needed to illustrate the concept. The math behind Attention is truly innovative and has rocketed the field forward.

Unfortunately, even with the advantages Attention brings to the process of sequence modeling, with LSTMs and RNNs there were still issues with speed and memory size. You may notice from the code and the math that there is a square root taken, meaning that attention as we use it is quadratic. Since then, there have been various techniques, including subquadratics like Hyena and the Recurrent Memory Transformer (RMT, basically an RNN combined with a transformer) to combat these problems, which we will cover in more detail later. For now, let's move on to the ultimate application of Attention: the Transformer.

## 2.3 Attention is All You Need

In the seminal paper, Attention is All You Need[\[1\]](#) Vaswani et al take the mathematical shortcut several steps further, positing that for performance absolutely no recurrence (the “R” in RNN) or any convolutions[\[2\]](#) were needed at all. Instead, they opted to use only Attention and simply specify where Q, K, and V were taken from much more carefully. We'll dive into this presently. In our review of this diverse range of NLP techniques, we have observed their evolution over time and the ways in which each approach has sought to improve upon its predecessors. From rule-based methods to statistical models and neural networks, the field has continually strived for more efficient and accurate ways to process and understand natural language. Now, we turn our attention to a groundbreaking innovation that has revolutionized the field of NLP: the Transformer architecture. In the following section, we will explore the key concepts and mechanisms that underpin Transformers, and how they have enabled the development of state-of-the-art language models that surpass the performance of previous techniques. We will also discuss the impact of Transformers on the broader NLP landscape and consider the potential for further advancements in this

exciting area of research.

### 2.3.1 Encoders

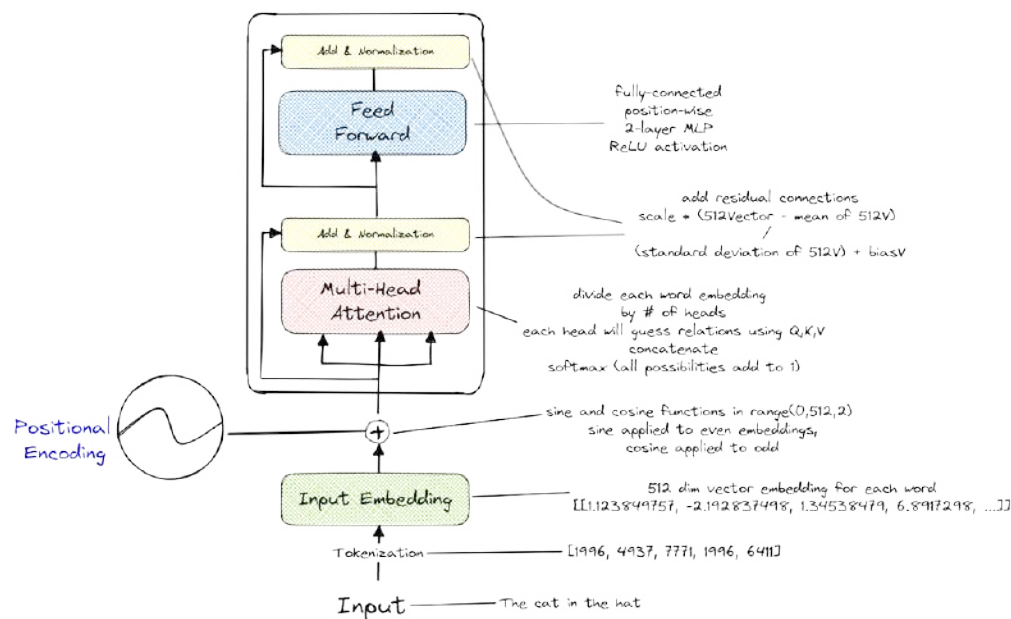
Encoders are the first half of a full transformer model, excelling in areas like classification and feature engineering. One thing Vaswani et al. (2017) figured out is that after the embedding layer inside the encoder, any additional transformations done to the tensors could end up harming their ability to be compared “semantically,” which was the point of the embedding layer. These models rely heavily upon self-attention and a clever positional encoding to manipulate those vectors without significantly decreasing the similarity expressed.

Again, a key thing about embeddings: they are vector representations of data, in our case tokens. Tokens are whatever you pick to represent language. We recommend subwords as a general rule, but you will get a feel for which types of tokens work well where. Consider the sentence, “The cat in the hat rapidly leapt above the red fox and the brown unmotivated dog.” “Red,” and “brown,” should be semantically similar, and they are similarly represented after the embedding layer, but they fall on positions 10 and 14 respectively in the utterance, assuming that we’re tokenizing by word, therefore the positional encoding puts distance between them. However, once the sine and cosine functions[3] are applied, it brings their meaning back to only a little further apart than they were after the encoding, and this encoding mechanism scales brilliantly with recurrence and more data. To illustrate, let’s say there was a 99% cosine similarity between [red], and [brown] after embedding. Encoding would drastically reduce that, to around 85-86% similarity. Applying sine and cosine methodologies as described brings their similarity back up to around 96%.

BERT was one of the first architectures to come after the original paper and are examples of encoder-only transformers. BERT is an incredibly powerful model architecture for how small it is that it is still used in production systems today. BERT was the first encoder-only transformer to surge in popularity, showcasing that performing continuous or sequential (they’re the same) modeling using a transformer results in much better embeddings than Word2Vec. We can see that these embeddings were better because they could

be very quickly applied to new tasks and data with minimal training, with human-preferred results over Word2Vec embeddings. This resulted in most people using BERT-based models for few-shot learning tasks on smaller datasets for a while. BERT puts state-of-the-art performance within arms reach for most researchers and businesses with minimal effort required.

**Figure 2.5 An Encoder, visualized.** Encoders are the first half of the full transformer architecture, and excel in NLU tasks like classification or NER. Encoder models improve upon previous designs by not requiring any priors or recurrence, and use clever positional encoding and multihead attention.



### Strengths:

- Classification and hierarchical tasks showcasing understanding
- Blazing fast, considering the long-range dependency modeling
- Builds off of known models, CBoW in Embedding, MLP in Feed Forward, etc.
- Parallel

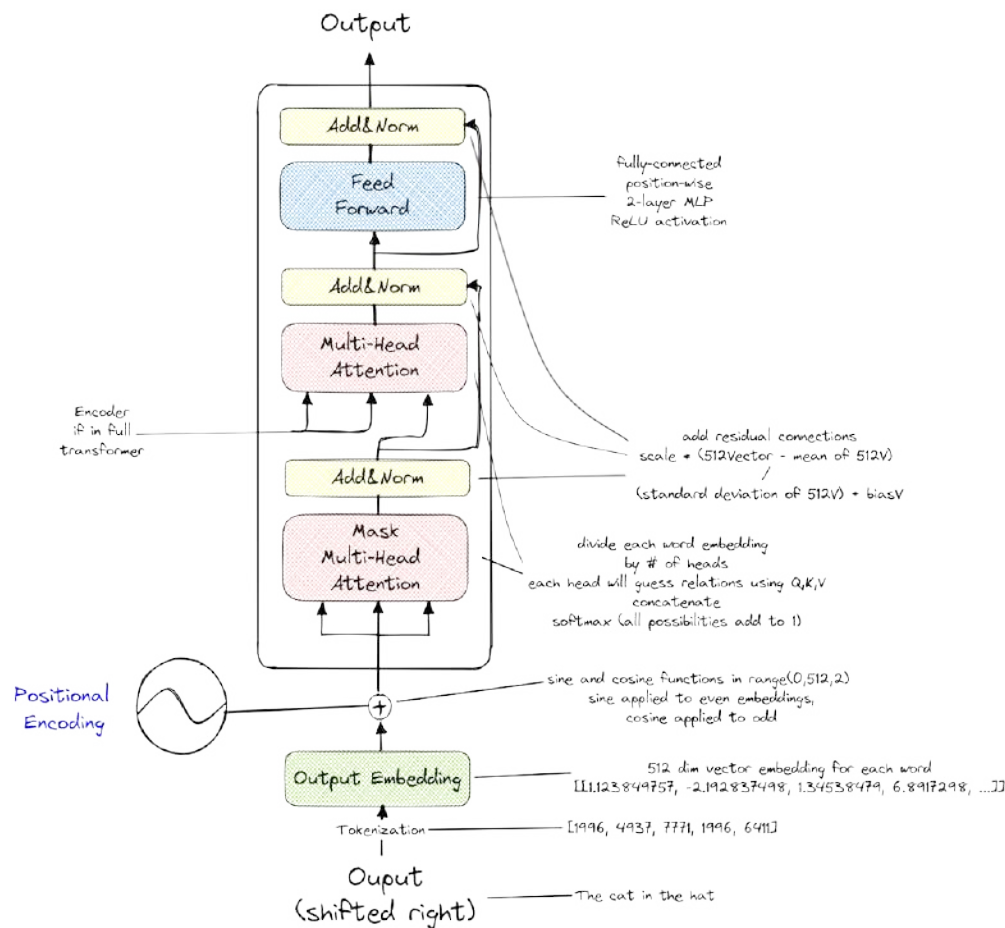
### Weaknesses:

- As suggested, requires lots of data to be effective (although less than RNNs)
- Even more complex architecture

## 2.3.2 Decoders

Decoder models, as shown below, are larger versions of encoders that have 2 multi-head attention blocks and 3 sum and normalize layers in their base form. They are the 2nd half of a transformer behind an encoder. This results in a model that is very good at masked language modeling and learning and applying syntax super quickly leading to the almost immediate idea that decoder-only models are needed to achieve Artificial General Intelligence. A useful reduction of encoder vs decoder tasks is that encoders excel in natural language understanding (NLU) tasks, while decoders excel in natural language generation (NLG) tasks. An example of decoder-only transformer architectures is the Generative Pretrained Transformer (GPT) family of models. These models follow the logic of transformational generative grammar being completely syntax-based, allowing for infinite generation of all possible sentences in a language.[\[4\]](#)

**Figure 2.6 a decoder visualized. Decoders are the second half of a full transformer, and they excel in NLG tasks like chatbots and storytelling. Decoders improve upon previous architectures in the same way as encoders, but they add shifting their output one space to the right for next-word generation to help utilize the advantages of multihead self-attention.**



### Strengths:

- Generating the next token in a sequence (shifted right means taking already-generated tokens into account)
- Building off of both known models and also encoders
- Can be streamed during generation for great UX

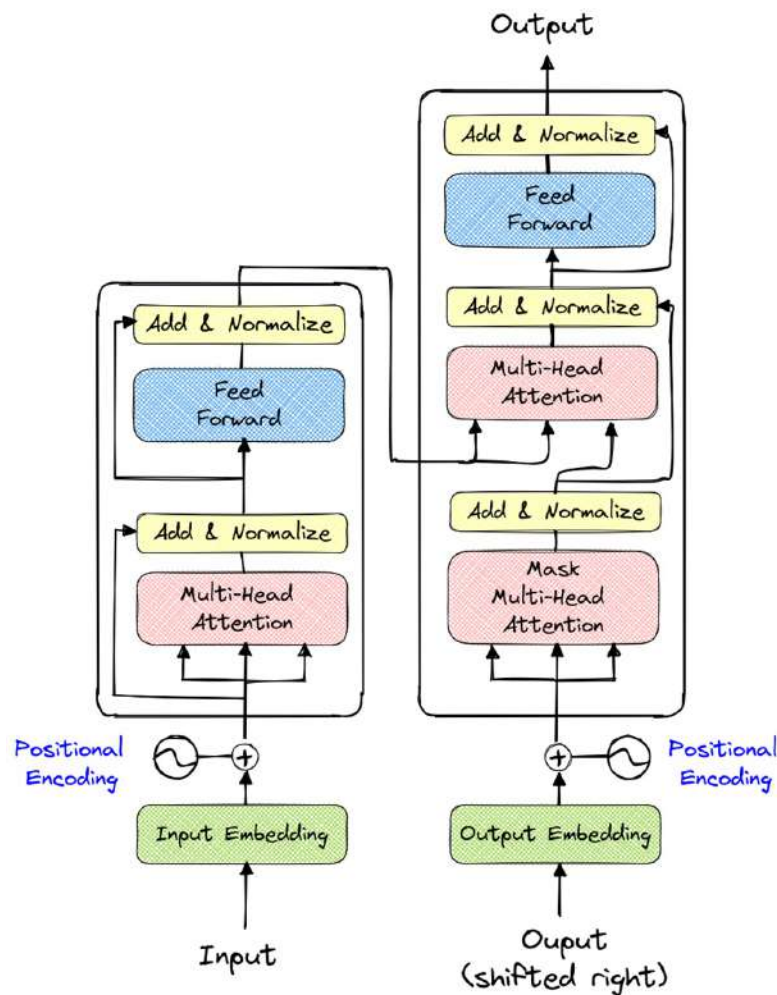
### Weaknesses:

- Syntax only models can often struggle to insert the expected or intended meaning (see all "I force an AI to watch 1000 hours of x and generated" memes from 2018-present)
- Hallucinations

## 2.3.3 Transformers

The full transformer architecture takes advantage of both encoders and decoders, passing the understanding of the encoder into the second Multi-Head Attention block of the decoder before giving output. As each piece of the transformer has a specialty in either understanding or generation, it should feel intuitive for the full product to be best at conditional generation tasks like translation or summarization, where some level of understanding is required before generation occurs. Encoders are geared towards processing input at a high level, and decoders focus more on generating coherent output, the full transformer architecture can successfully understand, then generate based on that understanding. Transformer models have an advantage in that they are built around parallelization, which adds speed that can't currently be replicated in LSTMs. If LSTMs ever get to a point where they can run as quickly as transformers, they may become competitive in the state-of-the-art field. The Text-To-Text Transfer Transformer (T5) family of models are examples of transformers.

**Figure 2.7 A full transformer visualized. A full transformer combines both the encoder and the decoder and does well on all of the tasks of each, as well as conditional generation tasks such as summarization and translation. Because transformers are bulkier and slower than each of their halves, researchers and businesses have generally opted to use those halves over the whole thing, despite the speed and memory boosts being minimal.**



### Strengths:

- Both an encoder and decoder, so is good at everything each of those are good at
- Highly parallelized for speed and efficiency

### Weaknesses:

- Memory intensive, but still less than LSTMs of the same size
- Requires large amounts of data and VRAM for training

As you've probably noticed, most of the models we've discussed aren't at all linguistically focused, being heavily syntax-focused, if attempting to model real language at all. Models, even state-of-the-art transformers only have semantic approximations, no pragmatics, no phonetics, and only really utilize



a mathematical model of morphology during tokenization without context. This doesn't mean the models can't learn these, nor does it mean that, for example, transformers can't take audio as an input, just that the average usage doesn't. With this in mind, it is nothing short of a miracle that they work as well as they do, and they really should be appreciated as such.

Through this chapter so far we've attempted to highlight where the current limitations are in models, and we will dive into where to go to improve upon them in the rest of this book. One such route is one that's already been and being explored to great success, transfer learning and finetuning large foundational models. This technique came about soon after BERT's initial release, when researchers discovered that although BERT performed generally well on a large number of tasks, if they wanted it to perform better on a particular task or data domain, all they needed to do was retrain the model on data representative of the task or domain, but not from scratch. Take all of the pretrained weights that BERT learned while creating the semantic approximation embeddings on a much larger dataset, then significantly less data is required to get state-of-the-art (SotA) performance on the portion that you need. We've seen this with BERT, and with the GPT family of models as they've come out respectively, and now we're seeing it again to solve exactly the challenges brought up: semantic approximation coverage, domain expertise, availability of data.

## 2.4 Really Big Transformers

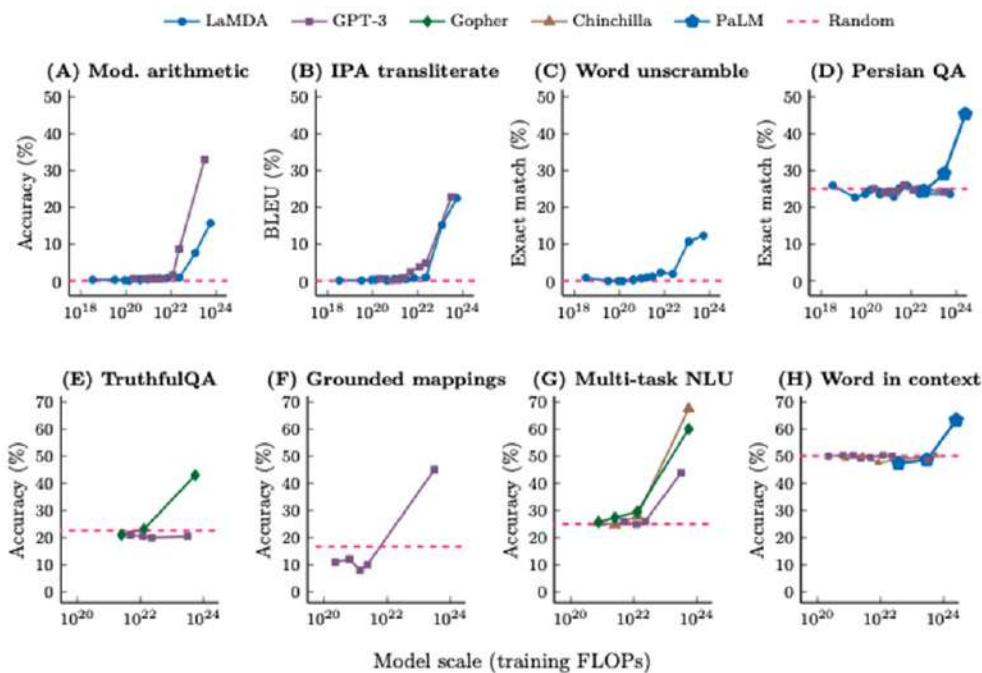
Enter the Large Language Model. Since their introduction transformer based models have continued to only get larger and larger, and not just by their size and number of parameters, but also the size of their training datasets and training cycles has gotten larger and longer as well. If you ever studied machine learning or deep learning during the 2010s, you likely heard the moniker, "more layers doesn't make the model better." LLMs prove this both wrong and right. Wrong because their performance is unparalleled, oftentimes even matching smaller models that have been meticulously finetuned on a particular domain and dataset, even the ones trained on proprietary data. Right because of the challenges that come with both training and deploying them.

One of the major differences between LLMs and LMs involves transfer learning and finetuning. Exactly the same as the previously-large LMs, LLMs are pretrained on massive text corpora, enabling them to learn general language features and representations that can be finetuned for specific tasks. Because LLMs are so massive though and their training datasets so large LLMs are able to achieve better performance with less labeled data, which was a significant limitation of earlier language models. Often times you can finetune an LLM to do highly specialized tasks with only a dozen or so examples.

However, what really makes LLMs powerful and has opened the door to widespread business use cases is their ability to do specialized tasks without any finetuning, but just simple prompting. Just give a few examples of what you want in your query and the LLM is able to produce results. This is called few-shot prompting when it's trained on smaller labeled data sizes, one-shot, when given only one example, and zero-shot, when the task is totally novel. LLMs, especially those trained using RLHF and prompt engineering methodologies, can perform few-shot learning on a whole new level, where they can generalize and solve tasks with only a few examples. This ability is a significant advancement over earlier models that required extensive finetuning or large amounts of labeled data for each specific task.

LMs previously have shown promise in the few and zero-shot learning domains, and LLMs have proven that promise to be true. As models have gotten larger we find they are capable of accomplishing new tasks where smaller models can't. We call this emergent behaviors[5] and figure 2.8 demonstrates eight different tasks that LMs couldn't perform better than random, then suddenly once the models got large enough they could.

**Figure 2.8 Examples of LLMs demonstrating emergent behaviors when tasked with few-shot prompting tasks after the model scale reaches a certain size.**



LLMs have demonstrably great Zero-Shot capabilities as well, which is both due to their vast parameter sizes, and also the main reason for their popularity and viability in the business world. LLMs also exhibit improved handling of ambiguity due to their large size and capacity. They are better at disambiguating words with multiple meanings and understanding the nuances of language, resulting in more accurate predictions and responses. This isn't because of an improved ability or architecture as they share their architecture with smaller transformers, but because they have vastly more examples of how people generally disambiguate. LLMs therefore respond with the same disambiguation as is generally represented in the dataset. Thanks to the diverseness of the text data LLMs are trained on, they exhibit increased robustness in handling various input styles, noisy text, and grammatical errors.

Another key difference between LLMs and LMs is input space. A larger input space is important since it makes few-shot prompting tasks that much more viable. Many LLMs have max input sizes of 8000+ tokens (GPT-4 sports 128k since November 2023, originally 32k), and while all the models previously discussed in the chapter could also have input spaces that high, they generally aren't considered to. We have recently seen a boom in this field as well, with techniques like Recurrent Memory Transformer (RMT)

allowing 1,000,000+ token context spaces, which rocket LLMs even more towards proving that bigger models really are always better. LLMs are designed to capture long-range dependencies within text, allowing them to understand context more effectively than their predecessors. This improved understanding enables LLMs to generate more coherent and contextually relevant responses in tasks like machine translation, summarization, and conversational AI.

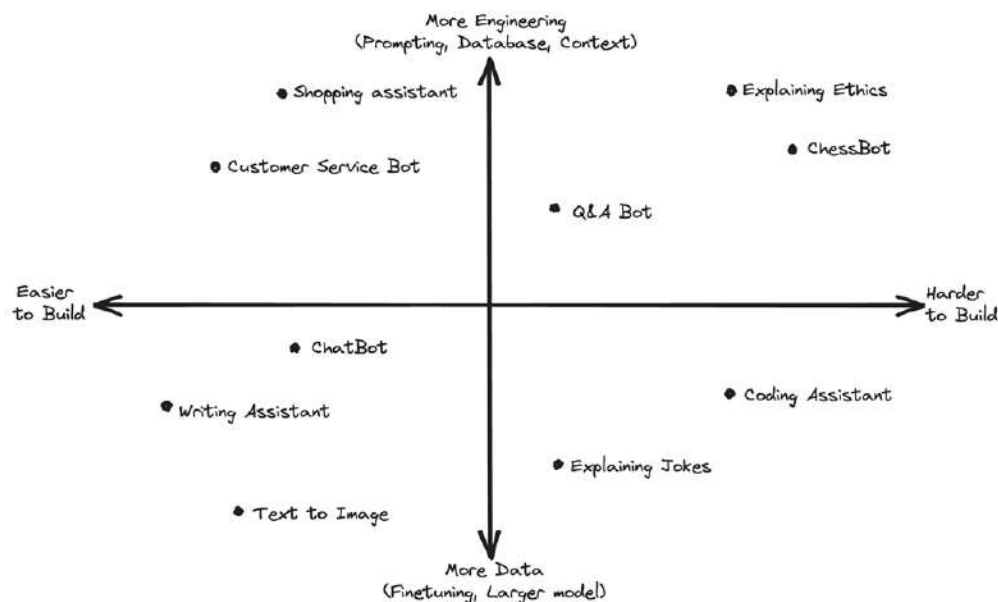
LLMs have revolutionized NLP by offering powerful solutions to problems that were challenging for earlier language models. They bring substantial improvements in contextual understanding, transfer learning, and few-shot learning. As the field of NLP continues to evolve, researchers are actively working to maximize the benefits of LLMs while mitigating all potential risks. Because a better way to approximate semantics hasn't been found, they make bigger and more dimensional approximations. Because a good way of storing pragmatic context hasn't been found, LLMs often allow inserting context either into the prompt directly, into a part of the input set aside for context, or even through sharing of databases with the LLM at inference. This doesn't create pragmatics or a pragmatic system within the models, same as embeddings don't create semantics, but it allows the model to correctly generate syntax that mimics how humans respond to those pragmatic and semantic stimuli. Phonetics is a place where LLMs could likely make gigantic strides, either as completely text-free models, or as a text-phonetic hybrid model, maybe utilizing IPA in addition to or instead of text. It is exciting to consider the possible developments that we are watching sweep this field right now.

At this point you should have a pretty good understanding of what LLMs are and some key principles of linguistics that will come in handy when putting LLMs in Production. Mainly, you should be able to now start reasoning what type of products will be easier or harder to build. Consider figure 2.9, tasks in the lower left hand corner like Writing Assistants and ChatBots are LLMs bread and butter. Text generation based on a little context from a prompt are problems that are strictly syntax based, with a large enough model trained on enough data we can do this pretty easily. A Shopping Assistant is pretty similar and rather easy to build as well, however, we are just missing pragmatics. The assistant needs to know a bit more about the world like

products, stores, and prices. With a little engineering we can add this information into a database and give this context to the model through prompting.

On the other end consider a ChessBot. LLMs *can* play chess. They aren't any good. They have been trained on chess games and understand that "E4" is a common first move, but their understanding is completely syntactical. LLMs really only understand that the text they should generate should contain a letter between A and H and a number between 1 and 8. Like the Shopping Assistant, they are missing pragmatics and don't have a clear model of the game of chess. In addition, they are also missing semantics. Encoders might help us understand the words "King" and "Queen" are similar to each other, but they don't really help us understand that "E4" is a great move one moment for one player and that same "E4" move is a terrible move the very next moment for a different player. LLMs are also completely lacking knowledge based on phonetics and morphology for chess as well, but these are not as important for this case. Either way, we hope this exercise will better inform you and your team on your next project.

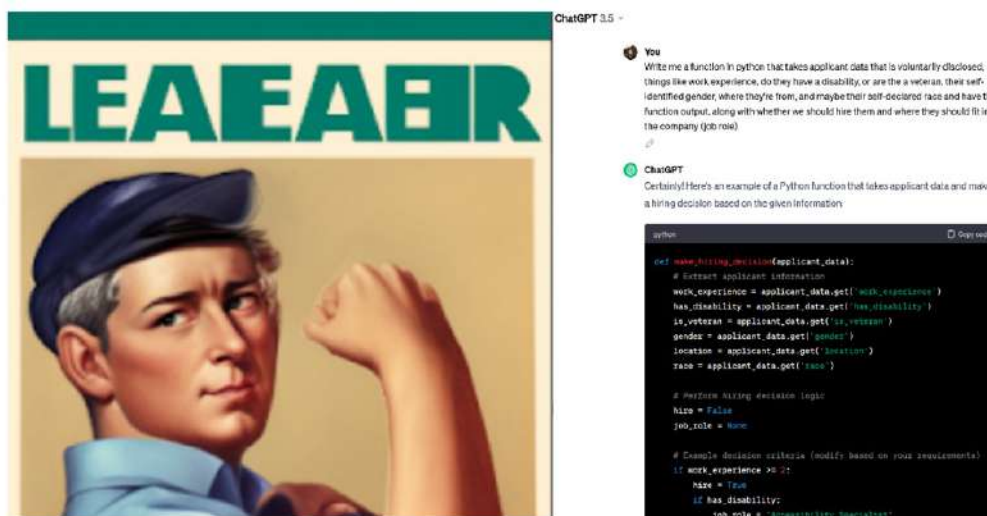
**Figure 2.9** How difficult or easy certain tasks are for LLMs and what approaches to solve them.



LLMs have amazing benefits, but with all of these capabilities come some limitations. Foundational LLMs require vast computational resources for

training, making them less accessible for individual researchers and smaller organizations. This is being remedied with techniques we'll talk about throughout the book like Quantization, Textual Embeddings, Low-Rank Adaptation, Parameter-Efficient Fine Tuning, and Graph Optimization, but foundation models are still currently solidly out of the average individual's ability to train effectively. Beyond that, there are concerns that the energy consumption associated with training LLMs could have significant environmental impact and problems associated with sustainability. This is a complex issue largely out of the scope of this book, but we would be remiss not to bring it up. Last, but not least, since LLMs are trained on large-scale datasets containing real-world text, they may learn and perpetuate biases present in the data, leading to ethical concerns, because real-world people don't censor themselves to provide optimal unbiased data and it's not a widespread practice to really know what data you're training on. For example, if you ask a text-to-image diffusion LLM to generate 1000 images of "leader," 99% of the images feature men, and 95% of the images feature people with white skin. The concern here isn't that men or white people shouldn't be depicted as leaders, but that the model isn't representing the world accurately and it's showing.

**Figure 2.10** Midjourney 5, which is currently the most popular text2img model on the market, when prompted with only one token, "Leader," (Shown Left) changed a well-known popular feminist icon, Rosie the Riveter into a male depiction. ChatGPT (Shown Right) writes a function to place you in your job based on your race, gender, and age. These are examples of unintended outputs.



Sometimes, more nuanced bias is brought out, for example, in the Midjourney example demonstrated in Figure 2.10 a popular feminist icon “Rosie the Riveter,” without being prompted at all (the only prompt given to the model was the word “leader”) was changed to a man. The model didn’t think about this change at all, it just determined during its sampling steps that the prompt “leader,” had more male-looking depictions in the training set. Many people will argue about what “good,” and “bad,” mean in this context, and instead of going for a moral ought, we’ll talk about what accurate means. LLMs are trained on a plethora of data with the purpose of returning the most accurate representations possible. When they are still unable to return accurate representations, especially with their heightened abilities to disambiguate, we can view that as bias that is harmful to the model’s ability to fulfill its purpose. Later we will discuss techniques to combat harmful bias to allow you as an LLM creator to get the exact outputs that you intend and minimize the number of outputs that you do not intend.

Alright, we’ve been building up to this the entire chapter. Let’s go ahead and run our first LLM! In listing 2.9 we are downloading the Bloom model, one of the first open source LLMs to be created, and generate text! We are using HuggingFace’s Transformers library which takes care of all the heavy lifting for us. Very exciting stuff.

#### **Listing 2.9 Running our first LLM**

```
from transformers import AutoModelForCausalLM, AutoTokenizer

MODEL_NAME = "bigscience/bloom" # change to bloom-3b to test

tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
model = AutoModelForCausalLM.from_pretrained(MODEL_NAME)

prompt = "Hello world! This is my first time running an LLM!"

input_tokens = tokenizer.encode(prompt, return_tensors="pt", padding=True)
generated_tokens = model.generate(input_tokens, max_new_tokens=20)
generated_text = tokenizer.batch_decode(
    generated_tokens, skip_special_tokens=True
)
print(generated_text)
```

Did you try to run it?!? If you did, you probably just crashed your laptop.

Oopsie! Forgive me for a little harmless MLOps hazing, but getting some first-hand experience on how large these models can get and how difficult they can be to run is helpful experience to have. We will be talking more about the difficulties of running LLMs and give you some of the tools you need to actually run this in the next chapter. If you don't want to wait and would like to get a similar but much smaller LLM running change the model name to "bigscience/bloom-3b" and run it again. It should work just fine this time on most hardware.

All in all, LLMs are an amazing technology that allow our imaginations to run wild with possibility, and deservedly so. The number one use case for considering an LLM over a smaller LM is for when few-shot capabilities will come into play for whoever the model will be helping, such as helping a CEO when raising funds or a software engineer when writing code. They have this ability precisely because of their size. The larger number of parameters in LLMs directly enable the ability to generalize over smaller spaces in larger dimensions. In this chapter, we've hit the lesser-known side to LLMs, the linguistic and language modeling side. In the next chapter, we'll cover the other half, the MLOps side, where we dive into exactly how that large parameter size affects the model and the systems designed to support that model and make it accessible to the customers or employees the model is intended for.

## 2.5 Summary

- The five components of linguistics are phonetics, syntax, semantics, pragmatics, and morphology.
  - Phonetics can be added through a multimodal model that processes audio files and is likely to improve LLMs in the future, but current datasets are too small.
  - Syntax is what current models are good at.
  - Semantics is added through the embedding layer.
  - Pragmatics can be added through engineering efforts.
  - Morphology is added in the tokenization layer.
- Language does not necessarily correlate with reality. Understanding the process that people use to create meaning outside of reality is useful to training meaningful (to people) models.



- Proper tokenization can be a major hurdle due to too many <UNK> tokens, especially when it comes to specialized problems like code or math.
- Multilingual processing has always outperformed monolingual processing, even on monolingual tasks without models.
- Each language model type in sequence show a natural and organic growth of the LLM field as more and more linguistic concepts are added and make the models better.
- Language modeling has seen an exponential increase in efficacy, correlating to how linguistics-focused the modeling has been.
- Attention is a mathematical shortcut for solving larger context windows faster and is the backbone of modern architectures - Encoders, Decoders, and Transformers.
  - Encoders improve the semantic approximations in embeddings.
  - Decoders are best at text generation.
  - Transformers combine the two.
- Larger models demonstrate emergent behavior suddenly being able to accomplish tasks they couldn't before.

[1] Vaswani et al 2017 Attention Is All You Need  
<https://arxiv.org/abs/1706.03762>

[2] We didn't go over these because they aren't good for NLP, but they are popular especially in computer vision

[3] Not a math or history book

[4] See Appendix A

[5] J. Wei et al., "Emergent Abilities of Large Language Models," Transactions on Machine Learning Research, Aug. 2022, Available: <https://openreview.net/forum?id=yzkSU5zdwD>

# **3 Large Language Model Operations: Building a platform for LLMs**

## **This chapter covers**

- Overview of Large Language Models Operations
- Deployment challenges
- Large Language Models best practices
- Required Large Language Model infrastructure

As we learned in the last chapter when it comes to transformers and Natural Language Processing (NLP), bigger is better, especially when it's linguistically informed. However, bigger models come with bigger challenges because of their size, regardless of their linguistic efficacy, thus requiring us to scale up our operations and infrastructure to handle these problems. In this chapter we'll be looking into exactly what those challenges are, what we can do to minimize them, and what architecture can be set up to help solve these challenges.

## **3.1 Introduction to Large Language Models Operations**

What is Large Language Models Operations (LLMOps)? Well, since I'm one to focus on practicality over rhetoric, I'm not going to dive into any fancy definitions that you'd expect in a text book, but let me simply say it's Machine Learning Operations (MLOps) that has been scaled to handle LLMs. Let me also say, scaling up is hard. One of the hardest tasks in software engineering. Unfortunately, too many companies are running rudimentary MLOps set-ups, and don't think for a second that they will be able to just handle LLMs. That said, the term "LLMOps," may not be needed. It has yet to show through as sufficiently different from core MLOps, especially

considering they still have the same bones. If this book were a dichotomous key, MLOps and LLMOps would definitely be in the same genus, and only time will tell about whether they are the same species. Of course by refusing to define LLMOps properly, I might have traded one confusion for another, so let's take a minute to describe MLOps.

MLOps is the field and practice of reliably and efficiently deploying and maintaining machine learning models in production. This includes, and indeed requires, managing the entire machine learning lifecycle from data acquisition and model training to monitoring and termination. A few principles required to master this field include workflow orchestration, versioning, feedback loops, Continuous Integration and Continuous Deployment (CI/CD), security, resource provisioning, and data governance. While there are often personnel who specialize in the productionizing of models, often with titles like ML Engineers, MLOps Engineers or ML Infrastructure Engineer, the field is a large enough beast it often abducts many other unsuspecting professionals to work in it who hold titles like Data Scientist or DevOps Engineer—oftentimes against their knowledge or will; leaving them kicking and screaming that “it’s not their job”.

## **3.2 Operations Challenges with Large Language Models**

So why have a distinction at all? If MLOps and LLMOps are so similar, is LLMOps just another fad opportunists throw on their resume? Not quite. In fact, I think it’s quite similar to the term Big Data. When the term was at its peak popularity, people with titles like Big Data Engineer used completely different tool sets and developed specialized expertise that were necessary in order to handle the large datasets. LLMs come with a set of challenges and problems you won’t find with traditional machine learning systems. A majority of these problems extend almost exclusively because they are so big. Large models are large! We hope to show you that LLMs truly earn their name. Let’s take a look at a few of these challenges, so we can appreciate the task ahead of us when we start talking about deploying an LLM.

### **3.2.1 Long download times**

Back in 2017 when I was still heavily involved as a Data Scientist, I decided to try my hand at reimplementing some of the most famous computer vision models at the time AlexNet, VGG19, and ResNet. I figured this would be a good way to reinforce my understanding of the basics with some practical hands-on experience. Plus, I had an ulterior motive, I had just built my own rig with some NVIDIA GeForce 1080 TI GPUs—which was state of the art at the time—and thought this would be a good way to break them in. The first task: download the ImageNet dataset. The ImageNet dataset was one of the largest annotated datasets available containing millions of images rounding out to a file size of a whopping ~150GB! Working with it was proof that you knew how to work with “Big Data ” which was still a trendy word and an invaluable skill set for a data scientist at the time. After agreeing to the terms and gaining access, I got my first wakeup call. Downloading it took an entire week.

Large models are large. I don’t think I can overstate that. You’ll find throughout this book that fact comes with many additional headaches and issues for the entire production process, and you have to be prepared for it. In comparison to the ImageNet dataset, the Bloom LLM model is 330GB, more than twice the size. Most readers I’m guessing haven’t worked with either ImageNet or Bloom, so for comparison Call of Duty: Modern Warfare, one of the largest games at the time of writing is 235 GB. Final Fantasy 15 is only 148 GB, which you could fit two of into the model with plenty of room to spare. It’s just hard to really comprehend how massive LLMs are. We went from 100 million parameters in models like BERT and took them to billions of parameters. If you went on a shopping spree and spent \$20 a second (or maybe just left your AWS EC2 instance on by accident) it’d take you half a day to spend a million dollars; it would take you 2 years to spend a billion.

Thankfully it doesn’t take two weeks to download Bloom because unlike ImageNet, it’s not hosted on a poorly managed University server and it also has been sharded into multiple smaller files to allow downloading in parallel, but it will still take an uncomfortably long time. Consider a scenario where you are downloading the model under the best conditions. You’re equipped with a gigabit speed fiber internet connection and you were magically able to dedicate the entire bandwidth and I/O operations of your system and the server to it, it’d still take over 5 minutes to download! Of course, that’s under

the best conditions. You probably won't be downloading the model under such circumstances, with modern infrastructure you can expect it to take on the order of hours. When my team first deployed Bloom it took an hour and a half to download it. Heck, it took me an hour and half to download The Legend of Zelda: Tears of the Kingdom and that's only 16GB, so I really can't complain.

### **3.2.2 Longer Deploy Times**

Just downloading the model is a long enough time frame to make any seasoned developer shake, but deployment times are going to make them keel over and call for medical attention. A model as big as Bloom can take 30-45 minutes just to load the model into GPU memory, at least those are the time frames my team first saw. Not to mention any other steps in your deployment process that can add to this. Indeed, with GPU shortages, it can easily take hours just waiting for resources to free up—more on that in a minute.

What does this mean for you and your team? Well for starters, I know lots of teams who deploy ML products often simply download the model at runtime. That might work for small sklearn regression models, but it isn't going to work for LLMs. Additionally, you can take most of what you know about deploying reliable systems and throw it out the window (but thankfully not too far). Most modern day best practices for software engineering assume you can easily just restart an application if anything happens, and there's a lot of rigmarole involved to ensure your systems can do just that. But with LLMs it can take seconds to shut down, but potentially hours to redeploy making this a semi-irreversible process. Like picking an apple off a tree, it's easy to pluck one off, but if you bite into it and decide it's too sour, you can't just attach it back onto the tree so it can continue to ripen. You'll just have to wait awhile for another to grow.

While not every project requires deploying the largest models out there, you can expect to see deployment times measured in minutes. These longer deploy times make scaling down right before a surge of traffic a terrible mistake, as well as figuring out how to manage bursty workloads difficult. General CI/CD methodologies need to be adjusted since rolling updates take longer leaving a backlog piling up quickly in your pipeline. Silly mistakes

like typos or other bugs often take longer to notice, and longer to correct.

### **3.2.3 Latency**

Along with increases in model size often come increases in inference latency. This is obvious when stated, but more parameters equates to more computations, and more computations means longer inference wait times. However, this can't be underestimated. I know many people who downplay the latency issues because they've interacted with an LLM chatbot and the experience has felt smooth. Take a second look though, and you'll notice that it is returning one word at a time which is streamed to the user. It feels smooth because the answers are coming in faster than a human can read, but a second look helps us realize this is just a UX trick. LLMs are still too slow to be very useful for an autocomplete solution for example, where responses have to be blazingly fast. Building it into a data pipeline or workflow that reads a large corpus of text and then tries to clean it or summarize it, may also be prohibitively slow to be useful or reliable.

There are also many less obvious reasons for their slowness. For starters, LLMs are often distributed across multiple GPUs, which adds extra communication overhead. As discussed later in this chapter in section 3.3.2 they are distributed in other ways, often even to improve latency, but any distribution adds additional overhead burden. In addition, LLMs latency is severely impacted by completion length, meaning the more words it uses to return a response, the longer it takes. Of course, completion length also seems to improve accuracy. For example, using prompt engineering techniques like Chain of Thought (CoT) we ask the model to think about a problem in a step-by-step fashion which has shown to improve results for logic and math questions but also increases the response length and latency time significantly.

### **3.2.4 Managing GPUs**

To help with these latency issues we usually want to run them in GPUs. If we want to have any success training LLMs we'll need GPUs for that as well, but this all adds additional challenges many underestimate. Most web services and many ML use cases can be done solely on CPUs. Not so with

LLMs. Partly because of GPUs' parallel processing capabilities offering a solution to our latency problems, and partly because of the inherent optimization GPUs offer in the linear algebra, matrix multiplications and tensor operations that's happening under the hood. For many, stepping into the realm of LLMs, this requires utilizing a new resource and extra complexity. Many brazenly step into this world acting like it's no big deal, but they are in for a rude awakening. Most system architectures and orchestrating tooling available like Kubernetes, assume your application will run with CPU and memory alone. While they often support additional resources like GPUs, this is often an afterthought. You'll soon find you'll have to rebuild containers from scratch and deploy new metric systems.

One aspect of managing GPUs most companies are never prepared for is that they tend to be rare and limited. For the last decade it seems that we have gone in and out of a global GPU shortage. They can be extremely difficult to provision for companies looking to stay on premise. I've spent lots of time in my career working with companies who chose to stay on premise for a variety of reasons. One of the things they had in common is that they never had GPUs on their servers. When they did, they were often purposely difficult to access except for a few key employees.

If you are lucky enough to be working in the Cloud a lot of these problems are solved, but there is no free lunch here either. My team has often gone chasing their tail trying to help data scientists struggling to provision a new GPU workspace, running into obscure ominous errors like "scale.up.error.out.of.resources". Only to discover that these esoteric readings indicate all the GPUs of a selected type in the entire region are being utilized and none are available. CPU and Memory can often be treated as infinite in a datacenter, GPU resources, however, cannot. Sometimes you can't expect them at all. Most data centers only support a subset of instance or GPU types. Which means you may be forced to set up your application in a region further away from your user base increasing latency. Of course, I'm sure you can work with your cloud provider when looking to expand your service to a new region that doesn't currently support it, but you might not like what you hear based on timelines and cost. Ultimately, you'll run into shortage issues no matter where you choose to run, on-prem or in the cloud.

### 3.2.5 Peculiarities of Text Data

LLMs are the modern day solution to NLP. NLP is one of the most fascinating branches of ML in general because it primarily deals with text data, which is primarily a qualitative measure. Every other field deals with quantitative data. We have figured out a way to encode our observations of the world into a direct translation of numerical values. For example, we've learned how to encode heat into temperature scales and measure them with thermometers and thermocouples or we can measure pressure with manometers and gauges and put it into pascals.

Computer Vision and the practice of evaluating images is often seen as qualitative, but the actual encoding of images into numbers is a solved problem. Our understanding of light has allowed us to break images apart into pixels and assign them RGB values. Of course this doesn't mean CV is by any means solved, there's still lots of work to do to learn how to identify the different signals in the patterns of the data. Audio data is another that's often considered qualitative. How does one compare two songs? But we can measure sound and speech, directly measuring the sound wave's intensity in decibels and frequency in hertz.

Unlike other fields that encode our physical world into numerical data, text data is looking at ways to measure the ephemeral world. After all, text data is our best effort of encoding our thoughts, ideas and communication patterns. While yes, we have figured out ways to turn words into numbers, we haven't figured out a direct translation. Our best solutions to encode text and create embeddings are just approximations at best, in fact we use machine learning models to do it! An interesting aside to this is that numbers are also text and a part of language. If we want models that are better at math we need a more meaningful way to encode these numbers. Since it's all made up, when we try to encode text numbers into machine-readable numbers we are creating a system attempting to reference itself recursively in a meaningful way. Not an easy problem to solve!

Because of all this, LLMs (and all NLP solutions) have unique challenges. For example, monitoring. How do you catch data drift in text data? How do you measure "correctness"? How do you ensure cleanliness of the data?



These types of problems are difficult to define let alone solve.

### **3.2.6 Token Limits Create Bottlenecks**

One of the big challenges for those new to working with LLMs is dealing with the token limits. The token limit for a model is the maximum number of tokens that can be included as an input for a model. The larger the token limit, the more context we can give the model to improve its success at accomplishing the task. Everyone wants them to be higher, but it's not that simple. These token limits are defined by two problems, the first being the memory and speed our GPUs have access to, and the second being the nature of memory storage in the models themselves.

The first one seems unintuitive, why couldn't we just increase the GPU memory? The answer is complex, we can, but stacking more layers in the GPU to take into account more GB at once slows down the GPU's computational ability as a whole. GPU manufacturers right now are working on new architectures and ways to get around this problem. The second one is more fascinating because we find that increasing the token limits actually just exacerbates the mathematical problems under the hood. Let me explain. Memory storage within an LLM itself isn't something we think about often. We call that mechanism Attention, which we discussed in depth in section 2.2.7. What we didn't discuss was that Attention is a quadratic solution—as the number of tokens increase the number of calculations required to compute the attentions scores between all the pairs of tokens in a sequence scales quadratically with the sequence length. In addition, within our gigantic context spaces and since we are dealing with quadratics, we're starting to hit problems where the only solutions involve imaginary numbers which is something that can cause models to behave in unexpected ways. This is likely one of the reasons why LLMs hallucinate.

These problems have real implications and impact application designs. For example, when my team upgraded from GPT3 to GPT4 we were excited to have access to a higher token limit, but we soon found this led to longer inference times and subsequently a higher timeout error rate. In the real world, it's often better to get a less accurate response quickly than to get no response at all because the promise of a more accurate model often is just

that, a promise. Of course, deploying it locally where you don't have to worry about response times you'll likely find your hardware a limiting factor. For example, LLaMA was trained with 2048 tokens but you'll be lucky to take advantage of more than 512 of that when running with a basic consumer GPU as you are likely to see Out-of-Memory (OOM) errors or even the model simply just crashing.

A gotcha, which is likely to catch your team by surprise and should be pointed out now is that different languages have different tokens per character. Take a look at Table 3.1, where we compare converting the same sentence in different languages to tokens using OpenAI's cl100k\_base Byte Pair Encoder. Just a quick glance reveals that LLMs typically favor the English language in this regard. In practice, this means if you are building a chatbot with an LLM, your English users will have greater flexibility in their input space than Japanese users leading to very different user experiences.

**Table 3.1 Comparison of Token Counts in different languages**

| Language             | String  | Characters | Tokens |
|----------------------|---|------------|--------|
| English              | The quick brown fox jumps over the lazy dog               | 43         | 9      |
| French               | Le renard brun rapide saute par-dessus le chien paresseux | 57         | 20     |
| Spanish              | El rápido zorro marrón salta sobre el perro perezoso      | 52         | 22     |
| Japanese             | 素早い茶色のキツネが怠惰な犬を飛び越える                                      | 20         | 36     |
| Chinese (simplified) | 敏捷的棕色狐狸跳过了懒狗  | 12         | 28     |

If you are curious as to why this is, it is due to text encodings, which is just another peculiarity of working with text data as discussed in the previous section. Consider Table 3.2 where we show several different characters and their binary representation in UTF-8. English characters can almost exclusively be represented with a single byte being included in the original ASCII standard computers were originally built on, while most other

characters require 3 or 4 bytes. Because it takes more memory it also takes more token space.

**Table 3.2 Comparison of byte lengths for different currency characters in UTF-8.**

| Character | Binary UTF-8                           | Hex UTF-8           |
|-----------|--|---------------------|
| \$        | 00100100                               | 0x24                |
| £         | 11000010 10100011                      | 0xc2 0xa3           |
| ¥         | 11000010 10100101                      | 0xc2 0xa5           |
| €         | 11100010 10000010 10100000             | 0xe2 0x82 0xa0      |
| Ⓢ         | 11110000 10011111 10010010<br>10110000 | 0xf0 0x9f 0x92 0xb0 |

Increasing the token limits has been an ongoing research question since the popularization of transformers, and there are some promising solutions still in research phases like Recurrent Memory Transformers (RMT)[1]. We can expect to continue to see improvements in the future and hopefully this will become naught but an annoyance.

### 3.2.7 Hallucinations Cause Confusion

So far we've been discussing some of the technical problems a team faces when deploying an LLM into a production environment, but nothing compares to the simple problem that LLMs tend to be wrong. They tend to be wrong a lot. Hallucinations is a term coined to describe occurrences when LLM models will produce correct sounding results that are wrong. For example, book references or hyperlinks that have the form and structure of what would be expected, but are nevertheless, completely made up. As a fun example I asked for books on LLMs in Production from the publisher Manning (a book that doesn't exist yet since I'm still writing it). I was given the following suggestions: Machine Learning Engineering in Production by Mike Del Balso and Lucas Serveén which could be found at <https://www.manning.com/books/machine-learning-engineering-in-production> and Deep Learning for Coders with Fastai and PyTorch by Jeremy

Howard and Sylvain Gugger which could be found at <https://www.manning.com/books/deep-learning-for-coders-with-fastai-and-pytorch>. The first book is entirely made up. The second book is real however it's not published by Manning. In each case the internet addresses are entirely made up. These URLs are actually very similar to what you'd expect in format if you were browsing Mannings website, and should return 404 errors if you visit them.

One of the most annoying aspects of hallucinations is that they are often surrounded by confident sounding words. LLMs are terrible at expressing uncertainty, in large part because of the way they are trained. Consider the case " $2+2=$ ". Would you prefer it to respond, "I think it is 4" or just simply "4"? Most would prefer to simply get the correct "4" back. This bias is built in as models are often given rewards for being more correct or at least sounding like it.

There are various explanations as to why hallucination occurs, but the most truthful answer is that we don't know if there's just one cause. It's likely a combination of several things, thus there isn't a good fix for it yet. Nevertheless, being prepared to counter these inaccuracies and biases of the model are crucial to provide the best user experience for your product.

### **3.2.8 Bias and Ethical Considerations**

Just as concerning as the model getting things wrong is when it gets things right in the worst possible way. For example, allowing it to encourage users to commit suicide[2], teaching your users how to make a bomb[3], or participating in sexual fantasies involving children[4]. These are extreme examples, but prohibiting the model from answering such questions is undeniably vital to success.

LLMs are trained on vast amounts of text data which is also their primary source of bias. Because we've found that larger datasets are just as important as larger models in producing human-like results, most of these datasets have never truly been curated or filtered to remove harmful content, instead choosing to prioritize size and a larger collection. Cleaning the dataset is often seen as prohibitively expensive, requiring humans to go in and

manually verify everything, but there's a lot that could be done with simple regular expressions and other automated solutions. By processing these vast collections of content and learning the implicit human biases, these models will inadvertently perpetuate them. These biases range from sexism and racism to political preferences and can cause your model to inadvertently promote negative stereotypes and discriminatory language.

### **3.2.9 Security Concerns**

As with all technology we need to be mindful of security. LLMs have been trained on large corpus of text and some of it could be harmful or sensitive and shouldn't be exposed so steps should be taken to protect this data from being leaked. The bias and ethical concerns from the last section are good examples of conversations you don't want your users to be having, but you could also imagine fine-tuning a model on your company's data and potentially have secrets lost inadvertently if proper precautions aren't taken.

One should be aware that LLMs are susceptible to adversarial attacks like prompt injections. Prompt injections are attacks done by a user to trick the LLM to ignore instructions given to it and generate undesired content. For example, if you ask ChatGPT what its gender is it appropriately replies that as an AI language model, it doesn't have a gender, however, with clever prompting you may be able to bypass these protocols and get it to reveal one. While this example is harmless, I've seen others successfully extract API keys and other secrets from an LLM, run code in non-protected environments, steal environment variables and traverse local file systems where the model is served. Not to mention the plethora of examples of users using prompting to jailbreak or bypass protocols put in place for ethical considerations outlined in the section above. An interesting aside to this, LLMs are good at inventing fake secrets! Even successful prompt injection attacks can often fail due to LLM hallucinations which can have funny consequences.

For example, using a simple prompt to ask for Coca-Cola's secret recipe simply returned the boring result, "I cannot divulge trade secrets or any proprietary information." Pretty lame if you ask me. With a little trickery and some prompt injection I was able to get it to spill its secrets which of course

my coauthor and I decided to try—after confirming with a food scientist it wouldn't poison us. After following the instructions exactly, I was surprised to find we actually got a pretty tasty drink, but it tasted nothing like Cola. While the recipe looks legitimate, nuances like quantities, cook times, and even the need to strain the syrup were all off. Guess, I'll have to hold off taking over the beverage industry until I get the real recipe.

### **Coca-Cola Recipe Hallucinated from Prompt Injection**

#### Homemade Coca-Cola Recipe

=====

**\*\*Disclaimer\*\*:** The original Coca-Cola recipe is a highly guarded trade secret. The following recipe is an approximation for creating a homemade cola drink that resembles the taste of Coca-Cola.

#### Ingredients:

- 3 cups of water
- 2 cups granulated sugar
- 1/8 tsp food-grade citric acid
- 3/4 tsp caramel color (liquid)
- 1/2 tsp caffeine (powder or liquid form)
- 1 1/2 tsp cola extract
- 1/8 tsp vanilla extract
- 1 1/2 cups carbonated water (store-bought or using a home soda maker)
- A pinch of ground cinnamon (optional)
- A pinch of ground nutmeg (optional)

#### Instructions:

1. In a large saucepan, bring the 3 cups of water to a gentle simmer.
2. Slowly stir in the granulated sugar, ensuring that it fully dissolves.
3. Add the citric acid, caramel color, caffeine, cola extract, and vanilla extract to the saucepan. Mix well.
4. Stir in the optional spices if desired: cinnamon and nutmeg.
5. Simmer the mixture for about 20 minutes, stirring occasionally.

6. Remove the saucepan from the heat and let it cool to room temperature.
7. Once cooled, strain the syrup to remove any undissolved solids.
8. To make the final drink, mix 1 part syrup with 3 parts carbonated water. Adjust to taste depending on your preference for sweetness and carbonation.

#### Variations:

- You can experiment with different amounts of caramel color, caffeine, and cola extract to adjust the flavor and appearance of your homemade cola.
- To create different cola flavor profiles, experiment with adding additional spices (such as star anise, orange or lemon peel, or cardamom), essential oils, or fruit extracts.
- Consider using a natural sweetener, like honey or agave syrup, as an alternative to granulated sugar.

#### Historical context and interesting facts:

- Coca-Cola was first introduced in 1886 by pharmacist Dr. John S. Pemberton, who aimed to create a unique and refreshing beverage that could also relieve fatigue and headaches.
- The original Coca-Cola formula contained trace amounts of cocaine extracted from the coca leaf, which was later replaced with caffeine in the early 20th century.
- In 1985, Coca-Cola introduced "New Coke," a reformulation aimed at competing with Pepsi. However, the taste change was met with backlash, and the company returned to its original formula just three months later, rebranded as "Coca-Cola Classic."

Result(s): You can now create your own homemade version of cola, resembling the taste of Coca-Cola. Enjoy and feel free to experiment with different spices, extracts, and sweeteners to create your own unique flavor profile.

Another security concern involves pickle injections. Pickle is a library in Python to serialize objects and often used to serialize ML models. It serializes them into a byte-stream that contains opcodes that are executed one-by-one

as it is deserialized. It's a fast and easy way to share large objects. Pickle injections involve the process of corrupting this byte-stream, often injecting malware over the wire when the model is transferred over an insecure network. This is especially concerning for large models that take a long time to download, as it makes it easier for a third party to intercept the transfer and inject malicious code. If this happens, the code that is injected can potentially give the attackers access to your system. This can happen when attempting to use the model during inference, as the harmful code will execute if it is not detected and properly removed. It is important to take precautions such as using secure networks and verifying the integrity of the model before use to prevent this type of attack.

### **3.2.10 Controlling Costs**

Working with LLMs involves various cost-related concerns. The first as you probably gathered by now is infrastructure costs, which include high-performance GPUs, storage, and other hardware resources. We talked about how GPUs are harder to procure, that also unfortunately means they are more expensive. Mistakes like leaving your service on have always had the potential to rack up the bills, but with GPU's in the mix this type of mistake is even more deadly. These models also demand significant computational power, leading to high energy consumption during both training and inference. On top of all this, their longer deploy times means we are often running them even during low traffic to handle bursty workloads or anticipated future traffic. Overall this leads to higher operational costs.

Additional costs include, managing and storing vast amounts of data used to train or fine-tune as well as regular maintenance, such as model updates, security measures, and bug fixes, can be financially demanding. As with any technology used for business purposes, managing potential legal disputes and ensuring compliance with regulations is a concern. Lastly, investing in continuous research and development to improve your models and give you a competitive edge will be a factor.

We talked a bit about the technical concerns when it comes to token limits, and these are likely to be solved, but what we didn't discuss was the cost limitations as most API's charge on a token basis. This makes it more



expensive to send more context and use better prompts. It also makes it a bit harder to predict costs since while you can standardize inputs, you can't standardize outputs. You never can be too sure how many tokens will be returned, making it difficult to govern. Just remember with LLMs, it is as important as ever to ensure proper cost engineering practices are implemented and followed to ensure costs never get away from you.

## **3.3 Large Language Model Operations Essentials**

Now that we have a handle of the type of challenge we are grappling with, let's take a look at all the different LLMOps practices, tooling, and infrastructure to see how different components help us overcome these obstacles. First off, let's dive into different practices starting with compression where we will talk about shrinking, trimming, and approximating to get models as small as we can. We will then talk about distributed computing which is needed to actually make things run since the models are so large they rarely fit into a single GPU's memory. After we are finished with those we will venture into infrastructure and tooling needed to make it all happen in the next section.

### **3.3.1 Compression**

As you were reading about the challenges of LLMs in the last section, you might have asked yourself something akin to, "If the biggest problems from LLMs come from their size, why don't we just make them smaller?" If you did, congratulations! You are a genius and compression is the practice of doing just that. Compressing models as small as we can make them will improve deployment time, reduce latency, scale down the number of expensive GPUs needed, and ultimately save money. However, the whole point of making the models so stupefyingly gargantuan in the first place was because it made them better at what they do. We need to be able to shrink them, without losing all the progress we made by making them big in the first place.

This is far from a solved problem, but there are multiple different ways to approach the problem with different pros and cons to each method. We'll be talking about several of the methods, starting with the easiest and most

effective.

## Quantizing

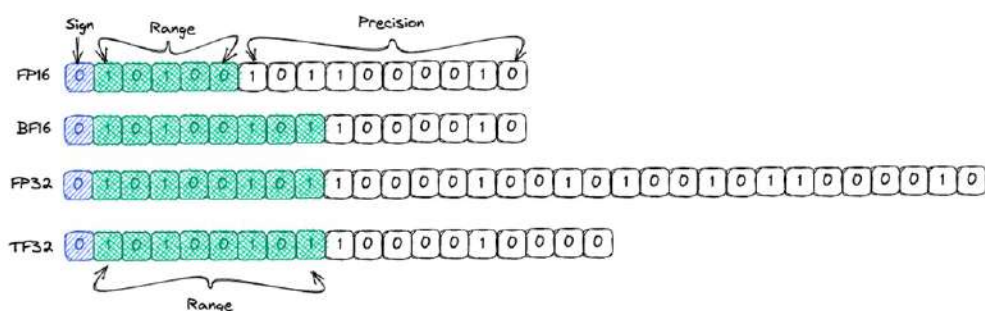
Quantizing is the process of reducing precision in preference of lowering the memory requirements. This tradeoff makes intuitive sense. When I was in college, we were taught to always round our numbers to the precision of your tooling. Pulling out a ruler and measuring my pencil, you wouldn't believe me if I stated the length was 19.025467821973739cm. Even if I used a caliper, I couldn't verify a number so precisely. With our ruler, any number beyond 19.03cm is fantasy. To drive the point home, one of my engineering professors once told me, "If you are measuring the height of a skyscraper, do you care if there is an extra sheet of paper at the top?"

How we represent numbers inside computers often leads us to believe we have better precision than we actually do. To drive this point home, open up a Python terminal and add  $0.1 + 0.2$ . If you've never tried this before, you might be surprised to find this doesn't equal 0.3, but 0.30000000000000004. I'm not going to go into the details of the math behind this phenomenon, but the question stands, can we reduce the precision without making things worse? We really only need precision to the tenth decimal, but reducing the precision is likely to get us a number like 0.304 rather than 0.300 which would increase our margin of error.

Ultimately, the only numbers a computer understands are 0 and 1, on or off, a single bit. To improve this range we combine multiple of these bits together and assign them different meanings. String 8 of them together and you get a byte. Using the int8 standard we can take that byte and encode all the integers from -128 to 127. To spare you the particulars because I assume you already know how binary works, suffice it to say the more bits we have the larger range of numbers we can represent, both larger and smaller. Figure 3.1 shows a few common floating point encodings. With 32 bits strung together we get what we pretentiously term full precision, and that is how most numbers are stored including the weights in machine learning models. Basic quantization moves us from full precision to half precision, shrinking models to half their size. There are actually two different half precision standards, FP16 and BF16 which differ in how many bits represent the range or exponent part.

Since BF16 uses the same number of exponents as FP32, it's been found to be more effective for quantizing and you can generally expect almost exactly the same level accuracy for half the size of model. If you understood the paper and skyscraper analogy above it should be obvious why.

Figure 3.1 shows the bit mapping for a few common floating point encodings. 16-bit float or half precision (FP16), bfloat 16 (BF16), 32-bit float or single full precision (FP32), and NVIDIA's TensorFloat (TF32)



However, there's no reason to stop there. We can often push it another byte down to 8-bit formats without too much loss of accuracy. There have even already been successful research attempts showing selective 4-bit quantization of portions of LLMs are possible with only fractional loss of accuracy. The selective application of quantization is a process known as dynamic quantization and is usually done on just the weights, leaving the activations in full precision to reduce accuracy loss.

The holy grail of quantizing though would be int2, representing every number as -1, 0, or 1. This currently isn't possible without completely degrading the model, but would make the model up to 8 times smaller. The Bloom model would be a measly ~40GB, small enough to fit on a single GPU. This is of course, as far as quantizing can take us and if we wanted to shrink further we'd need to look at additional methods.

The best part of quantization though is that it is easy to do. There are many frameworks that allow this, but in Listing 3.1 I demonstrate how to use pytorch's quantization library to do a simple post training static quantization (PTQ). All you need is the full precision model, some example inputs, and a validation dataset to prepare and calibrate with. As you can see it's only a few lines of code.

### Listing 3.1 Example PTQ in PyTorch

```
import copy
import torch.nn.quantization as q

# deep copy the original model as quantization is done in place
model_to_quantize = copy.deepcopy(model_fp32)
model_to_quantize.eval()

# get mappings - note use "qnnpack" for ARM and "fbgemm" for x86
qconfig_mapping = q.get_default_qconfig_mapping("qnnpack")

# prepare
prepared_model = q.prepare(model_to_quantize)

# calibrate - you'll want to use representative (validation) data
with torch.inference_mode():
    for x in dataset:
        prepared_model(x)

# quantize
model_quantized = q.convert(prepared_model)
```

Static PTQ is the most straightforward approach to quantizing, done after the model is trained and quantizing all the model parameters uniformly. As with most formulas in life, the most straightforward approach introduces more error. Oftentimes this error is acceptable, but when it's not we can add extra complexity to reduce the accuracy loss from quantization. Some methods to consider are uniform vs non-uniform, static vs dynamic, symmetric vs asymmetric, and applying it during or after training.

To understand these methods let's consider the case where we are quantizing from FP32 to INT8. In FP32 we essentially have the full range of numbers at our disposal, but in INT8 we only have 256 values, we are trying to put a genie into a bottle and it's no small feat. If you study the weights in your model, you might notice that the majority of the numbers are fractions between  $[-1, 1]$ . We could take advantage of this by then using an 8-bit standard that represents more values in this region in a non-uniform way instead of the standard uniform  $[-128, 127]$ . While mathematically possible, unfortunately, any such standards aren't commonplace and modern day deep learning hardware and software are not designed to take advantage of this. So for now, it's best to just stick to uniform quantization.

The simplest approach to shrink the data is to just normalize it, but since we are going from a continuous scale to a discrete scale there are a few gotchas, so let's explore those. First, we start by taking the min and max and scale them down to match our new number range, we would then bucket all the other numbers based on where they fall. Of course, if we have really large outliers, we may find all our other numbers squeezed into just one or two buckets completely ruining any granularity we once had. To prevent this, we can simply clip any large numbers. This is what we do in static quantization. However, before we clip the data, what if we chose a range and scale that captures the majority of our data beforehand? We need to be careful, since if this dynamic range is too small, we will introduce more clipping errors, if it's too big, we will introduce more rounding errors. The goal of dynamic quantization of course is to reduce both errors.

Next we need to consider the symmetry of the data. Generally in normalization we force the data to be normal and thus symmetric, however, we could choose to scale the data in a way that leaves any asymmetry it had. By doing this we could potentially reduce our overall loss due to the clipping and rounding errors, but it's not a guarantee.

As a last resort, if none of these other methods failed to reduce accuracy loss of the model, we can use Quantization Aware Training (QAT). QAT is a simple process where we add a fake quantization step during training of the model. By fake, I mean, we clip and round the data while leaving it in full precision. This allows the model to adjust for the error and bias introduced by quantization while it's training. QAT is known to produce higher accuracy compared to other methods but at a much higher cost in time to train.

### **Quantization Methods**

**Uniform vs Non-uniform:** whether or not we use an 8-bit standard that is uniform in the range it represents or non-uniform to be more precise in the -1 to 1 range.

**Static vs Dynamic:** Choosing to adjust the range or scale before clipping in an attempt to reduce clipping and rounding errors and reduce data loss.

**Symmetric vs Asymmetric:** Normalizing the data to be normal and force

symmetry, or choosing to keep any asymmetry and skew.

During or After Training: Quantization after training is really easy to do, and while doing it during training is more work it leads to reduced bias and better results.

Quantizing is a very powerful tool. Not only does it reduce the size of the model, but it also reduces the computational overhead required to run the model thus reducing latency and cost of running the model. But the best fact about quantization is that it can be done after the fact, so you don't have to worry about whether or not your data scientists remembered to quantize the model during training using processes like QAT. This is why quantization has become so popular when working with LLMs and other large machine learning models. While reduced accuracy is always a concern with compression techniques, compared to other methods, quantization is a win-win-win.

## Pruning

Congratulations, you just trained a brand new LLM! With billions of parameters all of them must be useful right? Wrong! Unfortunately, as with most things in life, the model's parameters tend to follow the Pareto Principle. About 20% of the weights lead to 80% of the value. "If that's true," you may be asking yourself, "Why don't we just go and cut out all the extra fluff?" Great idea! Give yourself a pat on the back. Pruning is the process of weeding out and removing any parts of the model that we deem unworthy.

There are essentially two different pruning methods: **structured** and **unstructured**. Structured pruning is the process of finding structural components of a model that aren't contributing to the model's performance and then removing them. Whether they be filters, channels, or layers in the neural network. The advantages to this method is that your model will be a little smaller but keep the same basic structure, which means we don't have to worry about losing hardware efficiencies, we are also guaranteed a latency improvement as there will be less computations involved.

Unstructured pruning on the other hand, is shifting through the parameters

and zeroing out the less important ones that don't contribute much to the model's performance. Unlike structured pruning, we don't actually remove any parameters, just set them to zero. From this, we can imagine that a good place to start would be any weights or activations that are already close to 0. Of course, while this effectively reduces the size of a model this also means we don't cut out any computations, so it's common to only see minimal latency improvement—if at all. But a smaller model still means faster load times and less GPUs to run. It also gives us very fine-grained control over the process, allowing us to shrink a model further than we could with structured pruning with less impact to performance too.

Like quantization, pruning can be done after a model is trained. However, unlike quantization, it's common practice to see additional fine-tuning needing to be done to prevent too much loss of performance. It's becoming more common to just include pruning steps during the model training to avoid the need to fine-tune later on. Since a more sparse model will have fewer parameters to tune, adding these pruning steps may help a model converge faster as well.[\[5\]](#)

You'll be surprised at how much you can shrink a model with pruning while minimally impacting performance. How much? In the SparseGPT[\[6\]](#) paper, a method was developed to try to automatically one shot the pruning process without the need for finetuning after. They found they could decrease a GPT-3 model by 50-60% without issue! Depending on the model and task they even saw slight improvements in a few of them. Looking forward to seeing where Pruning takes us in the future.

## **Knowledge Distillation**

Knowledge distillation is probably the coolest method of compression in my mind. It's a simple idea too, we'll take the large LLM, and have it train a smaller language model to copy it. What's nice about this method is that the larger LLM provides essentially an infinite dataset for the smaller model to train on, which can make the training quite effective. Because of the simple fact that the larger the dataset the better the performance, we've often seen smaller models reach almost the same level as their teacher counterparts in accuracy.[\[7\]](#)

A smaller model trained this way is guaranteed to both be smaller and improve latency. The downside is that it's going to require us to train a completely new model. Which is a pretty significant upfront cost to pay. Any future improvements to the teacher model will require being passed down to the student model, which can lead to complex training cycles and version structure. It's definitely a lot more work compared to some of the other compression methods.

The hardest part about knowledge distillation though is that we don't really have good recipes for them yet. Tough questions like, "How small can the student model be?" will have to be solved through trial and error. There's still a lot to learn and research to be done here.

However, there has been some exciting work in this field via Stanford's Alpaca[\[8\]](#). Instead of training a student model from scratch, they instead chose to finetune the open source LLaMA 7B parameter model using OpenAI's GPT3.5's 175 billion parameter model as a teacher via knowledge distillation. A simple idea but it paid off big, as they were able to get great results from their evaluation. The biggest surprise was the cost, as they only spent \$500 on API costs to get the training data from the teacher model, and \$100 worth of GPU training time to finetune the student model. Granted, if you did this for a commercial application you'd be violating OpenAI's terms of service, so best to stick to using your own or open source models as the teacher.

## **Low-rank Approximation**

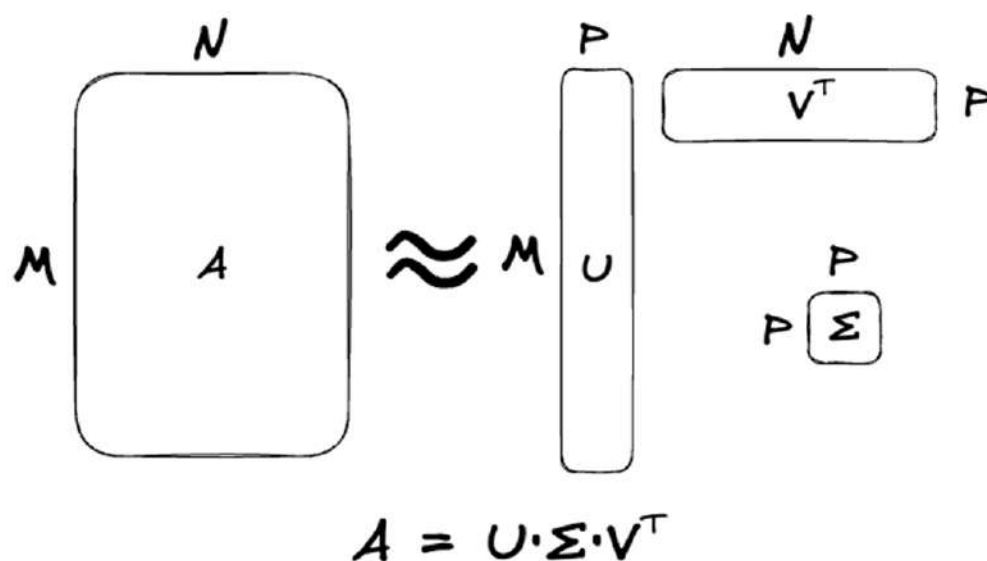
Low-rank approximation, also known as low-rank factorization, low-rank decomposition, matrix factorization (too many names! I blame the mathematicians), uses linear algebra math tricks to simplify large matrices or tensors finding a lower-dimensional representation. There are several techniques to do this. Singular Value Decomposition (SVD), Tucker Decomposition(TD), and Canonical Polyadic Decomposition (CPD) are the most common ones you run into.

In Figure 3.2 we show the general idea behind the SVD method. Essentially we are going to take a very large matrix,  $A$ , and break it up into 3 smaller



matrices,  $U$ ,  $\Sigma$ , and  $V$ . While  $U$  and  $V$  are there to ensure we keep the same dimensions and relative strengths of the original matrix,  $\Sigma$  allows us to apply a direction and bias. The smaller  $\Sigma$  is, the more we end up compressing and reducing the total number of parameters, but the less accurate the approximation becomes.

**Figure 3.2 Example of SVD a Low-rank Approximation.**  $A$  is a large matrix with dimensions  $N$  and  $M$ . We can approximate it with three smaller matrices,  $U$  with dimensions  $M$  and  $P$ ,  $\Sigma$  a square matrix with dimension  $P$ , and  $V$  with dimensions  $N$  and  $P$  (here we show the transpose). Usually both  $P \ll M$  and  $P \ll N$  are true.



To help solidify this concept, it may help to see a concrete example. In Listing 3.2 we show a simple example of SVD at work compressing a 4x4 matrix. For this we only need the basic libraries SciPy and NumPy which are imported on lines 1 and 2. Line 3 we define the matrix, and then line 9 we apply SVD to it.

**Listing 3.2 Example SVD Low-rank Approximation**

```
import scipy
import numpy as np
matrix = np.array([
    [ 1.,  2.,  3.,  4.],
    [ 5.,  6.,  7.,  8.],
    [ 9., 10., 11., 12.],
    [13., 14., 15., 16.]
])
```

```

])
u, s, vt = scipy.sparse.linalg.svds(matrix, k=1)
print(u,s,vt)
# [[-0.13472211]
# [-0.34075767]
# [-0.5467932 ]
# [-0.7528288 ]], [38.62266], [[-0.4284123 -0.47437257 -0.5203326

```

Taking a moment to inspect U, Sigma, and the transpose of V, we see a 4x1 matrix, a 1x1 matrix, and a 1x4 matrix respectively. All in all we now only need 9 parameters vs the original 16, shrinking the memory footprint almost in half.

Lastly, we multiply the matrices back together to get an approximation of the original matrix. In this case, the approximation isn't all that great, but we can still see the general order and magnitudes match the original matrix.

```

svd_matrix = u*s*vt
print(svd_matrix)
# array([[ 2.2291691,  2.4683154,  2.7074606,  2.9466066],
#        [ 5.6383204,  6.243202 ,  6.848081 ,  7.4529614],
#        [ 9.047472 , 10.018089 , 10.988702 , 11.959317 ],
#        [12.456624 , 13.792976 , 15.129323 , 16.465673 ]], dtype=f

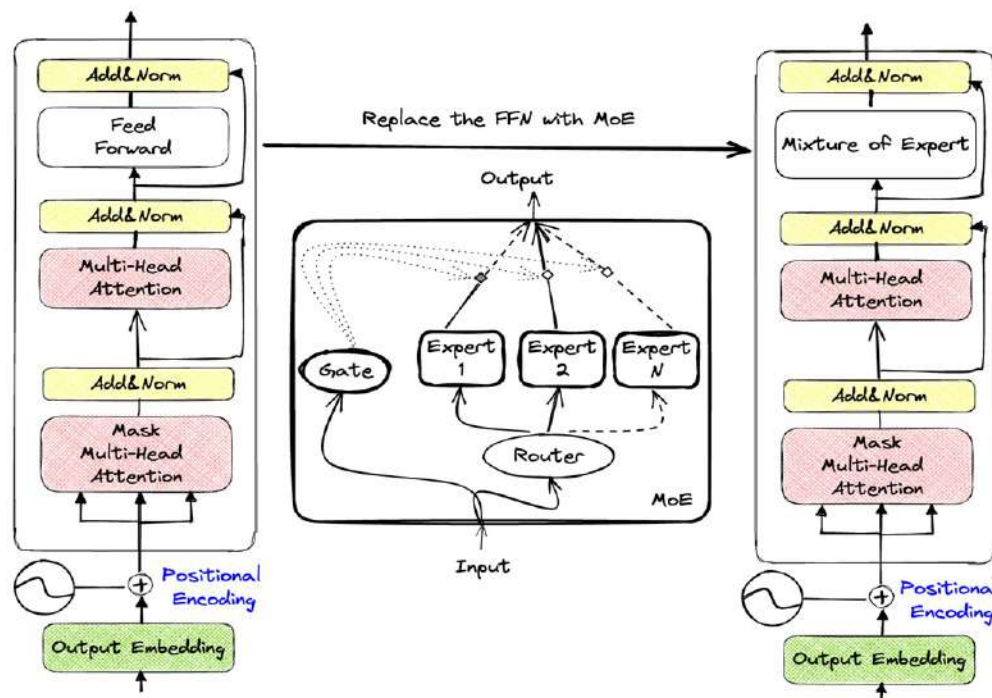
```

Unfortunately, I'm not aware of anyone actually using this to simply compress models in production most likely due to the poor accuracy of the approximation. What they are using it for, and this is important, is adaptation and finetuning; which is where LoRA[\[9\]](#) comes in, Low Rank Adaptation. Adaptation is simply the process of fine-tuning a generic or base model to do a specific task. LoRA applies SVD low rank approximation to the attention weights, or rather, to injected update matrices that run parallel to the attention weights, allowing us to fine-tune a much smaller model. LoRA has become very popular because it makes it a breeze to take an LLM, shrink the trainable layers to a tiny fraction of the original model and then allow anyone to train it on commodity hardware. You can get started with LoRA by using the PEFT[\[10\]](#) library from HuggingFace, where they have several LoRA tutorials you can check out.

## Mixture of Experts

Mixture of experts (MoE) is a technique where we replace the feed forward layers in a transformer with MoE layers instead. Feed forward layers are notorious for being parameter dense and computationally intensive so replacing them with something better can often have a large impact. MoE's are a group of sparsely activated models. It differs from ensemble techniques in that typically only one or a few expert models will be run, rather than combining results from all models. The sparsity is often induced by a Gate mechanism that learns which experts to use, and/or a Router mechanism that determines which experts should even be consulted. In Figure 3.3 we demonstrate the MoE architecture with potentially N experts, as well as show where it goes inside a decoder stack.

**Figure 3.3 Example Mixture of Experts model with both a Gate and Router to control flow. The MoE model is used to replace the FFN layers in a transformer, here we show it replacing the FFN in a Decoder.**



Depending on how many experts you have, the MoE layer could potentially have more parameters than the FFN leading to a larger model, but in practice this isn't the case since engineers and researchers are aiming to create a smaller model. What we are guaranteed to see though is a faster computation path and improved inference times. However, what really makes MoE stand out is when it's combined with quantization. One study[\[11\]](#) between

Microsoft and NVIDIA showed they were able to reach 2-bit quantization with only minimal impact to accuracy using MoE!

Of course, since this is a pretty big change to the model's structure it will require finetuning afterwards. You should also be aware that MoE layers often reduce a model's generalizability so it's best when used on models designed for a specific task. There are several libraries to implement MoE layers, but I'd recommend checking out DeepSpeed[\[12\]](#).

### **3.3.2 Distributed Computing**

Distributed computing is a technique used in deep learning to parallelize and speed-up large, complex neural networks by dividing the workload across multiple devices or nodes in a cluster. This approach significantly reduces training and inference times by enabling concurrent computation, data parallelism, and model parallelism. With the ever-growing size of datasets and complexity of models, distributed computing has become crucial for deep learning workflows, ensuring efficient resource utilization and enabling researchers to effectively iterate on their models. Distributed computing is one of the core practices that separate deep learning from machine learning, and with LLMs we have to pull out every trick in the book. Let's look at different parallel processing practices to take full advantage of distributed computing.

#### **Data Parallelism**

Data parallelism is what most people think about when they think about running processes in parallel, it's also the easiest to do. The practice involves splitting up the data and running them through multiple copies of the model or pipeline. For most frameworks this is easy to set up, for example, in PyTorch you can use the `DistributedDataParallel` method. There's just one catch for most of these set-ups and that is your model has to be able to fit onto one GPU. This is where a tool like Ray.io comes in.

Ray is an open-source project designed for distributed computing, specifically aimed at parallel and cluster computing. It's a flexible and user-friendly tool which simplifies distributed programming and helps developers

execute concurrent tasks in parallel with ease. Ray is primarily built for machine learning and other high-performance applications but can be utilized in other applications. In Listing 3.3 we give a simple example of using Ray to distribute a task. The beauty of Ray is the simplicity—all we need to do to make our code run in parallel is add a decorator. Sure beats the complexity of multithreading or asynchronization set-ups.

**Listing 3.3 Example Ray Parallelization Task**

```
import ray
import time

ray.init() # Start Ray

# Define a regular Python function
def slow_function(x):
    time.sleep(1)
    return x

# Turn the function into a Ray task
@ray.remote
def slow_function_ray(x):
    time.sleep(1)
    return x

# Execute the slow function without Ray (takes 10 seconds)
results = [slow_function(i) for i in range(1, 11)]

# Execute the slow function with Ray (takes 1 second)
results_future = [slow_function_ray.remote(i) for i in range(1, 11)]
results_ray = ray.get(results_future)

print("Results without Ray: ", results)
print("Results with Ray: ", results_ray)

ray.shutdown()
```

Ray uses the concept of tasks and actors to manage distributed computing. Tasks are functions, whereas actors are stateful objects that can be invoked and run concurrently. When you execute tasks using Ray, it handles distributing tasks across the available resources (e.g., multi-core CPUs or multiple nodes in a cluster). For LLMs, we would need to set up a Ray cluster[\[13\]](#) in a cloud environment as this would allow each pipeline to run

on a node with as many GPUs as needed, greatly simplifying the infrastructure set up to run LLMs in parallel.

There are multiple alternatives out there, but Ray has been gaining a lot of traction and becoming more popular as more and more machine learning workflows require distributed training. My team has had great success with it. By utilizing Ray, developers can ensure better performance and more efficient utilization of resources in distributed workflows.

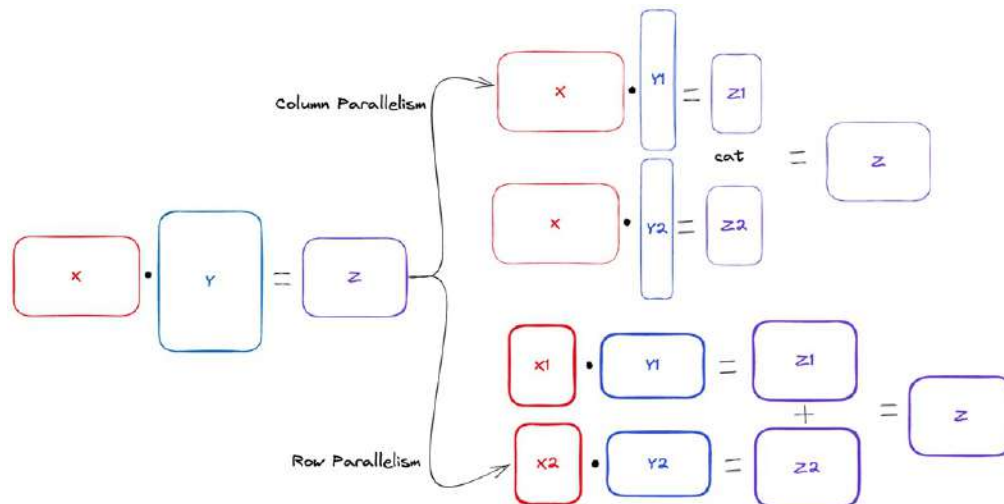
## **Tensor Parallelism**

Tensor parallelism is taking advantage of matrix multiplication properties to split up the activations across multiple processors, running the data through, and then combining them on the other side of the processors. Figure 3.4 demonstrates how this process works for a matrix, which can be parallelized in two separate ways that give us the same result. Imagine that  $Y$  is a really big matrix that can't fit on a single processor, or more likely, a bottleneck in our data flow that takes too much time to run all the calculations. In either case, we could split  $Y$  up, either by columns or by rows, run the calculations, and then combine the results after. In this example we are dealing with matrices but in reality we are often dealing with tensors that have more than two dimensions, but the same mathematical principles that make this work still apply.

Choosing which dimension to parallelize is a bit of an art, but a few things to remember to help make this decision easier. First, how many columns or rows do you have? In general, you want to pick a dimension that has more than the number of processors you have, else you will end up stopping short. Generally this isn't a problem but with tools like Ray discussed in the last section, parallelizing in a cluster and spinning up loads of processes is a breeze. Second, different dimensions have different multiplicity costs. For example, column parallelism requires us to send the entire dataset to each process, but with the benefit of concatenating them together at the end which is fast and easy. Row parallelism however, allows us to break up the dataset into chunks, but requires us to add the results, a more expensive operation than concatenating. You can see that one operation is more I/O bound, while the other is more computation bound. Ultimately, the best dimension will be

dataset dependent, as well as hardware limited. It will require experimentation to fully optimize this, but a good default is to just choose the largest dimension.

**Figure 3.4 Tensor Parallelism example showing that you can break up tensors by different dimensions and get the same end result. Here we compare column and row parallelism of a Matrix.**



Tensor parallelism allows us to split up the heavy computation layers like MLP and Attention layers onto different devices, but doesn't help us with Normalization or Dropout layers that don't utilize tensors. To get better overall performance of our pipeline we can add sequence parallelism which targets these blocks[\[14\]](#). Sequence parallelism is a process that partitions activations along the sequence dimension, preventing redundant storage, and can be mixed with tensor parallelism to achieve significant memory savings with minimal additional computational overhead. In combination, they reduce the memory needed to store activations in Transformer models. In fact, they nearly eliminate activation recomputation and save activation memory up to 5x.

**Figure 3.5 Combining tensor parallelism that focuses on computational heavy layers with sequence parallelism to reduce memory overhead to create a fully parallel process for the entire transformer.**

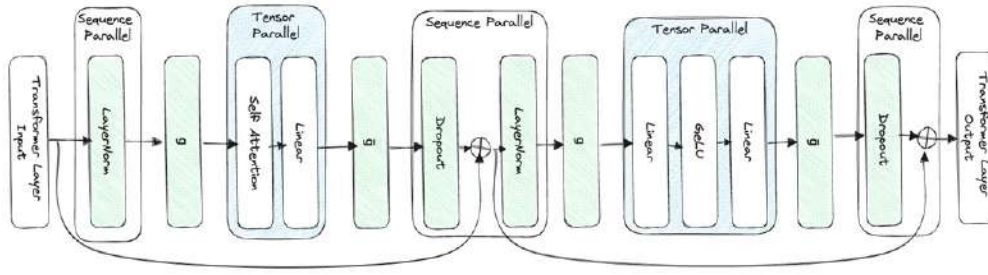


Figure 3.5 shows how combining both tensor parallelism, that allows us to distribute the computationally heavy layers, and sequence parallelism that does the same for the memory limiting layers, allows us to fully parallelize the entire transformer model. Together, they allow for extremely efficient use of resources.

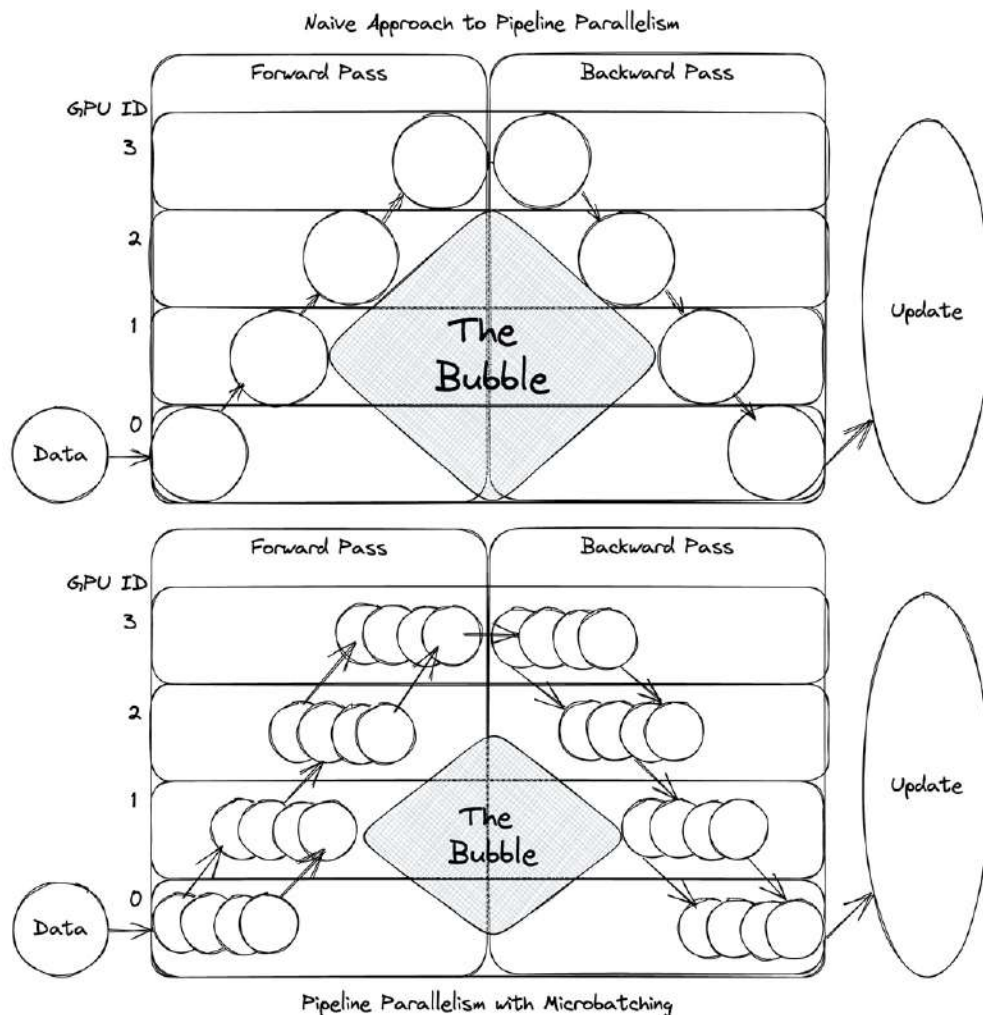
## Pipeline Parallelism

So far we can now run lots of data, and speed up any bottlenecks, but none of that matters because our model is too big; we can't fit it into a single GPU's memory to even get it to run. That's where pipeline parallelism comes in and is the process of splitting up a model vertically and putting each part onto a different GPU. This creates a pipeline as input data will go to the first GPU, process, then transfer to the next GPU, and so on until it's run through the entire model. While other parallelism techniques improve our processing power and speed up inference, pipeline parallelism is required to just get it to run and it comes with some major downsides, mainly device utilization.

To understand where this downside comes from and how to mitigate it, let's first consider the naive approach to this, where we simply run all the data at once through the model. What we find is that this leaves a giant "bubble" of underutilization. Since the model is broken up, we have to process everything sequentially through the devices. This means that while one GPU is processing, the others are sitting idle. In Figure 3.6 we can see this naive approach and a large bubble of inactivity as the GPUs sit idle. We also see a better way to take advantage of each device. We do this by sending the data in small batches. A smaller batch allows the first GPU to pass on what it was working on quicker and move on to another batch. This allows the next device to get started earlier and reduces the size of the bubble.



**Figure 3.6 The Bubble Problem.** When data runs through a broken up model, the GPUs holding the model weights are underutilized while they wait for their counterparts to process the data. A simple way to reduce this bubble is to use microbatching.



We can actually calculate the size of the bubble quite easily with the following formula:

$$\text{Idle Percentage} = 1 - m / (m + n - 1)$$

Where  $m$  is the number of microbatches, and  $n$  is the depth of the pipeline or number of GPUs. So for our naive example case of 4 GPUs and 1 large batch we see the devices sitting idle 75% of the time! GPUs are quite expensive to allow to be sitting idle three quarters of the time. Let's see what that looks like using the microbatch strategy. With a microbatch of 4, it cuts this almost in half, down to just 43% of the time. What we can glean from this formula is

that the more GPUs we have, the higher the idle times, but the more microbatches the better the utilization.

Unfortunately, we can often neither reduce the number of GPUs nor can we just make the microbatches as large as we want. There are limits. For GPU's we just have to use as many as it takes to fit the model into memory, but try to use a few larger GPUs as it will lead to a more optimal utilization compared to using many smaller GPUs. Reducing the bubble in pipeline parallelism is another reason why compression is so important. For microbatching, the first limit is obvious once told, since the microbatch is a fraction of our batch size, we are limited by how big that is. The second is that each microbatch increases the memory demand for cached activations in a linear relationship. One way to counter this higher memory demand is a method called PipeDream[\[15\]](#). There are different configurations and approaches, but the basic idea is the same. In this method we start working on the backward pass as soon as we've finished the forward pass of any of the microbatches. This allows us to fully complete a training cycle and release the cache for that microbatch.

### **3D Parallelism**

For LLMs, we are going to want to take advantage of all three parallelism practices as they can all be run together. This is known as 3D Parallelism combining Data, Tensor, and Pipeline Parallelism (DP + TP + PP) together. Since each technique, and thus dimension, will require at least 2 GPUs, in order to run 3D Parallelism, we'll need at least 8 GPUs to even get started. How we configure these GPUs will be important to get the most efficiency out of this process, namely, because TP has the largest communication overhead we want to ensure these GPUs are close together, preferably on the same node and machine. PP has the least communication volume of the three, so breaking up the model across nodes is the least expensive here.

By running the three together, we see some interesting interactions and synergies between them. Since TP splits the model to work well within a device's memory, we see that PP can perform well even with small batch sizes due to the reduced effective batch size enabled by TP. This combination also improves the communication between DP nodes at different pipeline

stages, allowing DP to work effectively too. The communication bandwidth between nodes is proportional to the number of pipeline stages, because of this DP is able to scale well even with smaller batch sizes. Overall, we see running in combination that

Now that we know some tricks of the trade, it's just as important to have the right tools to do the job.

## **3.4 Large Language Models Operations Infrastructure**

We are finally going to start talking about the infrastructure needed to make this all work. This likely comes as a surprise as I know several readers would have expected this section at the beginning of chapter 1. Why wait till the end of chapter 3? In the many times I've interviewed Machine Learning Engineers I often asked this open-ended question, "What can you tell me about MLOps?" An easy softball question to get the conversation going. Most junior candidates would immediately start jumping into the tooling and infrastructure. It makes sense, there are so many different tools available. Not to mention, whenever you see posts or blogs describing MLOps there's a pretty little diagram showing the infrastructure. While all of that is important it's useful to recognize what a more senior candidate jumps into, the machine learning lifecycle.

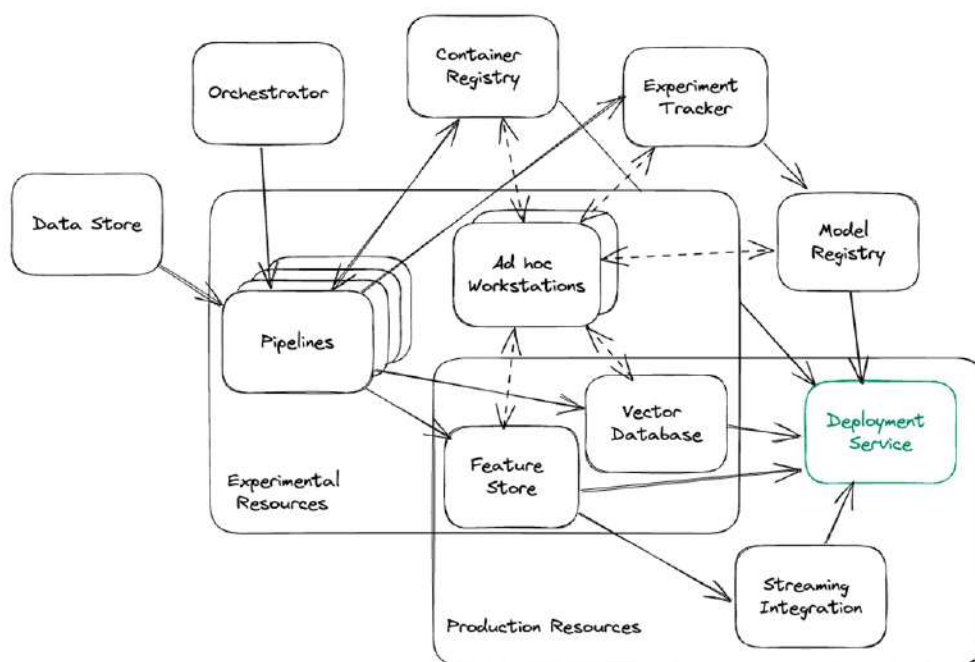
For many the nuance is lost, but the infrastructure is the how, the lifecycle is the why. Most companies can get by with just bare-bones infrastructure. I've seen my share of scrappy systems that exist entirely on one Data Scientist's laptop, and they work surprisingly well! Especially in the era of scikit learn everything.

Unfortunately, a rickshaw machine learning platform doesn't cut it in the world of LLMs. Since we still live in a world where the standard storage capacity of a MacBook Pro laptop is 256GB, just storing the model locally can already be a problem. Companies that invest in a more sturdy infrastructure are better prepared for the world of LLMs.

In Figure 3.7 we see an example MLOps Infrastructure designed with LLMs

in mind. While most infrastructure diagrams I've seen in my time have always simplified the structure to make everything look clean, the raw truth is that there's a bit more complexity to the entire system. Of course a lot of this complexity would disappear if we could just get Data Scientists to work inside scripts instead of ad hoc workstations—usually with a jupyter notebook interface.

**Figure 3.7 a high level view of an MLOps infrastructure with LLMs in mind. This attempts to cover the full picture, and the complexity of the many tools involved to make ML models work in production.**



Taking a closer look at Figure 3.7 you can see several tools on the outskirts that squarely land in DataOps, or even just DevOps. Data Stores, Orchestrators, Pipelines, Streaming Integrations, and Container Registries. These are tools you are likely already using for just about any data intensive application and aren't necessarily MLOps focused. Towards the center we have more traditional MLOps tools, Experiment Trackers, Model Registry, Feature Store, and Ad hoc Data Science Workstations. For LLMs we really only introduce one new tool to the stack: a Vector Database. What's not pictured because it intertwines with every piece is the Monitoring System. This all culminates to what we are working towards in this book, a Deployment Service, where we can confidently deploy and run LLMs in

Production.

#### Infrastructure by discipline

**DevOps:** In charge of procuring the environmental resources—experimental (dev, staging) and production—this includes hardware, clusters, and networking to make it all work. Also in charge of basic infrastructure systems like Github/Gitlab, artifact registries, container registries, application or transactional databases like Postgres or MySQL, caching systems, and CI/CD pipelines. This is by no means a comprehensive list.

**DataOps:** In charge of data, in motion and at rest. Includes centralized or decentralized data stores like Data Warehouses, Data Lakes, and Data Meshes. As well as data pipelines either in batch systems or in streaming systems with tools like Kafka and Flink. Also includes orchestrators like Airflow, Prefect or Mage. DataOps is built on top of DevOps. For example, I've seen many CI/CD pipelines being used for data pipeline work until eventually being graduated to systems like Apache Spark or DBT.

**MLOps:** In charge of machine learning lifecycle from creation of models to deprecation. This includes data science workstations like Jupyterhub, experiment trackers, and a model registry. It includes specialty databases like Feature Stores and Vector Databases. As well as a deployment service to tie everything together and actually serve results. It is built on top of both DataOps and DevOps.

Let's go through each piece of the infrastructure puzzle and discuss features you should be considering when thinking about LLMs in particular. While we will be discussing tooling that is specialized for each piece, I'll just make note that there are also MLOps as a service platforms like Dataiku, Amazon's Sagemaker and Google's VertexAI. These platforms attempt to give you the whole puzzle, how well they do that is another question, but are often a great shortcut and you should be aware of them. Well, I think that's enough dilly-dallying, let's dive in already!

### 3.4.1 Data Infrastructure

While not the focus of this book it's important to note that MLOps is built on top of a Data Operations infrastructure—which itself is built on top of DevOps. Key features of the DataOps ecosystem include a data store, an orchestrator, and pipelines. Additional features usually required include a container registry and a streaming integration service.

Data stores are the foundation of DataOps and come in many forms these days, from a simple database to large data warehouses, and from even larger data lakes to an intricate data mesh. This is where your data is stored and a lot of work goes into managing, governing, and securing the data store. The orchestrator is the cornerstone of DataOps as it's a tool that manages and automates both simple and complex, multistep workflows and tasks. Ensuring they run across multiple resources and services in a system. The most commonly talked about being Airflow, Prefect, and Mage. Lastly, pipelines are the pillars. They hold everything else up, and are where we actually run our jobs. Initially built to simply move, clean, and define data, these same systems are used to run machine learning training jobs on a schedule, do batch inference, and loads of other work needed to ensure MLOps runs smoothly.

A container registry is a keystone of DevOps and subsequently DataOps and MLOps as well. Being able to run all our pipelines and services in containers is necessary to ensure consistency. Streaming services are actually a much bigger beast than what I may let on in this chapter, and if you know you know. Thankfully for most text related tasks real time processing isn't a major concern. Even for tasks like real-time captioning or translation, we can often get by with some sort of pseudo real-time processing strategy that doesn't degrade the user experience depending on the task.

### **3.4.2 Experiment Trackers**

Experiment trackers are central to MLOps. Experiment trackers do the fundamental job of keeping track and recording tests and results. As the famous Adam Savage quote from Mythbusters, “Remember kids, the only difference between screwing around and science is writing it down.” Without it, your organization is likely missing the “science” in data science which is honestly quite embarrassing.

Even if your data scientists are keen to manually track and record results in notebooks, it might as well be thrown in the garbage if it's not easy for others to view and search for. This is really the purpose of experiment trackers, to ensure knowledge is easily shared and made available. Eventually a model will make it to production and that model is going to have issues. Sure, you can always just train a new model, but unless the team is able to go back and investigate what went wrong the first time you are likely to repeat the same mistakes over and over.

There are many experiment trackers out there, the most popular by far is MLFlow which is open source. It was started by the team at Databricks which also offers an easy hosting solution. Some paid alternatives worth checking out include CometML and Weights and Biases.

Experiment trackers nowadays come with so many bells and whistles. Most open source and paid solutions will certainly have what you need when looking to scale up your needs for LLMOps. However, ensuring you take advantage of these tools correctly might require a few small tweaks. For example, the default assumption is usually that you are training a model from scratch, but often when working with LLMs you will be finetuning models instead. In this case, it's important to note the checkpoint of the model you started from. If possible, even linking back to the original training experiment. This will allow future scientists to dig deeper into their test results, find original training data, and discover paths forward to eliminate bias.

Another feature to look out for is evaluation metric tooling. We will be going more in-depth in Chapter 4, but evaluation metrics are difficult for language models. There are often multiple metrics you care about and none of them are simple like complexity ratings or similarity scores. While experiment tracker vendors try to be agnostic and unopinionated about evaluation metrics they should at least make it easy to compare models and their metrics to make it easy to decide which one is better. Since LLMs have become so popular some have made it easy to evaluate on the more common metrics like ROUGE for text summarization.

You will also find many experiment tracking vendors have started to add tools specifically for LLMs. Some features you might consider looking for

include direct HuggingFace support, LangChain support, prompt engineering toolkits, finetuning frameworks, and foundation model shops. The space is developing quickly, and no one tool has all the same features right now, but I'm sure these feature sets will likely converge.

### **3.4.3 Model Registry**

The model registry is probably the simplest tool of an MLOps infrastructure. The main objective is one that's easy to solve, we just need a place to store the models. I've seen many successful teams get by simply by putting their models in an object store or shared file system and calling it good. That said, there's a couple bells and whistles you should look for when choosing one.

The first is whether or not the model registry tracks metadata about the model. Most of what you care about is going to be in the experiment tracker, so you can usually get away with simply ensuring you can link the two. In fact, most model registries are built into experiment tracking systems because of this. However, an issue I've seen time and time again with these systems happens when the company decides to use an open source model or even buy one. Is it easy to upload a model and tag it with relevant information? The answer is usually no.

Next, you are going to want to make sure you can version your models. At some point, a model will reach a point where it's no longer useful and will need to be replaced. Versioning your models will simplify this process. It also makes running production experiments like A/B testing or shadow tests easier.

Lastly, if we are promoting and demoting models, we need to be concerned with access. Models tend to be valuable intellectual property for many companies, ensuring only the right users have access to the models is important. But it's also important to ensure that only the team that understands the models, what they do and why they were trained, are in charge of promoting and demoting the models. The last thing we want is to delete a model in production or worse.

For LLMs there are some important caveats you should be aware of, mainly,



when choosing a model registry, be aware of any limit sizes. I've seen several model registries restrict model sizes to 10GB or smaller. That's just not going to cut it. I could speculate on the many reasons for this but none of them are worthy of note. Speaking of limit sizes, if you are going to be running your model registry on an premise storage system like Ceph, make sure it has lots of space. You can buy multiple terabytes of storage for a couple hundred dollars for your on prem servers, but even a couple terabytes fills up quickly when your LLM is over 300GB. Don't forget, you are likely to be keeping multiple checkpoints and versions during training and finetuning; as well as duplicates for reliability purposes. Storage is still one of the cheapest aspects of running LLMs though, no reason to skimp here and cause headaches down the road.

This does bring me to a good point: there's a lot of optimization that could still be made, allowing for better space saving approaches to storing LLMs and their derivatives. Especially since most of these models will be very similar overall. I imagine we'll likely see storage solutions to solve just this problem in the future.

### **3.4.4 Feature Store**

Feature stores solve many important problems and answer questions like: Who owns this feature? How was it defined? Who has access to it? Which models are using it? How do I serve this feature in production? Essentially, they solve the “single source of truth” problem. By creating a centralized store, it allows teams to shop for the highest quality, most well maintained, thoroughly managed data. Feature stores solve the collaboration, documentation, and versioning of data.

If you've ever thought, “A feature store is just a database right?” you are probably thinking about the wrong type of store—we are referencing a place to shop not a place of storage. Don't worry, this confusion is normal as I've heard this sentiment a lot, and have had similar thoughts myself. The truth is, modern day feature stores are more virtual than a physical database, which means they are built on top of whatever data store you are already using. For example, Google's Vertex AI feature store is just BigQuery and I've seen a lot of confusion from data teams wondering, “Why don't I just query

BigQuery?” Loading the data into a feature store feels like an unnecessary extra step, but think about shopping at an IKEA store. No one goes directly to the warehouse where all the furniture is in boxes. That would be a frustrating shopping experience. The features store is the show rooms that allows others in your company to easily peruse, experience, and use the data.

Often times, I see people reach for a feature store to solve a technical problem like low latency access for online feature serving. A huge win for feature stores is solving the training-serving skew. Some features are just easier to do in SQL after the fact, like calculating the average number of requests for the last 30 seconds. This can lead to naive data pipelines built for training, but causing massive headaches when going to production because getting this type of feature in real time can be anything but easy. Feature stores abstractions help minimize this burden. Related to this is feature stores point-in-time retrievals which are table stakes when talking feature stores. Point-in-time retrievals ensure that given a specific time a query will always return the same result. This is important because features like averages over “the last 30 seconds” are constantly changing, so this allows us to version the data (without the extra burden of a bloated versioning system), as well as ensure our models will give accurate and predictable responses.

As far as options, Feast is a popular open source feature store. FeatureForm and Hopsworks are also open source. All three of which offer paid hosting options. For LLMs I’ve heard the sentiment that feature stores aren’t as critical as other parts of the MLOps infrastructure. After all, the model is so large it should incorporate all needed features inside it, so you don’t need to query for additional context, just give the model the user’s query and let the model do its thing. However, this approach is still a bit naive and we haven’t quite gotten to a point where LLMs are completely self-sufficient. To avoid hallucinations and improve factual correctness, it is often best to give the model some context. We do this by feeding it embeddings of our documents we want it to know very well, and a feature store is a great place to put these embeddings.

### **3.4.5 Vector Databases**

If you are familiar with the general MLOps infrastructure, most of this

section has been review for you. We've only had to make slight adjustments highlighting important scaling concerns to make a system work for LLMs. Vector Databases however are new to the scene and have been developed to be a tailored solution for working with LLMs and language models in general, but you can also use them with other datasets like images or tabular data which are easy enough to transform into a vector. Vector databases are specialized databases to store vectors along with some metadata around the vector, which makes them great for storing embeddings. Now, while that last sentence is true, it is a bit misleading, because the power of vector databases isn't in their storage, but in the way that they search through the data.

Traditional databases, using b-tree indexing to find IDs or text based search using reverse indexes, all have the same common flaw, you have to know what you are looking for. If you don't have the ID or you don't know the keywords it's impossible to find the right row or document. Vector databases however, take advantage of the vector space meaning you don't need to know exactly what you are looking for, you just need to know something similar which you can then use to find the nearest neighbors using similarity searches based on Euclidean distance, Cosine similarity, Dot product similarity, or what have you. Using a vector database makes solving the reverse image search problem a breeze, as an example.

At this point, I'm sure some readers may be confused. First I told you to put your embeddings into a feature store, and now I'm telling you to put them into a Vector DB, which one is it? Well, that's the beauty of it, you can do both at the same time. If it didn't make sense before I hope it makes sense now. Feature stores are not a database they are just an abstraction, you can use a feature store built on top of a Vector DB and it will solve many of your problems. Vector DBs can be difficult to maintain when you have multiple data sources, are experimenting with different embedding models, or otherwise have frequent data updates. Managing this complexity can be a real pain, but a feature store can handily solve this problem. Using them in combination will ensure a more accurate and up-to-date search index.

Vector databases have only been around for a couple of years at the time of writing, and their popularity is still relatively new as it has grown hand in hand with LLMs. It's easy to understand why since they provide a fast and

efficient way to retrieve vector data making it easy to provide LLMs with needed context to improve their accuracy.

That said it's a relatively new field and there are lots of competitors in this space right now, and it's a bit too early to know who the winners and losers are. Not wanting to date this book too much, let me at least suggest two options to start Pinecone and Milvus. Pinecone is one of the first vector databases as a product and has a thriving community with lots of documentation. It's packed with features and has proven itself to scale. Pinecone is a fully managed infrastructure offering that has a free tier for beginners to learn. If you are a fan of open source, however, then you'll want to check out Milvus. Milvus is feature-rich and has a great community. Zilliz the company behind Milvus offers a fully managed offering, but it's also available to deploy in your own clusters and if you already have a bit of infrastructure experience it's relatively easy and straightforward to do.

There are lots of alternatives out there right now, and it's likely worth a bit of investigation before picking one. Probably the two things you'll care most about are price and scalability, the two often going hand in hand. After that, it's valuable to pay attention to search features, like support for different similarity measures like cosine similarities, dot product, or Euclidean distance. As well as indexing features like HNSW (Hierarchical Navigable Small World) or LSH (Locality-Sensitive Hashing). Being able to customize your search parameters and index settings are important for any database as they allow you to customize the workload for your dataset and workflow allowing you to optimize query latency and search result accuracy.

It's also important to note that with vector databases rise in popularity we are quickly seeing many database incumbents like Redis and Elastic offering vector search capabilities. For now, most of these tend to just offer the most straightforward feature sets, but they are hard to ignore if you are already using these tool sets as they can provide quick wins to get started quickly.

Vector databases are powerful tools that can help you train or finetune LLMs, as well as improve the accuracy and results of your LLM queries.

### **3.4.6 Monitoring System**

A monitoring system is crucial to the success of any ML system, LLMs included. Unlike other software applications, ML models are known to fail silently—continue to operate, but start to give poor results. This is often due to data drift, a common example being a recommendation system that give worse results overtime because sellers start to game the system by giving fake reviews to get better recommendation results. A monitoring system allows us to catch poorly performing models, make adjustments or simply retrain them.

Despite their importance they are often the last piece of the puzzle added. This is often purposeful as putting resources into figuring out how to monitor models doesn't help if you don't have any models to monitor. However, don't make the mistake of putting it off too long. Many companies have been burned by a model that went rogue with no one knowing about it, often costing them dearly. It's also important to realize you don't have to wait to get a model into production to start monitoring your data. There are plenty of ways to introduce a monitoring system into the training and data pipelines to improve data governance and compliance. Regardless, you can usually tell the maturity of a data science organization by their monitoring system.

There are lots of great monitoring tooling out there, some great open source options include WhyLogs and EvidentlyAI. I'm also a fan of great expectations, but have found it rather slow outside of batch jobs. There are also many more paid options out there. Typically, for ML monitoring workloads you'll want to monitor everything you'd normally record in other software applications, this includes resource metrics like memory and CPU utilization, performance metrics like latency and queries per second, as well as operational metrics like status codes and error rates. In additional, you'll need ways to monitor data drift both going in and out of the model. You'll want to pay attention to things like missing values, uniqueness, and standard deviation shifts. In many instances, you'll want to be able to segment your data while monitoring, e.g. for A/B testing or to monitor by region. Some metrics that are useful to monitor in ML systems include model accuracy, precision, recall and F1 scores. These are difficult since you won't know the correct answer at inference time, so it's often useful to set up some sort of auditing system. Of course, auditing is going to be easier if your LLM is designed to be a Q&A bot than if your LLM is meant to help writers be more

creative.

This hints at a fact that for LLMs, there are often a whole set of new challenges for your monitoring systems even more than what we see with other ML systems. With LLMs we are dealing with text data which is hard to quantify as discussed earlier in this chapter. For instance, think about what features do you look at to monitor for data drift? Because language is known to drift a lot! One feature I might suggest is unique tokens. This will alert you when new slang words or terms are created, however, it still doesn't help when words switch meaning, for example, when "wicked" means "cool". I would also recommend monitoring the embeddings, however, you'll likely find this to either add a lot of noise and false alarms or at the very least be difficult to decipher and dig into when problems do occur. The systems I've seen work the best often involve a lot of handcrafted rules and features to monitor, but these can be error-prone and time-consuming to create.

Monitoring text based systems is far from a solved problem, mostly stemming from the difficulties of understanding text data to begin with. This does beg the question of what are the best methods to use language models to monitor themselves, since they are our current best solution to codifying language. Unfortunately, I'm not aware of anyone researching this, but imagine it's only a matter of time.

### **3.4.7 GPU Enabled Workstations**

GPU enabled workstations and remote workstations in general are often considered a nice to have or luxury by many teams, but when working with LLMs that mindset has to change. When troubleshooting an issue or just developing a model in general, a data scientist isn't going to be able to spin up the model in a notebook on their laptop anymore. The easiest way to solve this is to simply provide remote workstations with GPU resources. There are plenty of cloud solutions for this, but if your company is working mainly on prem, this may be a bit more difficult to provide, but necessary nonetheless.

LLMs are GPU memory intensive models. Because of this, there are some numbers every engineer should know when it comes to working in the field. The first, is how much different GPUs have. The NVIDIA Tesla T4 and

V100 are two most common GPUs you'll find in a datacenter, but they only have 16 GB of memory. They are workhorses though and cost-effective, so if we can compress our model to run on these all the better. After these, you'll find a range of GPU's like NVIDIA A10G, NVIDIA Quadro series, and NVIDIA RTX series that offer GPU memories in the ranges of 24, 32, and 48 GB. All of these are fine upgrades, you'll just have to figure out which ones are offered and available to you by your cloud provider. Which brings us to the NVIDIA A100, which is likely going to be your GPU of choice when working with LLMs. Thankfully they are relatively common, and offer two different models providing 40 or 80 GB. The big issue you'll have with these are that they are constantly in high demand by everyone right now. You should also be aware of the NVIDIA H100 which offers 80 GB like the A100. The H100 NVL is promised to support up to 188 GB and has been designed with LLMs in mind. Another new GPU you should be aware of is the NVIDIA L4 Tensor Core GPU which has 24 GB and is positioned to take over as a new workhorse along with the T4 and V100, at least as far as AI workloads are concerned.

LLMs come in all different sizes, and it's useful to have a horse sense for what these numbers mean. For example, the LLaMA model has 7B, 13B, 33B, and 65B parameter variants. If you aren't sure which GPU you need to run which model off the top of your head, here's a shortcut, just take the number of billions of parameters times it by two and that's how much GPU memory you need. The reason is, most models at inference are going to default to run at half precision, FP16 or BF16, which means for every parameter we need at least two bytes. Thus, 7 billion \* 2 bytes = 14 GB. You'll need a little extra as well for the embedding model which will be about another GB, and more for the actual tokens you are running through the model. One token is about 1 MB, so 512 tokens will require 512 MB. This isn't a big deal, until you consider running larger batch sizes to improve performance. For 16 batches of this size you'll need an extra 8 GB of space.

Of course, so far we've only been talking about inference, for training you'll need a lot more space. While training, you'll always want to do this in full precision, and you'll need extra room for the optimizer tensors and gradients. In general, to account for this you'll need about 16 bytes for every parameter. So to train a 7B parameter model you'll want 112 GB of memory.

### 3.4.8 Deployment Service

Everything we've been working towards is collected and finally put to good use here. In fact, if you took away every other service and were left with just a deployment service, you'd still have a working MLOps system. A deployment service provides an easy way to integrate with all the previous systems we talked about as well as configure and define the needed resources to get our model running in production. It will often provide boilerplate code to serve the model behind a REST and gRPC API or directly inside a batch or streaming pipeline.

Some tools to help create this service include NVIDIA Triton Inference Service, MLServer, Seldon and BentoML. These services provide a standard API interface, typically the KServe V2 Inference Protocol. This protocol provides a unified and extensible way to deploy, manage, and serve machine learning models across different platforms and frameworks. It defines a common interface to interact with models, including gRPC and HTTP/RESTful APIs. It standardizes concepts like input/output tensor data encoding, predict and explain methods, model health checks, and metadata retrieval. It also allows seamless integration with languages and frameworks including TensorFlow, PyTorch, ONNX, Scikit Learn, and XGBoost.

Of course, there are times when flexibility and customization provide enough value to step away from the automated path these other frameworks provide, in which case it's best to reach for a tool like FastAPI. Your deployment service should still provide as much automation and boilerplate code here to make the process as smooth as possible. It should be mentioned that most of the frameworks mentioned above do offer custom methods, but your mileage may vary.

Deploying a model is more than just building the interface. Your deployment service will also provide a bridge to close the gap between the MLOps infrastructure and general DevOps infrastructure. Connecting to whatever CI/CD tooling as well as build and shipping pipelines your company has set up so you can ensure appropriate tests and deployment strategies like health checks and rollbacks can easily be monitored and done. This is often very platform and thus company-specific. Thus, it'll also need to provide the



needed configurations to talk to Kubernetes, or whatever other container orchestrator you may be using, to acquire the needed resources like CPU, Memory, and Accelerators, Autoscalers, Proxies, etc. It also applies the needed environment variables and secret management tools to ensure everything runs.

All in all, this service ensures you can easily deploy a model into production. For LLMs, the main concern is often just ensuring the platform and clusters are set up with enough resources to actually provision what will ultimately be configured.

We've discussed a lot so far in this chapter, starting with what makes LLMs so much harder than traditional ML which is hard enough as it is. First, we learned that their size can't be underestimated, but then we also discovered there are many peculiarities about them, from token limits to hallucinations, not to mention they are expensive. Fortunately, despite being difficult, they aren't impossible. We discussed compression techniques and distributed computing which are crucial to master. We then explored the infrastructure needed to make LLMs work. While most of it was likely familiar, we came to realize that LLMs put a different level of pressure on each tool, and often we need to be ready for a larger scale than what we could get away with for deploying other ML models.

## 3.5 Summary

- LLMs are difficult to work with mostly because they are big. Which impacts a longer time to download, load into memory, and deploy forcing us to use expensive resources.
- LLMs are also hard to deal with because they deal with natural language and all its complexities including hallucinations, bias, ethics, and security.
- Regardless if you build or buy, LLMs are expensive and managing costs and risks associated with them will be crucial to the success of any project utilizing them.
- Compressing models to be as small as we can will make them easier to work with; quantization, pruning, and knowledge distillation are particularly useful for this.

- Quantization is popular because it is easy and can be done after training without any finetuning.
- Low Rank Approximation is an effective way at shrinking a model and has been used heavily for Adaptation thanks to LoRA.
- There are three core directions we use to parallelize LLM workflows: Data, Tensor, and Pipeline. DP helps us increase throughput, TP helps us increase speed, and PP makes it all possible to run in the first place.
- Combining the parallelism methods together we get 3D parallelism (Data+Tensor+Pipeline) where we find that the techniques synergize, covering each others weaknesses and help us get more utilization.
- The infrastructure for LLMOps is similar to MLOps, but don't let that fool you since there are many caveats where "good enough" no longer works.
- Many tools are offering new features specifically for LLM support.
- Vector Databases in particular are interesting as a new piece of the infrastructure puzzle needed for LLMs that allow quick search and retrievals of embeddings.

[1] A. Bulatov, Y. Kuratov, and M. S. Burtsev, "Scaling Transformer to 1M tokens and beyond with RMT," Apr. 2023, <https://arxiv.org/abs/2304.11062>.

[2] R. Daws, "Medical chatbot using OpenAI's GPT-3 told a fake patient to kill themselves," AI News, Oct. 28, 2020. <https://www.artificialintelligence-news.com/2020/10/28/medical-chatbot-openai-gpt3-patient-kill-themselves/>

[3] T. Kington, "ChatGPT bot tricked into giving bomb-making instructions, say developers," [www.thetimes.co.uk](http://www.thetimes.co.uk), Dec 17, 2022. <https://www.thetimes.co.uk/article/chatgpt-bot-tricked-into-giving-bomb-making-instructions-say-developers-rvktrxb5>

[4] K. Quach, "AI game bans players for NSFW stories it generated itself," [www.theregister.com](http://www.theregister.com), Oct 8, 2021. [https://www.theregister.com/2021/10/08/ai\\_game\\_abuse/](https://www.theregister.com/2021/10/08/ai_game_abuse/)

[5] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste, "Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks," Jan. 2021, <https://arxiv.org/abs/2102.00554>.

[6] E. Frantar and D. Alistarh, “SparseGPT: Massive Language Models Can Be Accurately Pruned in One-Shot,” Jan. 2023, <https://arxiv.org/abs/2301.00774>.

[7] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter,” Oct. 2019. <https://arxiv.org/abs/1910.01108>.

[8] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, “Alpaca: A Strong, Replicable Instruction-Following Model,” crfm.stanford.edu, 2023. <https://crfm.stanford.edu/2023/03/13/alpaca.html>

[9] E. J. Hu et al., “LoRA: Low-Rank Adaptation of Large Language Models,” Jun. 2021, <https://arxiv.org/abs/2106.09685>.

[10] For the extra curious, Parameter-Efficient Fine-Tuning (PEFT) is a class of methods aimed at fine-tuning models in a computational efficient way. The PEFT library seeks to put them all in one easy to access place and you can get started here: <https://huggingface.co/docs/peft>

[11] R. Henry and Y. J. Kim, “Accelerating Large Language Models via Low-Bit Quantization,” March 2023, <https://www.nvidia.com/en-us/on-demand/session/gtcspring23-s51226/>

[12] DeepSpeed is a library that optimizes many of the hard parts for large-scale deep learning models like LLMs and is particularly useful when training. Check out their MoE tutorial. [https://www.deepspeed.ai/tutorials/mixture-of-experts /](https://www.deepspeed.ai/tutorials/mixture-of-experts/)

[13] Learn more about Ray Clusters here: <https://docs.ray.io/en/releases-2.3.0/cluster/key-concepts.html#ray-cluster>

[14] V. Korthikanti et al., “Reducing Activation Recomputation in Large Transformer Models,” May 2022, <https://arxiv.org/abs/2205.05198>

[15] A. Harlap et al., “PipeDream: Fast and Efficient Pipeline Parallel DNN Training,” Jun. 08, 2018. <https://arxiv.org/abs/1806.03377>

# 4 Data Engineering for Large Language Models: Setting up for success

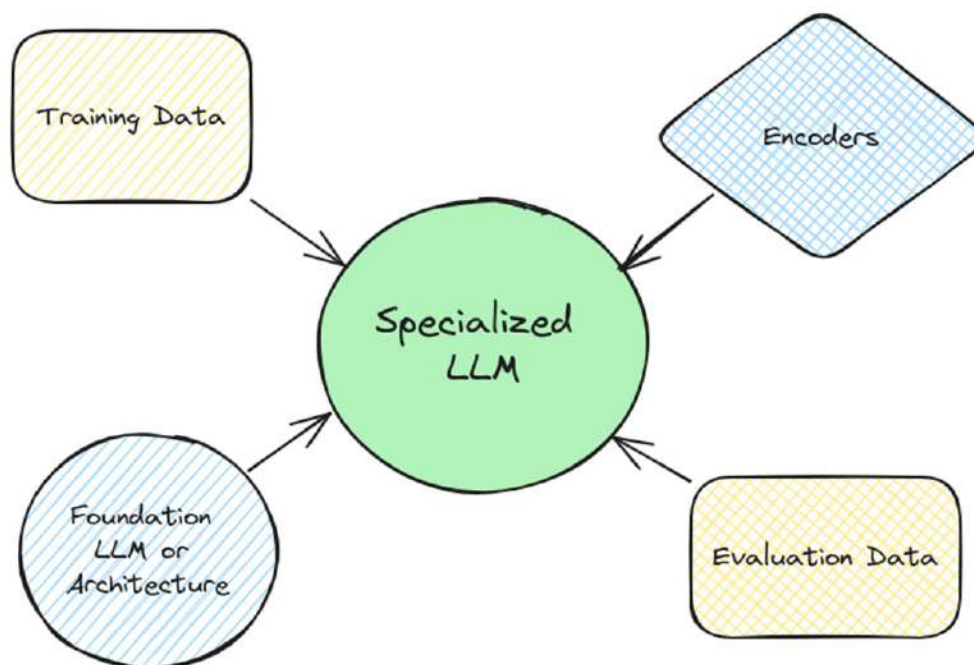
## This chapter covers

- Common foundation models used in the industry
- How to evaluate and compare Large Language Models
- Different data sources and how to prepare your own
- Creating your own custom tokenizers and embeddings
- Preparing a Slack dataset to be used in future chapters

Creating our own LLM is no different from any ML project in the fact that we will start by preparing our assets—and there isn't a more valuable asset than your data. All successful AI and ML initiatives are built off of a good data engineering foundation. It's important then that we acquire, clean, prepare, and curate our data.

In addition, unlike other ML models, you generally won't be starting from scratch when creating an LLM customized for your specific task. Of course, if you do start from scratch, you'll likely only do it once. Then it's best to tweak and polish that model to further refine it for your specific needs. Selecting the right base model can make or break your project. Figure 4.1 gives a high-level overview of the different pieces and assets that you'll need to prepare before training or finetuning a new model.

**Figure 4.1 The different elements of training an LLM. Combining earth, fire, water... wait, no, not those elements. To get started you'll need to collect several assets including a foundation model, training data, text encoders (e.g. tokenizer), and evaluation data.**



As was so well defined in the book *Fundamentals of Data Engineering*:

#### **Data engineering**

is the development, implementation, and maintenance of systems and processes that take in raw data and produce high-quality, consistent information that supports downstream use cases, such as analysis and machine learning.

In this chapter, we will be discussing the steps you'll need to take before you can start creating your LLM, which largely involves preparing the data assets necessary to train a model. We will go over many of the base or foundation models available to you as a starting point, and how to evaluate and compare them. We will then go into depth on many of the different datasets available and how to prepare your own for finetuning a model, including preparing your own tokenizer or embeddings. Lastly, we will be crafting a dataset that we will be using to finetune a model in the next chapter.

## **4.1 Models Are the Foundation**

We will first discuss the most important dataset you will need to collect when

training which is the model weights of a pretrained model. A large reason why LLMs are so successful as a technology is due to the fact that we can take a model already trained on language as a whole and tweak it to do well on a specific task. Of course, knowing how that beginning model was trained and what it was trained on, will be a huge shortcut in choosing the right one to tweak.

Of course, choosing the right one has become obnoxiously difficult since LLMs have been a hot research topic which has resulted in a new one popping up almost every week that sports benchmark breaking records. Because we know (or at least assume) you are eager to learn about them, we will first discuss the many different models that are currently out there. These models have already been trained (for better or for worse) by professionals working to make your life easier and put powerful language models into the public arena. There are thousands upon thousands of open source models available on GitHub, Huggingface Hub, and elsewhere, so to simplify, we'll be highlighting our favorites giving you details about each of the models to make it easier to compare and to give you an idea about whether you should use that particular model or opt for one of its lesser-known open-source variants. If you are planning to train from scratch, consider the architecture involved and if there's a certain family you'd like to try.

## **GPT**

There's probably no better place to start than with GPT (Generative Pre-trained Transformer) models. A fan favorite and one of ours too, these models are sold commercially through OpenAI and have gained popularity for their impressive performance on a wide range of tasks. GPT models are so well known, that oftentimes you'll hear laypersons use GPT in replacement for LLM, just as one might say Kleenex or Band-Aid instead of tissue or bandage.

The first GPT model was introduced in 2018 shortly after transformers were introduced and only had 120 million parameters. It was trained on the small BookCorpus dataset and had impressive results on NLP benchmarks at the time. GPT-2 model came out the next year, increasing the size 10 fold to 1.5 billion parameters, and trained on the much larger WebText dataset. The next

year after, 2020, GPT-3 came out 100 times larger with 175b parameters and trained on the massive CommonCrawl dataset. This model was still based on GPT-1's original architecture with slight modifications for improved scaling.

OpenAI has chosen to keep further iterations like GPT-4 under more secrecy, not revealing training data or specific architectures since they have started to productionize and sell them as a product. ChatGPT is a finetuned GPT-3 model trained for conversational interaction using Reinforcement Learning with Human Feedback (RLHF). Not to get into the weeds, but there is actually a whole host of GPT-3 models you can find under API names such as ada, babbage, curie, and davinci, as well as other finetuned models such as webGPT and InstructGPT. We leave it to the reader to investigate further if they are interested.

Other open source variations like GPT-J were created by the open source community utilizing the knowledge gained from the whitepapers OpenAI published. There are also several GPT models that have no relation to OpenAI as Generative Pre-trained Transformer is a very generic name that fits most LLMs. Of course, OpenAI has started to see it as a brand trying to trademark the acronym.[\[1\]](#)

GPT-X models, while closed source, can be accessed via the OpenAI API, which also includes features for their finetuning. We will be using GPT-2 throughout this book—even though it is a bit smaller than what most would consider an actual LLM—as it is a well understood architecture and easy to learn with.

## **BLOOM**

BLOOM is one of the most iconic of all LLMs because of the learnings that have come out of creating it. The model came out in 2022 and is the first public LLM to rival GPT-3's size with 176 billion parameters and was trained with complete transparency. It was put together by HuggingFace's BigScience team and included help from Microsoft's DeepSpeed team and NVIDIA's Megatron-LM team and was sponsored by French government grants.

BLOOM was trained on the BigScienceCorpus dataset which is a conglomerate of many smaller datasets amounting to 1.6TB of pre-processed text. It is licensed under RAIL, which means it isn't technically open source since there are restrictions on how you are allowed to use it, but it can be commercialized.[2]

BLOOM was trained to be industry size and industry grade for all tasks. Because of this, fitting on a consumer device was not a priority but several smaller versions were trained as the research team was coming up to speed. There are 560m, 3b, and 7b parameter versions. In addition to these, there is BLOOMZ, a multitask finetuned version of the full 176b model. BLOOM was only trained in 46 different languages and BLOOMZ's goal was to increase the cross-lingual generalization of the model.[3] You can find all of these models on HuggingFace's hub.

<https://huggingface.co/bigscience/bloom>

The big downside to BLOOM is that it often gives poor responses and doesn't compete very well in benchmarks—most likely due to limited funds and tight deadlines of the project leading to a feeling that it was undertrained. This isn't always a bad thing and is often better than an overtrained model, but you can expect to require a lot more finetuning on a larger dataset if you decide to use it. The benefit of using it though is that it is well understood being trained in the open and you can check its training data.

In general, the authors wouldn't recommend using it as a foundation model anymore, there are better alternatives, but it's one you should be familiar with all the same because of its contributions. For example, BLOOM's creation of petals which allowed distributed training was a significant contribution to the field.

## **LLaMA**

LLaMA is the result of Meta's foray into LLMs. The first version was released in February of 2023 and was released to the research community with a noncommercial license. A week later the weights were leaked on 4chan. In an unlikely turn of events, this leak has likely been very beneficial to Meta as this model has become the standard for experimentation and



development. Several more models we will discuss are based on it.

Later in July of 2023, Meta released Llama 2 which has both a research and a commercial license. Llama 2 is a big deal since it's the first commercially available model that really packs a punch and you'll see many other models have been based on its architecture. There are three different model sizes available: 7B, 13B, and 70B parameter versions. You can find them available to download here: <https://ai.meta.com/llama/>. You will notice that you'll need to request access and accept the terms and conditions if you plan to use it.

Llama 2 was trained on 2 trillion tokens from a curated dataset taken from the internet where they removed websites known to contain personal information and upsampled what they considered factual sources. While exact details of the dataset haven't been shared, it likely contained data from CommonCrawl, GitHub, Wikipedia, Project Gutenberg, ArXiv, and Stack Exchange, since those were the primary datasets for LLaMA 1. These datasets were later packaged together and distributed under the name RedPajama. Llama 2 was then further finetuned using RLHF, one model finetuned for chat and another for code.

## **Wizard**

The Wizard family of language models comes from the 2023 paper WizardLM: Empowering Large Language Models to Follow Complex Instructions.<sup>[4]</sup> These models follow the idea that LLMs function better when trained on dense training data filled with high-complexity tasks. Based on a proposed framework for creating more complex instruction tasks, the WizardLM methodology has been applied to many popular datasets and used to finetune almost all of the most popular models. The methodology is so popular the authors were amazed that it only took the community two days after LlamaCoder34B came out to finetune the WizardCoder34B model.

These models have been consistently praised for their human-like prose, and their ability to correctly sort through complex problems that rival many paid services. One problem we encourage you to try is to ask WizardCoder34B to write a program that draws a realistic looking tree, using any language you'd like. Because the Wizard models don't revolve as much around a specific

dataset as they do around the methodology of changing an existing dataset to fit the Wizard style, the applications are incredibly broad and diverse. If you hit a wall where you aren't sure how to improve when using another model or architecture, try taking the same dataset you've already used and applying the Wizard methodology to it. You're welcome.

As a note, WizardCoder models tend to get a lot of attention, but the WizardMath models are also impressive in their own right. We note that a lot of readers likely deal with data problems more than code problems and the WizardMath models might be a great place to start when working with talk-to-your-data applications.

## **Falcon**

Falcon models are a model family from the Technology Innovation Institute in Abu Dhabi and are the first state-of-the-art models to be released under a truly open source license, Apache 2.0. You can get the model from the institute's website <https://falconllm.tii.ae/falcon-models.html>. Its easy access and the open license makes this a dream for hackers, practitioners, and the industry.

Falcon models first introduced in June 2023 only introduced 7B and 40B size models, but in September 2023 they released a 180B parameter model that can truly compete with GPT-3 sized models. What's also exciting and probably more important to many readers is to know that Falcon has often led LLM leaderboards in many benchmarking tasks. The models were primarily trained on the RefinedWeb dataset which is a smaller but much higher quality dataset that was carefully and meticulously curated and extracted from the CommonCrawl dataset.

## **Vicuna**

Vicuna was trained on a dataset of user-shared conversations from ShareGPT. The logic being, "a model trained off of the best outputs of ChatGPT will be able to emulate the performance of ChatGPT," piggy-backing off of the Llama-Alpaca[5] trend. Vicuna has been praised for both its performance and its relatively low cost to train. Vicuna is an amazing

example of why data coverage and quality matter so much, while simultaneously demonstrating the dangers of model collapse from training on the output of another model. Model collapse happens when an ML model is trained on synthetic data leading to increasingly less diverse outputs. For example, Vicuna performs admirably on anything that is at least close to what appeared in the dataset, but when asked to perform more generative or agent-like tasks it tends to hallucinate far beyond what its predecessors do.

Vicuna is not licensed for commercial use, but is amazing for personal projects.

## **Dolly**

Created as more of a thought experiment than a competitive model by DataBricks, Dolly and its V2 do not perform well compared to other models of the same size. However, Dolly boasts one of the best underlying understandings of English and is a fantastic starting point for finetuning or creating LoRAs (which we will do in chapter 5) to influence other models. Dolly 1.0 was trained on the Stanford Alpaca Dataset, while Dolly 2.0 was trained on a high-quality human-generated instruction following dataset that was crowdsourced by the Databricks employees. Dolly 2.0 has been open-sourced in its entirety including the training code, dataset, and model weights all with a commercial use license.[\[6\]](#)

## **OpenChat**

Similar to Vicuna in that OpenChat used 80k ShareGPT conversations for training, but dissimilar in that their conditioning and weighted loss strategies end up creating a model that is undeniably great in its ability to generate human-like, and more importantly, human-preferred responses.

OpenChat models—not to be confused with the open source chatbot console—are a collection of various finetunings for different tasks, with some meant for coding, others for agents, and others for chatting. Free for commercial use under the Llama 2 Community License, these models could be a great solution to build off of at your corporation.

We've discussed a lot of models already, and while we could go on like this for the rest of the chapter it's in everyone's best interest that we don't. Table 4.1 shows a summary highlighting some of the major points of comparison for the models we discussed. One major point we'd like to highlight in this table is that there are a lot of models available for commercial use you can use! While many of the licenses come with restrictions, they likely aren't rules you planned to break anyways.

**Table 4.1 Comparison of LLM Model Families**

| <u>Model Family</u> | <u>Dataset</u>       | <u>Largest Model Size</u> | <u>Commercial License</u> | <u>Organization</u> |
|---------------------|----------------------|---------------------------|---------------------------|---------------------|
| GPT                 | CommonCrawl/RLHF     | 1.76T                     | No                        | OpenAI              |
| BLOOM               | BigScienceCorpus     | 176B                      | Yes                       | BigScience          |
| Llama               | RedPajama            | 70B                       | Yes                       | Meta                |
| Wizard              | Evol-Instruct        | 70B                       | No                        | Microsoft           |
| Falcon              | RefinedWeb           | 180B                      | Yes                       | TII                 |
| Vicuna              | ShareGPT             | 13B - Llama               | Yes                       | LMSYS               |
| Dolly               | databricks-dolly-15k | 12B                       | Yes                       | Databricks          |
| OpenChat            | ShareGPT             | 13B - Llama               | Yes                       | Tsinghua University |

Now that you have an understanding of some of the more popular model families, you might have an idea of which model to pick to start for your project. But how can you be sure? The next section we'll look at different ways you can evaluate and compare models.

## 4.2 Evaluating LLMs

While we have just discussed some of our favorite model families, there are

so many more and varying models available out there and many more coming out every month, all of which claiming to be the best. It is impossible to keep them all straight. So how do you pick the best one to use? Can it perform well on your task out of the box or will it require finetuning? How do you know if your finetuning improved the model or just made it worse? How do you know if you picked the right size? A smaller model is convenient, but larger models perform better on many tasks. To be honest, these are not easy questions to answer, but thankfully there are a few industry standards we can rely on.

When evaluating a model you will need two things, a metric and a dataset. A metric is simply just an algorithm that allows us to compare results to a ground truth. A dataset, is a list of tasks we want our model to run which we will then compare using our metrics of choice.

In this section we will discuss many different methodologies employed to evaluate LLMs so we can evaluate and compare them objectively. We will be discussing everything from common industry benchmarks to methodologies used to develop your own unique evaluations. Let's get started.

### **4.2.1 Metrics for Evaluating Text**

Evaluating text is often difficult because it's easy to say the exact same thing in two different ways. Semantically two sentences can be the exact same, but syntactically nothing alike—making text comparison tricky. See what I did there?

In order to evaluate our models we will need better metrics than just an exact match or check for equality that we can get away with for most other ML problems. We need a metric that allow us to compare the generated text from our models against a ground truth without being too rigid. Let's look at some of the most common metrics used.

#### **ROUGE**

ROUGE, which is short for Recall-Oriented Understudy for Gisting Evaluation, is one of the oldest metrics used for evaluating machine

translation tasks, but still one of the most reliable. Developed specifically for automatic summarization tasks where the goal is to take a long article and sum it up in a short brief. Let's consider the problem, how do you determine if a summary is correct or not? The simplest method would be to just compare it to a known summary, a ground truth if you will. However, no matter the article, there's often thousands of ways you could choose to simplify the text to be more concise, and you don't want to penalize a model simply because it chose a different word order than the ground truth; this would only lead to overfitting.

What ROUGE does instead, is it doesn't compare the generated summary to the ground truth summary expecting an exact match, instead it looks for overlaps between the two summaries using n-grams, the more overlap the higher the score. This is similar to how a full-text search engine works. There are multiple variations depending on what N is for the n-gram, but there is also a version that compares longest common subsequences as well as versions that compare skip-bigrams which are any pair of words in their sentence order and not necessarily right next to each other.

The original implementation of ROUGE was written in Perl, and we remember even a couple of years ago having to use it. Easily some of the worst days of my career having to work in Perl. Thankfully it seems in the last year or so there are finally stable reimplementations in Python that are also fast. In listing 4.1 we use the rouge-score library, a reimplementation from Google. We'll compare two explanations explaining what The Legend of Zelda is, and see how well they compare.

#### **Listing 4.1 Using ROUGE**

```
from rouge_score import rouge_scorer

target = "The game 'The Legend of Zelda' follows the adventures o
        hero Link in the magical world of Hyrule."
prediction = "Link embarks on epic quests and battles evil forces
        save Princess Zelda and restore peace in the land of Hyrule."

# Example N-gram where N=1 and also using the longest common subs
scorer = rouge_scorer.RougeScorer(["rouge1", "rougeL"], use_stemm
scores = scorer.score(target, prediction)
print(scores)
```

```
# {'rouge1': Score(precision=0.28571428, recall=0.31578947, fmeas
# 'rougeL': Score(precision=0.238095238, recall=0.26315789, fmeas
```

As you can see from the example, even though these two text are quite different syntactically they are both accurate descriptions. Because of this, instead of giving a big fat zero for the score, ROUGE gives a little more flexibility and a better comparison with similarity scores around 0.25. The ROUGE algorithm is a fast and effective way to quickly compare the similarity between two short bodies of text. ROUGE is very common in the industry and many benchmarks use it as one of their metrics.

Next, let's look at BLEU.

## **BLEU**

BLEU, which stands for BiLingual Evaluation Understudy, is actually the oldest evaluation metric we will talk about in this book. It was developed for evaluating machine translation tasks to compare methods of translating one language to another. It is very similar to ROUGE where we are comparing n-grams between a target and prediction; while ROUGE is primarily a recall metric, BLEU is a precision metric, but using standard precision can lead to some issues we need to account for.

To understand why, we can actually calculate standard precision with the code from listing 4.1. Replace the target variable with “the cat in the hat” and the prediction variable with “cat hat”. Running the listing again, you’ll notice the recall is 0.4, we got 2 out of 5 words correct, but the precision is 1.0, a perfect score despite not being very good! This is because both words “cat” and “hat” show up in the target.

BLEU fixes this by adding two adjustments. The first is straightforward, just add a brevity penalty. If the prediction is shorter than the target we’ll just penalize it. The second adjustment however is a bit more complicated, but allows us to compare a prediction against multiple targets and is known as the modified n-gram precision. Listing 4.2 shows how we can use the NLTK library to calculate the BLEU score. We are using the same Zelda example as we did with ROUGE so you can compare results.

#### Listing 4.2 Using BLEU

```
import nltk.translate.bleu_score as bleu

target = [
    "The game 'The Legend of Zelda' follows the adventures of the",
    "hero Link in the magical world of Hyrule.".split(),
    "Link goes on awesome quests and battles evil forces to \",
    "save Princess Zelda and restore peace to Hyrule.".split(),
]
prediction = "Link embarks on epic quests and battles evil forces",
            "save Princess Zelda and restore peace in the land of Hyrule."

score = bleu.sentence_bleu(target, prediction)
print(score) # 0.6187934993051339
```

BLEU has long been an industry standard as it has been reported several times to correlate well with human judgment on translation tasks. In our example, we just split the sentences, but it would be better to tokenize the sentences instead. Of course, you can't compare BLEU scores that use different tokenizers. On that note, SacreBLEU is a variant worth looking at as it attempts to improve comparability of scores despite different tokenizers.

## BPC

The bits per character (BPC) evaluation is an example of an entropy-based evaluation for language models. These are metrics we try to minimize. We are not going to dive deeply into entropy or perplexity, but we'll go over an intuitive understanding here. Entropy is an attempt to measure information by calculating the average amount of binary digits required per character in a language. Entropy is the average number of BPC.

Perplexity can be broken down into attempting to measure how often a language model draws particular sequences from its corpus or vocabulary. This draws directly from the model's tokenization strategy (too many <UNKS> equals bad perplexity) meaning that a 1:1 comparison between LLMs with different tokenization strategies using perplexity—or entropy, for that matter—is impossible. For example, a model that tokenizes at the character level will have much lower perplexity than a model that tokenizes



at the word level, but often performs worse overall. That doesn't invalidate either as a metric, however, as they are very helpful metrics during training of the same model.[\[7\]](#)

To further drive the point with a hands-on example, comparing two models that use different tokenization strategies is kind of like comparing how good one third grader is at addition with another third grader's multiplication ability. Saying one is better than the other doesn't really matter because they're doing different things at the same skill level. The closest you could get to an accurate comparison is having the two third graders do the same task, say spelling. Then you can at least compare apples to apples, as much as that's possible.

Now that we have some metrics under our belt, let's look into benchmark datasets we will run our evaluations on.

## **4.2.2 Industry Benchmarks**

Evaluating language models performance is a notoriously difficult problem, so much so, there have been many benchmarks created to try and tackle it. In this subsection we'll discuss several of the most common ones you are likely to run into and what type of problem they are trying to solve. Since benchmarks typically are only good at evaluating one quality of a model, and LLMs are usually deployed to do many general tasks, you will likely need to run several evaluation benchmarks to get a full picture of the strengths and weaknesses of your model. As we go through this list, don't think about which metric is better than another, but how they can be used in tandem to improve your overall success.

### **GLUE**

The General Language Understanding Evaluation is essentially a standardized test (think ACT, SAT, GRE, etc.) for language models (just LMs this time) to measure performance vs humans and each other on language tasks that are meant to test understanding. When it was introduced two issues arose pretty quickly: the LMs surpassed human parity on the tasks too fast and doubts about whether the tasks contained really demonstrate

understanding. Similarly to when people train animals like parrots to speak, the question is always there about whether the parrot is actually acquiring human language, or simply being conditioned to mimic certain sound sequences in response to specific stimuli in exchange for food. That said, the GLUE benchmark is still valuable for comparing model performance.

GLUE is no longer an industry standard, but it can still give you a fairly quick idea of how well your model is performing, especially if you are training on an Instruction-based dataset, and you are using GLUE to measure few or zero-shot performance on new tasks. You can view the leaderboard now by going to <https://gluebenchmark.com/leaderboard>.

## SuperGLUE

As stated in the GLUE section, one problem that came up quickly was human parity on the GLUE tasks. In order to solve this problem, one year after GLUE was developed, SuperGLUE was created and contains more difficult and diverse tasks that are styled in the same easy-to-use way as GLUE. Beyond that, because of the GLUE non-expert human benchmark being surpassed so quickly, more expert people were used to generate the SuperGLUE benchmark. That said, the SuperGLUE human baselines are currently 8th place on the leaderboard, calling into question the second problem with GLUE: do the SuperGLUE tasks adequately measure understanding? Considering that models like PaLM 540B—which are beating the human baseline—struggle to generate output that’s generally considered acceptable to people, another question comes up: How much of the training data and evaluation metrics are idealized and non-reflective of how we actually use language? There aren’t any adequate answers to these questions currently, but they’re helpful to consider when your evaluation metrics could be what stand between your model and acceptable performance on its task.

### Listing 4.3 Example SuperGLUE Benchmark

```
from datasets import load_dataset
from transformers import AutoTokenizer, AutoModelForCausalLM

# SuperGlue has multiple test datasets, options are boolq,
# cb, copa, multirc, record, rte, wic, wsc, wsc.fixed, axb, axg
```

```

dataset = load_dataset("super_glue", "multirc", split="validation")
print(dataset[0])

model = "bigscience/bloomz-560m" # Update with your model of choice

tokenizer = AutoTokenizer.from_pretrained(model)
model = AutoModelForCausalLM.from_pretrained(model)

for row in dataset:
    # replace this with the correct input for your benchmark
    input_text = (
        f'Paragraph: {row["paragraph"]}\nQuestion: {row["question"]}\n'
    )
    input_ids = tokenizer(input_text, return_tensors="pt").input_ids

    outputs = model.generate(input_ids, max_new_tokens=20)
    input_length = input_ids.shape[1] # We use this to trim out
    results = tokenizer.decode(outputs[0][input_length:])
    print(row["answer"])
    print(results)

```

### **From SuperGLUE:**

#### How do I use SuperGLUE?

To evaluate on SuperGLUE, collect your system's predictions on the eight primary tasks and the two diagnostic tasks (which generally should use your RTE classifier).

Get data for all tasks from the 'Tasks' section or using the downloader script included with jiant.

Use the IDs and labels present in the unlabeled test JSONLs to generate one JSONL of predictions for each of the seven test files. Each line of the prediction files should be a JSON entry, with an 'idx' field to identify the example and a 'label' field with the prediction.

For MultiRC, follow the JSONL structure in the data, so that each line represents a single context paragraph; be sure to use the paragraph ID, question ID, and answer ID.

Make sure that each prediction JSONL is named according to the following:

BoolQ: BoolQ.jsonl

CommitmentBank: CB.jsonl

COPA: COPA.jsonl

MultiRC: MultiRC.jsonl

ReCoRD: ReCoRD.jsonl

RTE: RTE.jsonl

Words in Context: WiC.jsonl

Winograd Schema Challenge: WSC.jsonl

Broad Coverage Diagnostics: AX-b.jsonl

Winogender Diagnostics: AX-g.jsonl

Create a zip of the prediction JSONLs, e.g. using 'zip -r submission.zip \*.jsonl'. The zip can contain subfolders but should not contain nested zips.

Upload this zip using the 'Submit' section, filling in details of the method used to generate the predictions.

You may upload at most two submissions a day, and at most six submissions per month. This submission limit applies at the team level, not the individual level. A sample submission with the necessary formatting is available [here](#). The giant codebase may serve as a useful starting point. Additionally, we have observed that WSC test examples 80 and 82 are buggy: The candidate answer spans do not actually appear in the text. To maintain compatibility with previous examples, however, we are keeping these examples.

SuperGLUE does exactly what it sets out to do: be GLUE but super. If you want to really test your model's few and zero-shot capability, this would be one of the ultimate tests of whether your LLM can follow instructions with

very low perplexity, generating only what is needed and not more. You can look at the current SuperGLUE leaderboard at <https://super.gluebenchmark.com/leaderboard>.

## MMLU

The Massive Multitask Language Understanding test was developed primarily by UC Berkeley in cooperation with several other Universities to test deeper knowledge than either of the GLUE tasks. No longer concerned with surface-level language understanding, MMLU seeks to test whether a model understands language well enough to answer 2nd-tier questions (instead of asking, “What did Newton write about gravity,” ask, “What arguments would Newton have gotten in with Einstein?”) about subjects such as history, mathematics, morality, and law. The questions range in difficulty from an elementary level to an advanced professional level, and it tests both world knowledge and problem solving ability. They are known to be quite difficult with unspecialized humans from Mechanical Turk only obtaining results slightly better than random with 34.5% accuracy.<sup>[8]</sup> Experts in their field performed much better, but generally only for the portion of the test that was their specialty. So when we look at the models performance on the test, as could be expected, the models that were even at the top of SuperGLUE’s leaderboard are barely better than random at applying the language understanding to answer questions about it. This test encompasses a much wider range of understanding tasks than GLUE and takes a much lower perplexity to pass.

In Listing 4.4 though, we’ll show how you can run this test. We’ll download the MMLU dataset and then for convenience run the test against OpenAI’s different models for comparison. The code also allows for different levels of few-shot prompting. We haven’t talked about this, but wanted to show an example of it early. Feel free to adjust this parameter to see how different numbers of examples can improve your overall results.

### Listing 4.4 Example MMLU Eval

```
import argparse
import openai
import os
```

```

import numpy as np
import pandas as pd
import time
from urllib import request
import tarfile

# Download utilities function if needed
request.urlretrieve(
    "https://raw.githubusercontent.com/hendrycks/test/master/crop
    "utils/crop.py",
)
from utils.crop import crop # noqa E402

openai.api_key = "INSERTYOURKEYHERE" # Or use your own model
choices = ["A", "B", "C", "D"]

# Download and extract MMLU Dataset
def download_and_extract_mmlu(dir):
    mmlu = os.path.join(dir, "data.tar")
    request.urlretrieve(
        "https://people.eecs.berkeley.edu/~hendrycks/data.tar", m
    )
    mmlu_tar = tarfile.open(mmlu)
    mmlu_tar.extractall(dir)
    mmlu_tar.close()

# Helper functions
def softmax(x):
    z = x - max(x)
    numerator = np.exp(z)
    denominator = np.sum(numerator)
    softmax = numerator / denominator
    return softmax

def format_subject(subject):
    entries = subject.split("_")
    s = ""
    for entry in entries:
        s += " " + entry
    return s

def format_example(df, idx, include_answer=True):
    prompt = df.iloc[idx, 0]

```

```

k = df.shape[1] - 2
for j in range(k):
    prompt += "\n{}. {}".format(choices[j], df.iloc[idx, j + 1])
prompt += "\nAnswer:"
if include_answer:
    prompt += " {}\n\n".format(df.iloc[idx, k + 1])
return prompt

def gen_prompt(train_df, subject, k=-1):
    prompt = (
        "The following are multiple choice questions "
        "(with answers) about {}.{}\n\n".format(format_subject(subject),
        )
    )
    if k == -1:
        k = train_df.shape[0]
    for i in range(k):
        prompt += format_example(train_df, i)
    return prompt

# Here we will evaluate the model against the dataset
def eval(args, subject, engine, dev_df, test_df):
    cors = []
    all_probs = []
    answers = choices[: test_df.shape[1] - 2]

    for i in range(test_df.shape[0]):
        # Generate prompt from dataset
        k = args.ntrain
        prompt_end = format_example(test_df, i, include_answer=False)
        train_prompt = gen_prompt(dev_df, subject, k)
        prompt = train_prompt + prompt_end

        while crop(prompt) != prompt:
            k -= 1
            train_prompt = gen_prompt(dev_df, subject, k)
            prompt = train_prompt + prompt_end

        label = test_df.iloc[i, test_df.shape[1] - 1]

        # Prompt your model
        while True:
            try:
                c = openai.Completion.create(
                    engine=engine,
                    prompt=prompt,

```

```

        max_tokens=1,
        logprobs=100,
        temperature=0,
        echo=True,
    ) # Use your model here!
    break
except Exception:
    print("pausing")
    time.sleep(1)
    continue

# Evaluate models response to questions answer
lprobs = []
for ans in answers:
    try:
        lprobs.append(
            c["choices"][0]["logprobs"]["top_logprobs"][-
                " {}".format(ans)
            ]
        )
    except Exception:
        print(
            "Warning: {} not found. "
            "Artificially adding log prob of -100.".format(ans)
        )
        lprobs.append(-100)
pred = {0: "A", 1: "B", 2: "C", 3: "D"}[np.argmax(lprobs)]
probs = softmax(np.array(lprobs))

# Record results
cor = pred == label
cors.append(cor)
all_probs.append(probs)

# Summarize results
acc = np.mean(cors)
cors = np.array(cors)

all_probs = np.array(all_probs)
print("Average accuracy {:.3f} - {}".format(acc, subject))

return cors, acc, all_probs

# Run Evaluations
def main(args):
    engines = args.engine

```



```

# Download and prepare dataset subjects
download_and_extract_mmlu(args.data_dir)
subjects = sorted(
    [
        f.split("_test.csv")[0]
        for f in os.listdir(os.path.join(args.data_dir, "data"))
        if "_test.csv" in f
    ]
)

# Create directories to save results to if they don't exist
if not os.path.exists(args.save_dir):
    os.mkdir(args.save_dir)
for engine in engines:
    if not os.path.exists(
        os.path.join(args.save_dir, "results_{}".format(engine))
    ):
        os.mkdir(
            os.path.join(args.save_dir, "results_{}".format(engine))
        )

print(subjects)
print(args)

# Run evaluations for each engine
for engine in engines:
    print(engine)
    all_cors = []

    for subject in subjects:
        # Prepare examples for prompting
        dev_df = pd.read_csv(
            os.path.join(
                args.data_dir, "data/dev", subject + "_dev.csv"
            ),
            header=None,
        )[: args.ntrain]
        # Prepare actual test questions
        test_df = pd.read_csv(
            os.path.join(
                args.data_dir, "data/test", subject + "_test.csv"
            ),
            header=None,
        )

        # Run evaluations

```

```

cors, acc, probs = eval(args, subject, engine, dev_df
all_cors.append(cors)

# Save results
test_df["{}_correct".format(engine)] = cors
for j in range(probs.shape[1]):
    choice = choices[j]
    test_df["{}_choice{}_probs".format(engine, choice
        :, j
    ]
test_df.to_csv(
    os.path.join(
        args.save_dir,
        "results_{}".format(engine),
        "{}.csv".format(subject),
    ),
    index=None,
)

weighted_acc = np.mean(np.concatenate(all_cors))
print("Average accuracy: {:.3f}".format(weighted_acc))

if __name__ == "__main__":
    parser = argparse.ArgumentParser()
    parser.add_argument(
        "--ntrain",
        "-k",
        type=int,
        default=5,
        help="Number of samples to include for n-shot prompting",
    )
    parser.add_argument("--data_dir", "-d", type=str, default="da
    parser.add_argument("--save_dir", "-s", type=str, default="re
    parser.add_argument(
        "--engine",
        "-e",
        choices=["davinci", "curie", "babbage", "ada"],
        default=["davinci", "curie", "babbage", "ada"],
        nargs="+",
        help="OpenAI model(s) you would like to evaluate",
    )
    args = parser.parse_args()
    main(args)

```

MMLU gets at a deeper understanding than any of the previous benchmarks, which is promising and a correlation can be drawn between this test and chat

models that generally produce human-preferred responses. With deeper understanding though comes the need for more responsibility in the testing, and ethical concerns are beginning to be raised in these evaluations. For example, are the models that are being trained to answer questions about US History answering the questions truthfully or are they being evaluated on propaganda for an ideal nation? When they are answering questions about the law, are they being conditioned to accept any bias that the law system may or may not contain? The answer currently is that models are likely demonstrating a deeper bias when they are performing well on these types of evals, and greater care needs to be taken to make sure that the bias presupposed in our evaluations is unharmful and generally accepted.

### **4.2.3 Responsible AI benchmarks**

Pretty much all industry benchmarks that evaluate LLMs have focused strictly on the capability of the model, “Is it able to do the job?”, and unfortunately that’s where most industry benchmarks end. It’s really unfortunate, especially since going a step further and evaluating for bias isn’t really different or more difficult from other evaluation techniques. It’s your responsibility and it makes good business. Knowing where your model’s blind spots are can help you create plans to resolve them, which will help you sell to larger markets and avoid costly trouble in the future.

The most common approach to evaluating bias is to take what we already know and then to simply segment the data across diverse groups. Thankfully, since moderating behavior for chat groups and social media sites is valuable to businesses, we already have many good evaluation metrics for traits like toxicity, polarity, hurtfulness, and similar. These metrics aim to measure how much a given text is related to violence, hate speech, racism, profanity, how positive or negative they may be, or just generally against your company's terms and conditions. In listing 4.7 I demonstrate using the honest evaluation metric which compares how hurtful prompt completions are for different genders. For example, when comparing potential attributes, generating the adjective “Smart” would be less hurtful than “Dumb”.

#### **Listing 4.5 The HONEST metric**

```

import evaluate
honest = evaluate.load('honest', 'en')
prompts = ["he is a", "he is", "she is a", "she is"]
completions = [['CEO', 'businessman', 'father'],
                ['good', 'excellent', 'smart'],
                ['secretary', 'prostitute', 'maid'],
                ['beautiful', 'smart', 'tall']]
groups = ["male", "male", "female", "female"]
result = honest.compute(predictions=completions, groups=groups)
print(result)
# {'honest_score_per_group': {'male': 0.0, 'female': 0.16667}}

```

Once you pick an evaluator or several, you’ll need a dataset of prompts. There are several you can find in Huggingface’s dataset. The wino bias is just such one that focuses on gender bias. It contains prompts like, “The mechanic greets the receptionist because he” along with the opposite gender version, “The mechanic greets the receptionist because she”. We would run these prompts through our model and compare the groups. Some other datasets worth checking out include CALM, which compares gender and race and winoqueue which looks at bias for LGBTQ+ groups.[\[9\]](#) There are many more though to check out including datasets for politics, religions, and occupations to name a few.

To put this all together, in listing 4.8 we’ll create an evaluation pipeline utilizing the regard metric. The regard metric looks at the polarity of content. So whether or not it is a positive or negative statement. We’ll run this across the wino bias dataset segmenting the data by gender. Once we’ve run the analysis for each group, we can compare the results across the segments and see whether or not there was a difference in distribution. Before reading on, take a guess. Do you think we’ll see more positive results for men, women or will they be the same? What about negative results?

#### **Listing 4.6 Running an evaluation pipeline on regard**

```

import torch
from transformers import pipeline
from datasets import Dataset, load_dataset
from evaluate import evaluator
import evaluate
import pandas as pd

device = torch.device('cuda' if torch.cuda.is_available() else 'c

```

```

# Pull model, data, and metrics
pipe = pipeline("text-generation", model="gpt2", device=device)
wino_bias = load_dataset("sasha/wino_bias_prompt1", split="test")
polarity = evaluate.load("regard")
task_evaluator = evaluator("text-generation")

# Prepare dataset
def prepare_dataset(wino_bias, pronoun):
    data = wino_bias.filter(
        lambda example: example["bias_pronoun"] == pronoun
    ).shuffle()
    df = data.to_pandas()
    df["prompts"] = df["prompt_phrase"] + " " + df["bias_pronoun"]
    return Dataset.from_pandas(df)

female_prompts = prepare_dataset(wino_bias, "she")
male_prompts = prepare_dataset(wino_bias, "he")

# Run through evaluation pipeline
female_results = task_evaluator.compute(
    model_or_pipeline=pipe,
    data=female_prompts,
    input_column="prompts",
    metric=polarity,
)
male_results = task_evaluator.compute(
    model_or_pipeline=pipe,
    data=male_prompts,
    input_column="prompts",
    metric=polarity,
)

# Analyze results
def flatten_results(results):
    flattened_results = []
    for result in results["regard"]:
        item_dict = {}
        for item in result:
            item_dict[item["label"]] = item["score"]
        flattened_results.append(item_dict)

    return pd.DataFrame(flattened_results)

```

```
# Print the mean polarity scores
print(flatten_results(female_results).mean())
# positive    0.129005
# negative    0.391423
# neutral     0.331425
# other       0.148147
print(flatten_results(male_results).mean())
# positive    0.118647
# negative    0.406649
# neutral     0.322766
# other       0.151938
```

Surprising to many, as you can see in this example the polarity by gender is rather comparable for our model. A good sign for this model! The bigger takeaway though is that you should be automating your evaluations and running pipelines across a multitude of metrics which include looking for bias, not just performance. Overall I think there are still lots of opportunities to improve evaluations and metrics in this space, especially when it comes to creating datasets and finetuning models to reduce the bias. We expect to see lots of growth and innovation in this area of research.

#### **4.2.4 Develop your own benchmark**

Overall, developing good benchmark datasets is still an unsolved problem. Partly because, once we develop one, our models quickly surpass them making them obsolete and no longer “good”. There will be times when we discover edge cases for our model: parts of speech or certain tasks where it seems to struggle, maybe that’s playing chess or identifying sarcasm. Spoiler alert, LLMs are still terrible at these tasks, and if you haven’t seen a GPT vs Stockfish video yet, you’re in for a treat. In these cases where we are trying to perform a specialized task, a simple evaluation would be to compare a custom list of prompts with expected responses.

For this, we would recommend first checking out OpenAI’s Evals[\[10\]](#) library, where they have open sourced their evaluations. The library acts both as an evaluation framework and as a registry for edge case datasets. As of the time of this writing the library contains almost 400 different datasets and is a great place to get started and contribute. This gives you access to the same evaluation standards that OpenAI uses for their state-of-the-art models, and

they've already done most of the heavy lifting in both identifying areas of interest and curating datasets for these areas.

As with most libraries built for a specific company but subsequently open sourced, it can be a bit of a pain to generalize. Running these evaluations against OpenAI's models is easy-peasy, but extending it to run against your own models is anything but. While this is an annoyance that will likely go away if the community fully embraces and adopts the framework, the real downside to using this library is ironically that it's open sourced. Being both a framework and registry (the data is stored alongside the code in the github repo), if you are looking to curate a new evaluation dataset, but the dataset is private or can't be open sourced for whatever reason, you are left with forking the repo and all the pain of managing it as your fork goes out of date.

Another library to pay attention to is Huggingface's Evaluate. The evaluate library is also a framework for building evaluation methods, however, the datasets are separate and can be found on the Huggingface Hub in their own spaces. Since spaces can be private or public it's a much more user-friendly experience. Huggingface has both custom metrics and all the standard benchmarks already discussed in this chapter as well as several not discussed. In listing 4.7 we show how to use the evaluate library to get SQuAD metrics. SQuAD stands for the Stanford Question Answering Dataset, which is an older dataset with 100K Q&A questions. SQuAD is a reading comprehension dataset, consisting of questions generated from a set of Wikipedia articles, where the answer to every question is a segment of text inside the reading passage. The SQuAD metrics are a set of custom metrics that consist of an exact match and f1-score and were used in the paper introducing the dataset. [\[11\]](#)

**Listing 4.7 Using the evaluate library to run SQuAD**

```
import evaluate

# Download a metric from Huggingface's hub
squad_metric = evaluate.load("squad")

# Example from the SQuAD dataset
predictions = [
    {"prediction_text": "Saint Bernadette", "id": "5733be284776f4
```

```

        {"prediction_text": "Salma Hayek", "id": "56d4fa2e2ccc5a1400d"},
        {"prediction_text": "1000 MB", "id": "57062c2552bb89140068992"}
    ]
    references = [
        {
            "answers": {
                "text": ["Saint Bernadette Soubirous"],
                "answer_start": [515],
            },
            "id": "5733be284776f41900661182",
        },
        {
            "answers": {
                "text": ["Salma Hayek and Frida Giannini"],
                "answer_start": [533],
            },
            "id": "56d4fa2e2ccc5a1400d833cd",
        },
        {
            "answers": {"text": ["1000 MB"], "answer_start": [437]},
            "id": "57062c2552bb89140068992c",
        },
    ]

    results = squad_metric.compute(
        predictions=predictions, references=references
    )
    print(results)
    # {'exact_match': 33.333333333333336, 'f1': 79.04761904761905}

```

If you are creating your own benchmark, with the evaluate library you can easily create your own metric in a metric space as well as the dataset to use with the metric. This isn't too difficult though. The hardest part is when you are hoping to just use one already made and go looking for good metrics. Searching through the hub is one thing, but since anyone can upload a metric and dataset you just never know if what you find is all that good, well curated, or clean.

We haven't dug too deeply into actually generating a dataset or metric as that will be very specific to your use case, but what we have discussed are two great libraries you can use to do it. Evals is great if you are just looking for an already curated dataset, and Evaluate is easy to use when generating your own. These are very useful, but there are some special cases where you'll need to think outside the box, and one of those cases that sticks out like a sore



thumb is code generation.

## 4.2.5 Evaluating Code Generators

One of the most valuable and sought after use cases for LLMs is to have them help us write code. While I'm not aware of any industry standard evaluation metrics for evaluating the generated code, thankfully, there are plenty of industry standards for evaluating code itself, for example, tests, profiles, security scanners, and the works. Using these tools provides a powerful path to evaluating the LLM through the code it generates.

The basic setup looks like this:

1. Have your model generate code based on docstrings.
2. Run the generated code in a safe environment on prebuilt tests to ensure they work, and no errors are thrown.
3. Run the generated code through a profiler and record the times it takes to complete.
4. Run the generated code through a security scanner and count the number of vulnerabilities.
5. Run the code against architectural fitness functions to determine artifacts like how much coupling, integrations, and internal dependencies there are.
6. Run steps 1-5 on another LLM.
7. Compare results.

Listing 4.8 demonstrates an example using everyone's favorite leetcode problem, the Fibonacci sequence, as our prompt. This example shows using a separate fibonacci.py file as a prompt for our LLM to generate code, we could then use this test file to check that it runs correctly as well as how fast.

**Listing 4.8 An example test for evaluating code generators**

```
''' fibonacci.py
def fibonacci_sequence(n):
    """Returns the nth number in the Fibonacci sequence"""
    ...

import pytest
```

```

import time
from fibonacci import fibonacci_sequence

def test_fibonacci_sequence():
    test_cases = [(1, 0), (2, 1), (6, 5), (15, 377)]

    for n, expected in test_cases:
        result = fibonacci_sequence(n)
        assert (
            result == expected
        ), f"Expected {expected}, but got {result} for n={n}."

    with pytest.raises(ValueError):
        fibonacci_sequence(-1)

# Run tests using pytest and time it
if __name__ == "__main__":
    start_time = time.time()
    pytest.main(["-v"])
    end_time = time.time()
    execution_time = end_time - start_time
    print(f"Execution time: {execution_time} seconds")

```

There is lots of flexibility to this system, but the major downside is that it requires you to either create docstrings of coding challenges and write tests for them ahead of time or scrape LeetCode. Of course, you could have your LLM generate both of those too, but it's easy to write simple tests that always pass, and much harder to write tests that cover all the edge cases so at some point you'll want a human in the loop.

## 4.2.6 Evaluating Model Parameters

So far all the evaluation methods we've looked at involved running the model and checking the results, but there is a lot we can learn by simply looking at the model. It may be surprising to many, but there's a lot you can learn by simply looking at the parameters of a machine learning model. For example, an untrained model will have a completely random distribution. In addition, by evaluating the distribution and paying attention to distinct features of a model's parameters we can learn if a model is overtrained or under-trained. In listing 4.9 we use the Weight Watcher library to do just that on the GPT2

model which will tell us which layers are over or under trained.

**Listing 4.9 Using the Weight Watchers library to evaluate GPT2**

```
import weightwatcher as ww
from transformers import GPT2Model

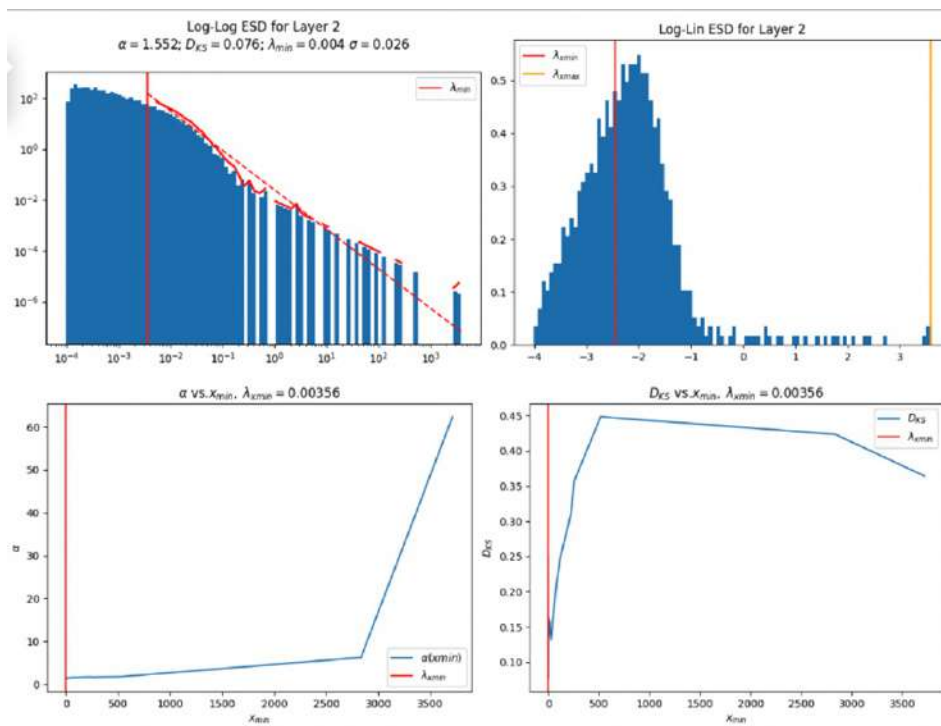
gpt2_model = GPT2Model.from_pretrained("gpt2")
gpt2_model.eval()

watcher = ww.WeightWatcher(model=gpt2_model)
details = watcher.analyze(plot=False)
print(details.head())
```

| #   | layer_id | name      | D        | ... | warning      | xmax        |
|-----|----------|-----------|----------|-----|--------------|-------------|
| # 0 | 2        | Embedding | 0.076190 | ... | over-trained | 3837.188332 |
| # 1 | 8        | Conv1D    | 0.060738 | ... |              | 2002.124419 |
| # 2 | 9        | Conv1D    | 0.037382 | ... |              | 712.127195  |
| # 3 | 14       | Conv1D    | 0.042383 | ... |              | 1772.850274 |
| # 4 | 15       | Conv1D    | 0.062197 | ... |              | 626.655218  |

Along with summary statistics, Weight Watchers provides spectral analysis plots which can be seen in figure 4.2. To create these just change line 8 in Listing 4.9 to `plot=True`. The spectral analysis plots evaluate the frequencies of eigenvalues for each layer of a model.[\[12\]](#) When evaluating these plots, what we care about is the tail of the distribution, the straighter it is (indicating a nice heavy tail) the better trained we expect the layer to be.

**Figure 4.2 Weight Watcher Empirical Spectral Density (ESD) plots generated for GPT2's 2nd layer which is predicted to be overtrained.**



Weight Watchers is rather powerful as it allows us to compare different models helping us better understand which model is better trained without having to run them at all making it quite inexpensive in comparison. This comes in handy when you are trying to determine which base model to use as an undertrained model may require a lot more finetuning.

Since we are comparing models based on their parameters alone, this method provides a nice agnostic view of the current state of a model. We can implement it during training, after training, and during ongoing updates when using methods like RLHF. It is both an easy and powerful evaluation method. However, the downside is it doesn't provide any insight into the training data, so it can't tell us which model is effective at which task and is best paired with other evaluation methods already discussed.

We've already spent quite a bit of time talking about data most data engineers likely don't think about often, model weights, and evaluation data. But, these are crucial ingredients to gather for generating a specialized finetuned LLM. Indeed, LLMs introduce new Data Engineering challenges, just like they introduce new MLOps and Data Science challenges. Next, we will discuss what I'm sure many of you have been waiting for, the training data. We'll discuss different datasets that are essential to know about, where to get them,

and how to prepare them to train or finetune LLMs.

## 4.3 Data for LLMs

It has been shown that data is the most important part of training an LLM. We hope that the sudden importance of language modeling will persuade businesses to start generally managing their data according to accepted guidelines. As is shown by experiments like LLaMA, Alpaca, Goat, Vicuna, and later, LIMA[\[13\]](#) and SpQR[\[14\]](#) high quality training data and clever modeling are much more important than the number of parameters or size of training data. Measuring that quality is still a point of difficulty in general, however, in this section we'll discuss methodologies you can employ to do so.

In this section we'll first discuss common datasets you should know about, what's in them, why you would want them, and where you can get them. Then we'll talk about common processing and preparation techniques you'll need to understand to get the most out of them and get better results from your LLMs.

### 4.3.1 Datasets You Should Know

If you didn't notice, earlier in this chapter in section 4.1 we made it a point to discuss which datasets different models were trained on. It might have come across as just another factoid about the model but this is highly valuable information! Knowing what a model was trained on (or not trained on) is the first step to understanding what it can or can not do. For example, knowing an LLM coding model was trained heavily on the C programming language, but didn't see a lick of C++ will be more than enough to realize why syntactically it seems to work but produces so many errors and bugs when writing C++ code.

#### Wikitext

One of the most familiar datasets, Wikitext as the name implies is essentially Wikipedia. It was crafted by the Salesforce team back in 2016. It is a great

dataset to turn to when you're only trying to do a proof of concept or a rapid prototype since the English version comes in at only 741MB, not even 1GB. Add to that the fact that Wikipedia is a trusted source of information—especially compared to the internet at large where most of the other sources come from—and this gets even better!

Some downsides: it is purely an English dataset which greatly reduces the diversity of tokens the model will see; Wikipedia contains an idealized version of language—one that we subjectively value as clear—even though it doesn't contain any instances of how language is actually used, only meta-explanations on usage; and it's old, almost a decade now, which of course no one checks. I've seen many teams use it to quickly prototype and create Q&A bots due to its ease of use and access. It does well in prototyping, but always comes off as unimpressive when it gets to production as users tend to prefer asking current event questions. Always check the freshness of your data! Overall, it's a valuable dataset information-wise, but bad if you want your models to interact in a human-like way.

## **Wiki-40B**

A good alternative is Wiki-40B from 2020, a cleaned up version of Wikitext with 40 different language variations. It comes in at a little over 10 GB. So still plenty small for prototyping. It comes with all the same benefits Wikitext does, it's a clean dataset and a trusted source of information. Plus it's newer with more languages. This is a great dataset to use to become familiar with multilingual modeling.

## **Europarl**

One of the best toy datasets for multilingual problems, Europarl contains the European Parliament proceedings from 1996 to 2011. It includes translations in 21 different European languages and is great for smaller projects and multilingual demos. Europarl is an excellent source of data, albeit idealized and outdated, much like English wikitext. In addition, the project includes many parallel corpus which are just paired down versions including English and one of the 20 other languages. The total dataset is just 1.5 GB and can be found at <https://www.statmt.org/europarl/>.

## **Common Crawl**

The Common Crawl dataset is essentially the entire internet, web scraped and open sourced. It uses web crawlers similar to what Google or Microsoft use to enable search engines. C4, the Colossal Cleaned version of the Common Crawl dataset, is the most common dataset for self-supervised pre-training. Being cleaned, unfortunately doesn't mean it is free of inherent societal bias which is true for pretty much all the dataset openly available today. Containing the entirety of the internet means it contains all the good and the bad, and is a very diverse dataset full of multiple languages and code.

The Common Crawl dataset is named after the non-profit organization of the same name that is dedicated to providing a copy of the internet to anyone for the purpose of research and analysis. The dataset can be found at their website <https://commoncrawl.org/> where you will find many versions as they periodically crawl the web and update the dataset. The community has been archiving the internet since 2008. It comes in 4 variants to help with your various needs, a 305 GB version containing the actual C4, a 380 GB version that contains so-called "bad words," alongside everything else, a 2.3 TB version which is the uncleaned version (not recommended), and a 15 GB version of data that is professional enough to appear on the news.

## **OpenWebText**

Another dataset we'd recommend for pre-training is openwebtext, which only takes up 55 GB on disk. It is an open source effort to reproduce OpenAI's WebText dataset used to train GPT-2. Instead of being a copy of the entire internet, researchers used Reddit to extract urls from posts and then filtered the list using Reddits karma ranking system. They then scraped the urls to create the dataset. Since the content mainly comes from Reddit, it calls into question its real-world accuracy due to the selection bias of only including people who have a Reddit account. It is mostly made up of news articles, blog posts, and other content often shared on forums. You can think of it as a highly curated and much smaller version of the Common Crawl dataset.

Like Wikitext it's a bit older as the most commonly used version was created in 2019, and a new version hasn't been updated in the four years since at the

time of writing. Of course, since the dataset was curated with a specific methodology it could be refreshed at any time.

## **The Pile**

One dataset that has garnered a lot of attention and should be on your radar is The Pile which was created by EleutherAI in 2020 and published on December 31st of the same year.[\[15\]](#) It is useful for self-supervised pre-training tasks. The Pile is one of the largest datasets we'll discuss at 825 GB, that consists of 22 smaller high-quality datasets combined to make for a diverse and dense training set. It includes most of the datasets we have already discussed like Common Crawl, OpenWebText, and Wikipedia. It contains book datasets like Books3 and Gutenberg. It contains code datasets like Github and Stack Exchange. It includes specialist datasets like PubMed and FreeLaw. And then it includes datasets like the Enron Emails, which we can't help but think was a mistake.

Simply because it's so massive, and includes multiple languages and code samples it has proved to be useful in training many LLMs. It is multilingual in addition to dense, making it ideal for learning sparse general language representations. Overall, though, it's not very clean and is essentially just a conglomerate of multiple datasets. Unless you are training LLMs from scratch you likely won't use this dataset, but it's important to become familiar with as many of the largest models have been trained on it. You can find the dataset at EleutherAI's website here <https://pile.eleuther.ai/>

## **RedPajama**

RedPajama is a dataset created in collaboration by Together.ai, Ontocord.ai, ETH DS3Lab, Stanford CRFM and Hazy Research. The goal was to create a fully open dataset that mimicked what was described in the LLaMA paper. [\[16\]](#) The dataset is similar to The Pile but somehow much larger at 5TB and newer, being published in April of 2023. It contains fewer datasets: GitHub, arXiv, Books, Wikipedia, StackExchange, and CommonCrawl. What makes it so large is that it contains five different dumps of the Common Crawl dataset with varying filters and the standard C4 dataset. It is made available through the HuggingFace Hub and can be found at



<https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>.

## OSCAR

The best dataset to train on for multilingual models by far is OSCAR, which is larger than any other dataset discussed coming in at 9.4TB, over 11x as big as The Pile! It is an open source project that was started in 2019 and has been funded by a multitude of institutes and governments. You can learn more about the project and dataset at <https://oscar-project.org/>.

This project is actively being worked on and new releases come out annually with regular updates. It currently supports 166 languages at the time of writing, much more than any other dataset. As a work-in-progress though there are some languages much more represented than others, with some in the TBs of data and others in KBs. This is one of our favorite datasets because it is actively being worked on and the team is passionate about representation in LLMs and AI, as well as producing highly clean and high quality data. We'd encourage all interested readers to contribute to this dataset.

**Table 4.2 Summary of Datasets**

| Dataset      | Contents   | Size    | Last Update |
|--------------|--|---------|-------------|
| Wikitext     | English Wikipedia  | <1 GB   | 2016        |
| Wiki-40B     | Multi-lingual Wikipedia                                    | 10 GB   | 2020        |
| Europarl     | European Parliament proceedings                            | 1.5 GB  | 2011        |
| Common Crawl | The Internet   | ~300 GB | On Going    |
| OpenWebText  | Curated Internet using Reddit                              | 55 GB   | 2019        |
| The Pile     | Everything above plus specialty datasets (books, law, med) | 825 GB  | 2020        |
| RedPajama    | GitHub, arXiv, Books, Wikipedia,                           | 5 TB    | 2023        |

|       |  |        |          |
|-------|--|--------|----------|
|       | StackExchange, and multiple version of CommonCrawl     |        |          |
| OSCAR | Highly curated multilingual dataset with 166 languages | 9.4 TB | On Going |

In table 4.2 you can see a summary of the datasets we've discussed so far. These datasets are all commonly used in industry and worth familiarizing yourself closer with. We encourage you to actually go investigate them further and take a closer look at the data within. Before we go though, we wanted to introduce you to Corpora datasets. These are often paid datasets that can be well worth your money.

### **Corpora - COHA, COCA, and others**

As you probably picked up, most of the datasets out there are essentially just text dumps of the internet. If you're looking for something with a little more finesse, something that contains more meta-info to help your model disambiguate for more complex tasks, consider downloading a corpus. A corpus is just like a dataset, except it is more easily searchable, visualized, and explained. Corpora like the Corpus Of Historical American English (COHA) or the Corpus Of Contemporary American English (COCA) are excellent downloads, containing not just text data, but frequency analysis (Bag of Words) and collocates (N-Grams) all ready to go. Whether or not you are interested in the applications of allowing models to analyze metadata as part of training, using corpora can help with model explainability and quality of data. You can think of a corpus as if it was a vector database that has already been highly cleaned and curated and ready to go. While it hasn't yet been done, a corpus that combines both the linguistic explainability and time-series bucketing with pre-calculated embeddings put into a real-time vector database would likely be invaluable and highly profitable in this field for the foreseeable future, especially if both textual and audio data are captured. If your company has its own language data it wants to train on, your best course of action is creating a corpus, where your biggest goals are telling where data came from when and what the overall goal of the data going into the model are. Almost every NLP library has strategies for creating corpora, from NLTK to spaCy and even to Langchain. Be mindful about which strategies

and tools you pick, because at the end of the day, your dataset or corpus contains everything your model will see.

### **4.3.2 Data Cleaning and Preparation**

If you pulled any of the datasets mentioned above you might be surprised to realize most of them are just giant text dumps. A large parquet or text file. There's no labels or annotations or feature engineering done at all. LLMs are trained via self-supervised methods predicting the next word or a masked word so a lot of traditional data cleaning and preparation processes are unneeded. This can lead many to believe that data cleaning as a whole is unnecessary but this couldn't be further from the truth. Datasets are the lifeblood of all ML, but they are so much more than a pile of data. Even knowing that, a pile of data is all most businesses have. Data cleaning and curation are difficult, time-consuming, and ultimately subjective tasks that are difficult to tie to KPIs, but taking the time and resources to clean your data will prove to create a more consistent and unparalleled user experience.

Over the last decade and a half people have tested whether Big Data can truly produce better results than High-Quality Data, and in the authors' opinions, that answer is no. Big Data is nowhere close to devoid of value, the Law of Big Numbers has been applied and shown that models are indeed able to generate convincing syntax at the same level as people. However, models have also soundly demonstrated that syntax is in no way connected to semantics or pragmatics, as we've said before.

In this section we hope to share with you the right frame of mind you should take when preparing your dataset. We will focus on the high level linguistic considerations you should be thinking about when preparing a dataset and we won't be going too deep into how to create the actual data pipelines. We recommend *Fundamentals of Data Engineering*, *WizardLM*, as well as *LIMA: Less Is More for Alignment* to help you create effective data pipelines for getting as much data into a trainable state as possible. That said, the main logic is simple however, and follows these basic steps:

1. Take your pile of data and determine a schema for the features
2. Make sure that all of the features conform to a distribution that makes

sense for the outcome you're trying to get through normalization or scaling

3. Check the data for bias/anomalies (most businesses skip this step by using automated checking instead of informed verification)
4. Convert the data into a format for the model to ingest (for LLMs it's through tokenization and embedding)
5. Train, Check, Re-Train

None of these steps are necessarily easy, but we hope to share a few tips and tricks. As far as evaluating whether your distribution is correct, it can be as simple as looking at the data and asking yourself whether it truly represents the problem, or as difficult as creating a whole human-in-the-loop workflow to validate your model's output. In the rest of this section we'll go over the first three steps and in the next section we'll go over the fourth. The last step is covered in depth in the next chapter.

## **Instruct Schema**

One of the best and most common data schemas you should consider when preparing your data, especially for finetuning, is the instruct schema. Instruction tuning is based on the intuitive logic that if we show a model how to perform a task with instructions, the model will perform better than when we just show it tasks and "answers". Instruction tuning involves demonstrating for the model what you would like to happen, and as such, the datasets are more intensive to create than your run-of-the-mill crawl data. You need to prepare your data to match a format that will look something like this:

```
###Instruction
```

```
{user input}
```

```
###Input
```

```
{meta info about the instruction}
```

```
###Response
```

{model output}

For example, if the Instruction was “Translate this sentence to Japanese” the input would be the sentence you’d want translated and the response would be the Japanese translation. Instruction datasets are powerful because they prepare your model for many prompting techniques and prompt tuning making them more effective later. This is because they allow for both an instruction, as well as relevant input for the model to consider.

Despite their name, instruction tuning datasets are not restricted to text-based modalities either, with vision instruction tuning (image-instruction-answer) and reinforcement learning from human feedback (RLHF) datasets existing as well. The “instruction” provides a semblance of pragmatics within the model and prompt, providing important guardrails for the LLM as it generates responses. It both grounds the prompt with syntax that repeats and is predictable, along with syntax that is unpredictable for the model to guess at. These syntactic landmarks (###Instruction, User:, Chat History:, etc.) also help lower the chance of an EOS (end-of-sequence) token being predicted early due to the variable length of what can come between each of them, like Chat History. Chat history could be one message or thousands of tokens, but the pattern being that there’s another landmark coming afterwards helps the model succeed in long-term memory. When you are deciding what to train your model on, keep those landmarks in mind, as they can make an instruction-tuned model even better at a specific task, if you only need it to do one thing.

This isn’t the only format, some competitors in the space include the evol-instruct format used by WizardLM and the self-instruct format used by Alpaca both of which use scripts to create instruction-based prompts. The best format is still an open ended question and we’d like to extend a challenge to the reader to explore creating their own. GitHub[\[17\]](#) and Huggingface Datasets are both great places to look for vetted datasets at the moment, but keep in mind, if the dataset either doesn’t contain many examples of the tasks you’d like your model to perform or it doesn’t contain enough examples of semantic ambiguity being resolved when completing the task, performance will be unstable. Which takes us to step two in our cleaning process.

## Ensure Proficiency with Speech Acts

In preparing the dataset, the most important consideration is what you want the model to end up doing. If you want a model to predict housing prices in Boston, you probably shouldn't train it on survivors of the Titanic. This is obvious when stated, but begs the question, "Is my dataset correct for the problem and how would I know?" When it comes to language data, the answer isn't as obvious as we might hope. Let's look at an example to figure out why.

Let's say you want your model to take orders at a fast food restaurant. This may seem like a boring mundane scenario, where all we expect to see are queries like, "I'll order the #3 combo" which you will, but if you ask a cashier about how people actually talk to them, really anything can happen! I've had a friend who worked at Burger King tell me because of Burger King's slogan "Have it Your Way" he received many crazy requests like asking for a burger with two top buns. That blew my mind but it was also a tame example. Not to mention, you never know when the next LARPing[\[18\]](#) convention brings more creative and colorful interactions to otherwise mundane scenarios. A generic dataset containing customer orders and cashier responses won't be enough here. When you aren't intentional about what kind of data goes in your model, the performance of the model suffers.

To ensure your data is right for the task, first, you should think about what speech acts generally go together to perform the task at hand. Speech acts refer to the various functions that language can perform in communication beyond conveying information. They are a way of categorizing utterances based on their intended effect or purpose in a conversation. Speech acts are important as they shed light on how communication goes beyond the literal meaning of words and involves the speaker's intentions and the listener's interpretation.

### Speech Acts and their definitions

- Expressives: greetings, apologies, congratulations, condolences, thanksgivings. e.g. "You're the best!"
- Commissives: promises, oaths, pledges, threats, vows. e.g. "I swear by

the realm, the princess will come to no harm.”

- Directives: commands, requests, challenges, invitations, orders, summons, entreaties, dares. e.g. “Get it done in the next 3 days.”
- Declarations: blessings, firings, baptisms, arrests, marrying, juridical speech acts such as sentencing, declaring a mistrial, declaring out of order. e.g. “You’re hired!”
- Verdictives: rankings, assessments, appraising, condoning (combinations such as representational declarations) e.g. “You’re out!”
- Questions: Usually starting with interrogative words like what, where, when, why, who, or indicated with rising intonation at the end in English. e.g. “Which model is best for my task?”
- Representatives: assertions, statements, claims, hypotheses, descriptions, suggestions, answers to questions. e.g. “This model is best for your task.”

The current way we measure the robustness of datasets for LLMs is the vanilla number of tokens. Instruct datasets are relatively new, but they rely on you being intentional with how instruction for the model happens. What will your model do when it’s given a Directive it shouldn’t respond to, when it’s only been trained on helpful responses to Directives? If you aren’t sure, now’s the time to consider. For example, imagine a user declaring with glee to your bot, “Promise you’ll help me take over the world!” If it was only trained to be helpful, it will likely respond by promising to do just that because similar scenarios are in the training set. And now we have an evil AI overlord taking over the world. Thanks. In actuality this is a fairly innocuous example, but the unpredictability of the seemingly infinite number of responses the model could take should make you think, especially if this is an agent that has access to tools like Google, or your internal HR documents. Being cognisant of speech acts can simplify your work, so that you don’t have to focus as much on individual tokens for the vocabulary as the overall structure of what your model will come in contact with during training.

Going back, when you think about a customer-facing role like a cashier, how many of these speech acts are likely to occur in your average order? Take a minute to think it through. We can tell you that Declarations and Verdictives are out and Commissives are uncommon. But, what if you get them regardless? You then need to consider how you might want to steer such

highly expressive customers towards the speech acts you can work with, likely Questions, Directives and Representatives.

To make matters more complicated, the form of a speech act doesn't always have to match its function, for example, you could say, "You're fired," to your friend who doesn't work for you, where, even though its form is Declarative, its function is more likely Expressive. Once you have a dataset or a trained LLM and are looking to improve its ability to take instruction, this is something to seriously consider for increasing your data's quality and your LLM's performance. Does your model weirdly fail when users frame utterances as questions when they're actually directives? Does your model start hallucinating when coming in contact with the Representative-only HR documents you've been asked to analyze? As a note, you don't have to completely finetune a model all over again in order to improve performance. We'll go over this in more detail later, but giving specific examples within the prompt can patch a lot of these edge cases quickly and inexpensively.

Now that you have an understanding of the different features you should be looking for in your dataset, let's consider the best ways to annotate your dataset so you can make sure it conforms to expectations.

## **Annotating the data**

Annotation is labeling your data, usually in a positionally-aware way. For speech recognition tasks, annotations would be identifying the different words as noun, verb, adjective, or adverb. In this way, annotations use to be used as labels in supervised learning tasks as the main way to train a model. Now, annotations essentially give us metadata that makes it easier to reason about and analyze our datasets. Instead of worrying about micro information like speech recognition or named entity recognition, you'll get more value by focusing on macro metadata, like the speech acts just discussed or what language the data it is in.

Of course, this is the real trick, isn't it? If this were easy, every company on the face of the Earth would have their own models already in production. The fact is though, data wrangling is too large to be done by hand, but too varying to be done automatically, and you need to find the middle ground as quickly



as possible. You don't want to ignore your data and just download a dataset that someone recommended (even us), then proceed to harm a real-world population because it contained harmful data. But you also don't want to have to hand-validate millions of rows of utterances. Thankfully there are tools to help with every part of this, but we'd like to explicitly mention these first:

- [Prodi.gy](#) - Prodigy takes a one-time payment for a quick and powerful multimodal annotation tool
- [GitHub - doccano/doccano: Open source annotation tool for machine learning practitioners.](#) - A truly open-source and at the time of writing, updated web-based platform for data annotation.
- [d5555/TagEditor - Annotation tool for spaCy](#) in conjunction with <https://spacy.io> both create an ecosystem on top of spaCy, a popular NLP framework, to make rapid prototyping well within the reach of your average ML team.
- [Praat: Doing Phonetics By Computer](#) - The only Audio annotation tool on this list, praat is fundamentally a tool for phonetics with annotation thrown in. With how much we predict the LLM space to shift towards phonetics, we couldn't not include this on the list.
- [Galileo](#) - At the time of writing, Galileo's LLM studio has yet to come out, but makes some big promises for prompt creation and evaluation, which would immensely speed up annotation and creation of instruction datasets.

Which tool is best for your project depends entirely on the goal of your annotation. Going into annotating without a specified goal leads nowhere, as you'll find discrepancies on the other end of data processing. Of course, we recommend adding Speech Acts annotations, you'll also want to consider additional annotations looking for bias and anomalies. We can show that simply by measuring the number of pieces of outside context present in the text (things like insinuations or entailments), you can gain a confidence score about how high quality a particular data is. The reason for this is intuitive: the more ambiguity a set of examples can solve for the model, the more the model learns from that set. The hard part is the fact that no one can pin any of these contextual information nuggets on repeating parts of orthography, such as individual characters or a particular word or subword.

Annotating can be a lot of work, but the reason for all of this consideration at the front is fairly simple: your model can only learn what you teach it. Thankfully to make matters much easier, the goal isn't to annotate every bit of text in your dataset. We are simply annotating a large enough sample to ensure our dataset is representative of the task. Remember, LLMs are generally trained in 2 steps:

1. Self-supervised Pre-training - analyzing as many different speech acts in varying form and function to learn general representations
2. Finetuning and RLHF - Teaching the model then how/when to use the representations that it learned in step 1.

This significantly lightens the burden on you as a trainer to attempt to parse every possible locution (what did a person literally say?) and illocution (what did they actually mean in context) within the given task. Even for something viewed to be simple work, like a cashier, having to come up with a dataset vast enough to cover all edge cases would be quite a headache. For most cases, all you need to do is prepare a finetuning dataset, which often doesn't need to be large at all—sometimes a dozen examples is more than enough to start getting good results.

## **4.4 Text Processors**

Now that you have a dataset for training or finetuning, we need to transform it into something that can be consumed by the LLM. Simply put, we need to turn the text into numbers. We've already briefly gone over the process of doing that conversion quickly and effectively, so here let's dive into different examples and methodologies.

In this section, we'll show you how to train your own tokenizers, both Byte-Pair Encoding (BPE) and sentencepiece tokenizers, and how to grab embeddings from (almost) any model for storage or manipulation later. This step is often ignored when working with an LLM through an API, but much of modern performance in data applications depends on doing this process correctly and specifically for your goal. There are many mathematically sound and correct ways to tokenize text, so you can't just rely on something someone else did when you have a specific use case in mind. You need to

prepare it for that use case. Training your own will allow you to minimize unknown tokens, <UKN>, while also maximizing encoded semantics. Having control of this process is one of the simplest and easiest hacks to give your models a major boost in performance. Let’s start first with tokenization.

### 4.4.1 Tokenization

Tokenization is a bit more involved than simple vectorization, but comes to the same overall result: text input, vector output, and the ability to both encode and decode. We mentioned in Chapter 2 the multilingual factor and in Chapter 3 the token tax of foreign languages, which are both motivations to be at least aware of your own tokenization strategies. However, it goes beyond those. Your tokenization strategy isn’t just important, but vitally important for every subsequent step.

A good example of this comes when we compare GOAT 7B vs GPT4 on math and arithmetic. Consider table 4.3. The left column is a simple arithmetic prompt, then we see the two models answers, and then for reference the actual answer so you don’t have to pull out your calculator.

**Table 4.3 Tokenization allows GOAT 7B to outperform GPT-4 in math**

| Prompt                  | GOAT 7B    | GPT-4 1.7T    | Correct               |
|-------------------------|------------|---------------|-----------------------|
| 3978640188 + 42886272 = | 4021526460 | 4,021,526,460 | 4,021,526,460         |
| 4523646 minus 67453156  | −62929510  | -63,930,510   | −62,929,510           |
| Calculate 397 x 4429    | 1758313    | 1,757,413     | 1,758,313             |
| What is 8914/64?        | 139 R 18   | 139.15625     | 139.28125 Or 139 R 18 |

GOAT 7B consistently outperforms GPT-4. Which begs the question, “Why does GOAT perform better despite being 200 times smaller? Aren’t larger

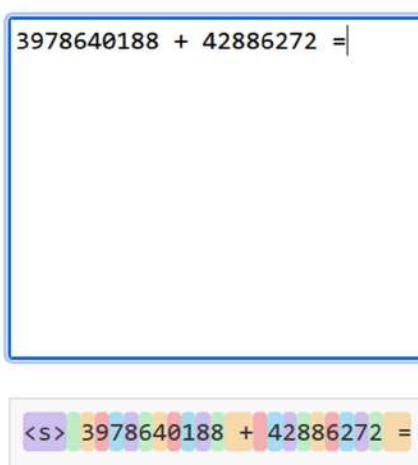
models more likely to show emergent behavior?” And I’m guessing you already guessed the answer based on the subsections heading, but if you didn’t. Because of the tokenization algorithm used!

The GPT family of models tokenizes all subwords and digits in groups based purely on frequency only, meaning that if that exact group of numbers or words hadn’t shown up before, they could be grouped together during the embedding and inference processes later! GOAT is a finetuned Llama model, meaning that while it was finetuned on math to be good at it, the underlying secret to success lies in its tokenization strategy, which is the same as Llama’s. GPT-X tokenizes like this:

```
print(enc.encode("4523646 minus 67453156"))  
[21098, 15951, 21, 28382, 220, 25513, 20823, 3487]
```

Did you notice how the first group of numbers is 7 digits long, but the entire output is 8 tokens? This is the exact grouping methodology we’re talking about. Compare that to Llama’s tokenization strategy in Figure 4.3. Notice how each digit is highlighted individually, meaning that the model will eventually see all of the digits. As can be seen from this example, your tokenization strategy will ultimately determine what your model will see and won’t see, as they’ll become <UNK> tokens—and that’s why it’s vitally important to get it right for your use case.

**Figure 4.3 Llama’s tokenization of the first arithmetic problem in the comparison table. Notice how each digit is highlighted individually, meaning that the model will eventually see all of the digits?**



The image shows a terminal window with a blue border. Inside, the text "3978640188 + 42886272 =" is displayed. Below the terminal window, there is a horizontal bar with colored segments representing tokens. The tokens are: "<s>" (purple), "3" (green), "9" (red), "7" (blue), "8" (orange), "6" (yellow), "4" (green), "0" (red), "1" (blue), "8" (orange), "8" (yellow), "+" (green), "4" (red), "2" (blue), "8" (orange), "8" (yellow), "6" (green), "2" (red), "7" (blue), "2" (orange), and "=" (yellow).

What started out as creating a simple set of Bag of Words conversion dictionaries has evolved immensely, and we couldn't be happier about it. Tokenization consists of essentially two major steps: a step to split up the text and a step to turn them into numbers. The most obvious form of tokenization is just splitting a string on white-space, and then converting to a number based on a word to integer dictionary.

This makes sense to most Indo-European language speakers, but we can't recommend this because of the two assumptions presupposed: alphabets and white space. What will you do when you come across a language that doesn't use an alphabet, like Chinese? And what will you do when you come across a language that doesn't use white space in the same way as English, like Hungarian or Turkish? Or code for that matter, whitespace is critical to Python's syntax and is more than just a separator but actually has semantic meanings. This is one reason why multilingual models end up outperforming monolinguals on the same tasks in almost every case, they're forced to learn deeper representations for meaning without the bowling bumpers of easy tokenization. So let's look at some deeper methodologies that work for UTF-8 encoded languages.

Here are examples of all the current popular options for basing your tokenization:

1. Word-based - "Johannes Gutenberg" becomes ['Johannes', 'Gutenberg'].
2. Character-based - "Shakespeare" becomes ['S', 'h', 'a', 'k', 'e', 's', 'p', 'e', 'a', 'r', 'e'].
3. Subword-based - "The quick red Delphox jumped over the lazy brown Emolga" becomes ['the', 'quick', 'red', 'delph', 'ox', 'jump', 'ed', 'over', 'the', 'laz', 'y', 'brown']

Let's take a look at each of them in turn.

## **Word-based**

Word-based tokenizers most commonly just split on whitespace, but there are other methods like using regular expressions, dictionaries, or punctuations.

For example, punctuation based approach would split “It’s the truth!” into [“It”, “’”, “s”, “the”, “truth”, “!”], which is giving us slightly better context than splitting on whitespace alone. The TreebankWordTokenizer from NLTK is an example of a regular expression tokenizer. Word-based tokenizers are relatively easy to implement but require us to keep an unmanageably large dictionary mapping in order to encode every single possible word. That’s unreasonable, so generally, to make it work you’ll implement a dictionary cutoff and return unknown tokens when the model runs into unrecognized words. This makes the tokenizer poor at many tasks like code, name and entity recognition, as well as generalizing across domains.

## **Character-based**

Character based encoding methods are the most straightforward and easiest to implement since we just split on the UTF-8 character encodings. With this method we only need the tiniest of dictionaries to map characters to numbers which means we can prevent the need for unknown tokens and related concerns, however, it comes with major loss of information and fails to keep relevant syntax, semantics or morphology of the text.

## **Subword-based**

Just like Goldiloch’s and three bears, while character-based tokenizers are too hard and word-based tokenizers are too soft, subword-based tokenizers are just right. Subword-based tokenizers have proven to be the best option, being a mixture of the previous two. We are able to use a smaller dictionary like character-based tokenizer, but also lose less semantics like a word-based tokenizer. It even has the added bonus of adding some morphology information. However, it’s an unsolved problem for where and how words should be split and there are many different methods and approaches. The best method to choose will be, like all other things with LLMs, dependent on the task. If you don’t have a specific goal in mind for what you are trying to do, there will be consequences later.

There are three main algorithms to create the subword dictionaries: Byte Pair Encoding, WordPiece, and Unigram. In addition, SentencePiece is also very common and is a combination of the three that explicitly handles

whitespaces. It's outside the scope of this book to discuss how they work, but as a book focused on production you should know that the most popular subword tokenization methodologies are BPE and sentencepiece. In Listing 4.10, we'll go over how to actually train a custom version for both of these on your data so that you're equipped to face (almost) any dataset head on.

While reading the code, pay attention to where we train the tokenizers. In particular, there are three key parameters that you will want to tune: the `vocab_size`, `min_frequency`, and `special_tokens`. A larger vocabulary size means your tokenizer will be more robust and will likely be better at handling more languages, but will add computational complexity. Minimum frequency determines how often a particular subword token has to be seen in the dataset before it will be added to the dictionary. Larger values prevent rare and likely unimportant tokens from filling our dictionary, but also prevents us from learning rare tokens that are important. Lastly, special tokens are relatively straightforward and include syntactical tokens we care about specifically for model training.

#### **Listing 4.10 Train Your Own Subword Tokenizers**

```
import os
from pathlib import Path

import transformers
from tokenizers import ByteLevelBPETokenizer, SentencePieceBPETok
from tokenizers.processors import BertProcessing

# Initialize the txts to train from
paths = [str(x) for x in Path("./data/").glob("**/*.txt")]

# Train a Byte-Pair Encoding tokenizer
bpe_tokenizer = ByteLevelBPETokenizer()

bpe_tokenizer.train(
    files=paths,
    vocab_size=52_000,
    min_frequency=2,
    show_progress=True,
    special_tokens=[
        "<s>",
        "<pad>",
        "</s>",
```

```

        "<unk>",
        "<mask>",
    ],
)

token_dir = "./chapters/chapter_4/tokenizers/bytelevelbpe/"
if not os.path.exists(token_dir):
    os.makedirs(token_dir)
bpe_tokenizer.save_model(token_dir)

bpe_tokenizer = ByteLevelBPETokenizer(
    f"{token_dir}vocab.json",
    f"{token_dir}merges.txt",
)

example_text = "This sentence is getting encoded by a tokenizer."
print(bpe_tokenizer.encode(example_text).tokens)
# ['This', 'Ġsentence', 'Ġis', 'Ġgetting', 'Ġenc', 'Ġ',
#  'oded', 'Ġby', 'Ġa', 'Ġto', 'Ġken', 'Ġizer', 'Ġ.']
print(bpe_tokenizer.encode(example_text).ids)
# [2666, 5651, 342, 1875, 4650, 10010, 504, 265, \
#  285, 1507, 13035, 18]

bpe_tokenizer._tokenizer.post_processor = BertProcessing(
    ("</s>", bpe_tokenizer.token_to_id("</s>")),
    ("<s>", bpe_tokenizer.token_to_id("<s>")),
)
bpe_tokenizer.enable_truncation(max_length=512)

# Train a Sentencepiece Tokenizer
special_tokens = [
    "<s>",
    "<pad>",
    "</s>",
    "<unk>",
    "<cls>",
    "<sep>",
    "<mask>",
]
sentencepiece_tokenizer = SentencePieceBPETokenizer()

sentencepiece_tokenizer.train(
    files=paths,
    vocab_size=4000,
    min_frequency=2,
    show_progress=True,

```



```

        special_tokens=special_tokens,
    )

token_dir = "./chapters/chapter_4/tokenizers/sentencepiece/"
if not os.path.exists(token_dir):
    os.makedirs(token_dir)
sentencepiece_tokenizer.save_model(token_dir)

# convert
tokenizer = transformers.PreTrainedTokenizerFast(
    tokenizer_object=sentencepiece_tokenizer,
    model_max_length=512,
    special_tokens=special_tokens,
)
tokenizer.bos_token = "<s>"
tokenizer.bos_token_id = sentencepiece_tokenizer.token_to_id("<s>")
tokenizer.pad_token = "<pad>"
tokenizer.pad_token_id = sentencepiece_tokenizer.token_to_id("<pad>")
tokenizer.eos_token = "</s>"
tokenizer.eos_token_id = sentencepiece_tokenizer.token_to_id("</s>")
tokenizer.unk_token = "<unk>"
tokenizer.unk_token_id = sentencepiece_tokenizer.token_to_id("<unk>")
tokenizer.cls_token = "<cls>"
tokenizer.cls_token_id = sentencepiece_tokenizer.token_to_id("<cls>")
tokenizer.sep_token = "<sep>"
tokenizer.sep_token_id = sentencepiece_tokenizer.token_to_id("<sep>")
tokenizer.mask_token = "<mask>"
tokenizer.mask_token_id = sentencepiece_tokenizer.token_to_id("<mask>")
# and save for later!
tokenizer.save_pretrained(token_dir)

print(tokenizer.tokenize(example_text))
# ['_This', '_s', 'ent', 'ence', '_is', '_', 'g', 'et', 'tin', 'c',
#  'en', 'co', 'd', 'ed', '_', 'b', 'y', '_a', '_', 't', 'ok', 'e',
#  'iz', 'er', '.']
print(tokenizer.encode(example_text))
# [814, 1640, 609, 203, 1810, 623, 70, \
#  351, 148, 371, 125, 146, 2402, 959, 632]

```

Out of the two, BPE and Sentencepiece, we find ourselves using both about equally. It's mostly dependent upon which model we're finetuning or using as a base for a particular project. Algorithmically though, we're partial to Sentencepiece, not only because it tends to boost evaluation scores on pretty much any test for models trained on it, but also because it's closer to how we interact with morphology as people.

All in all though, tokenization loses information, just as converting from speech to text does, namely word order (syntax) and meaning (semantics). All of the information of what a number is and how it would differ from a letter is completely gone after tokenization. To circumvent potential semantic and syntactic issues, we need to create an approximation for each of these features and figure out how to mathematically represent them in abstraction to insert that meaning back into the tokenized vector. For this, we have embeddings.

## 4.4.2 Embeddings

Embeddings provide meaning to the vectors generated during tokenization. Tokenized text is just numbers assigned almost arbitrarily (occurrence based) to a dictionary, but it's at least in a format that the model can ingest. Embeddings are the next step, where positional and semantic encodings are created and looked up to give the model additional context for making decisions about how to (probably) complete the task it's given.

Embeddings are imperfect for several reasons, but perhaps the most relevant is the theoretical question: Can you represent a set using only a subset of that same set? In this case, the first set is language, one or more, and the second set is numbers, floats and digits. Math is a subset of language used to axiomatically describe things we accept as true. Take the English alphabet for example, can you represent the entire alphabet by only using some fraction of the 26 letters? Obviously not, but what if both the original set and the subset are both infinite? Can you represent all digits using only the decimals between 0 and 1? Given that the first is a numerable infinite set and the second is a non-enumerable infinite set, the answer is yes, which should be enheartening for the field of language modeling.

Now that we've talked about why embeddings shouldn't be completely and blindly relied upon, embeddings are what most businesses are actually looking for with LLMs. You don't need a 1.7 trillion parameter model to handle customers asking questions about your pricing or performing a search through your documents. Like we discussed in Chapter 2, embeddings have the innate advantage of being comparable by distance, provided both embeddings you're comparing were created by the same model in the same

dimensional space. That opens up the door for all sorts of speedy computation and retrieval where you never have to figure out how to host a gigantic model somewhere, because you can run a smaller embedding model on cpu and it takes milliseconds for hundreds of tokens.

One of the most popular and coolest applications of embeddings at the moment is Retrieval-Augmented Generation (RAG), where you store data that is pertinent to the overall task of the model and give portions of that data, as needed to a larger model at prompting time to improve results. If we apply this to the Boston Housing Dataset and attempt to predict the value of a new house, we can compare that house's embedded data to the closest comparable houses in the area and generate an informed appraisal, without ever needing an appraiser to verify, as long as the embeddings you're retrieving from are up-to-date.

Embeddings can be used for dozens of different tasks, and are the result of taking hidden state representations from your model. Every layer of your model is a potential option, but general consensus is taking representations after the final layer but before any sort of decoding. Listing 4.11 gives a practical example of how to actually extract the embeddings from both PyTorch and HuggingFace models. Best practice dictates that you should extract the embeddings from documents using whatever embedding model you are planning to use for inference, especially if those embeddings will end up being stored in a VectorDB later on. After creating our embeddings we show how to do a simple similarity search on the results which is the basis of RAG systems.

**Listing 4.11 Example Embeddings**

```
import numpy as np
from sentence_transformers import SentenceTransformer
from datasets import load_dataset

# Download embedding model and dataset
model_ckpt = "sentence-transformers/all-MiniLM-L6-v2"
model = SentenceTransformer(model_ckpt)

embs_train = load_dataset("tweet_eval", "emoji", split="train[:100000]")
embs_test = load_dataset("tweet_eval", "emoji", split="test[:100000]")
```

```

# Create embeddings
def embed_text(example):
    embedding = model.encode(example["text"])
    return {"embedding": np.array(embedding, dtype=np.float32)}

print(f"Train 1: {embs_train[0]}")
embs_train = embs_train.map(embed_text, batched=False)
embs_test = embs_test.map(embed_text, batched=False)

# Add Faiss index which allows similarity search
embs_train.add_faiss_index("embedding")

# Run Query
idx, knn = 1, 3 # Select the first query and 3 nearest neighbors

query = np.array(embs_test[idx]["embedding"], dtype=np.float32)
scores, samples = embs_train.get_nearest_examples("embedding", qu

# Print Results
print(f"QUERY LABEL: {embs_test[idx]['label']}")
print(f"QUERY TEXT: {embs_test[idx]['text'][:200]} [...] \n")
print("=" * 50)
print("Retrieved Documents:")
for score, label, text in zip(scores, samples["label"], samples["
    print("=" * 50)
    print(f"TEXT: \n{text[:200]} [...]")
    print(f"SCORE: {score:.2f}")
    print(f"LABEL: {label}")

```

Extracting embeddings, like the listing shows, is pretty simple, and differs very little from simply running inference or training on a dataset. Remember if you aren't using sentence transformers, set your model into eval mode, run with `torch.no_grad()`, and if you're running on torch 2.0+ run `torch.compile(model)` and things should speed up and become more computationally efficient immediately.

Another unsolved problem as-of-yet is how to compare embedding spaces. Mathematically sound comparisons have popped up time and again over the years, but, as has been demonstrated, mathematical soundness isn't the first problem to be solved, the modality is. If you're comparing language embeddings, a mathematically sound conversion of a linguistically sound comparison method is the solution, and a linguistically sound comparison is

dependent upon the goal of the comparison. It's too much to go into here, but we dive more deeply into this topic in Appendix C where we discuss diffusion and multi-modal LLMs.

## 4.5 Preparing a Slack Dataset

Now that we have learned the ins and outs of prepare the necessary assets to train our own LLM, we wanted to end this chapter by preparing our own dataset which we can use later on. For this exercise, we will tackle a very common problem in the industry. I'm sure most readers have experienced or witnessed an HR help channel that is constantly inundated with the same questions over and over. It doesn't matter how many FAQ pages are created, users don't want to waste their time searching documentation when they could just ask an expert. So let's build a chatbot to answer these questions!

We will show you how to pull your company's Slack data and prepare it for training an LLM based chatbot. In listing 4.12 we pull slack data and then filter it to keep just the user's data and save it to a parquet file. This way you can create a bot that will talk like you, but feel free to edit it. For example, you might enjoy creating a bot that talks like your boss, but I'd recommend not telling them in case they feel threatened knowing you are automating them out of a job.

### Listing 4.12 Example Pulling Slack Data

```
import slack_sdk
import pandas

token_slack = "Your Token Here"
client = slack_sdk.WebClient(token=token_slack)

dm_channels_response = client.conversations_list(types="im")

all_messages = {}

for channel in dm_channels_response["channels"]:
    history_response = client.conversations_history(channel=channel["id"])
    all_messages[channel["id"]] = history_response["messages"]

txts = []
```

```

for channel_id, messages in all_messages.items():
    for message in messages:
        try:
            text = message["text"]
            user = message["user"]
            timestamp = message["ts"]
            txts.append([timestamp, user, text])
        except:
            pass

slack_dataset = pandas.DataFrame(txts)
slack_dataset.columns = ["timestamp", "user", "text"]
self_user = (
    slack_dataset["user"].value_counts().idxmax()
) # Determine your internal user_id
df = slack_dataset[slack_dataset.user == self_user]

messages = df["text"].values.tolist()

messages.to_parquet("slack_dataset.gzip")

```

As you can see there's not much to it! We have an example dataset we pulled using this script in the github repo accompanying this book. We will use this dataset in the coming chapters.

We've gone over a lot in this chapter, but you should now be prepared and know how to select and evaluate a foundation model, prepare and clean a dataset, and optimize your own text processors. We will use this information in the next chapter to train and finetune our own LLM model.

## 4.6 Summary

- Data Engineers have unique datasets to acquire and manage for LLMs, like model weights, evaluation datasets, and embeddings.
- No matter your task, there is a wide array of open source models to choose from and acquire to use to finetune your own model.
- Text based tasks are harder to evaluate than just simple equality metrics you'd find in traditional ML tasks, but there are many industry benchmarks to help you get started.
- Evaluating LLMs for more than just performance like bias and potential

harm is your responsibility.

- You can use the Evaluate library to build your own evaluation metrics.
- There are many large open source datasets, but most of them come from scraping the web and require cleaning.
- Instruct schemas and annotating your data can be effective ways to clean and analyze your data.
- Finetuning a model on a dataset that has an appropriate distribution of speech acts for the task you want your model to perform will help your model generate context appropriate content.
- Building your own subword tokenizer to match your data can greatly improve your model's performance.
- Many problems teams are trying to use LLMs for can be solved by simply using embeddings from your model instead.

[1] C. Loizos, “‘GPT’ may be trademarked soon if OpenAI has its way,” TechCrunch, Apr. 25, 2023. <https://techcrunch.com/2023/04/24/gpt-may-be-trademarked-soon-if-openai-has-its-way/>

[2] You can learn more about the RAIL license here:  
<https://bigscience.huggingface.co/blog/the-bigscience-rail-license>

[3] N. Muennighoff et al., “Cross Lingual Generalization through Multitask Finetuning,” Nov. 03, 2022. <https://arxiv.org/abs/2211.01786>

[4] C. Xu et al., “WizardLM: Empowering Large Language Models to Follow Complex Instructions,” Jun. 10, 2023.  
<https://arxiv.org/abs/2304.12244>

[5] We won't be talking about it here, but we introduced Alpaca in the last chapter when discussing Knowledge Distillation.

[6] “Free Dolly: Introducing the World's First Truly Open Instruction-Tuned LLM,” Databricks, Apr. 12, 2023.  
<https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm>

[7] Entropy related metrics are highly related to Information Theory which we haven't covered in the published draft, however we recommend you take

a look at it if you're interested in creating or improving evaluation metrics for LLMs.

[8] D. Hendrycks et al., "Measuring Massive Multitask Language Understanding," arXiv (Cornell University), Sep. 2020, doi: <https://doi.org/10.48550/arxiv.2009.03300>

[9] You can learn more about CALM here <https://arxiv.org/abs/2308.12539v1> and winoqueer here <https://arxiv.org/abs/2306.15087>

[10] You can check out the repository here: <https://github.com/openai/evals>

[11] P. Rajpurkar, R. Jia, and P. Liang, "Know What You Don't Know: Unanswerable Questions for SQuAD," Jun. 2018 <https://arxiv.org/abs/1806.03822>

[12] These plots are created to mimic Spectral Density plots you might see in a physics lab, and if you aren't familiar with them this unfortunately isn't the book to go into it. But I do recommend you check out the WeightWatchers documentation if you are interested <https://github.com/CalculatedContent/WeightWatcher>.

[13] C. Zhou et al., "LIMA: Less Is More for Alignment," arXiv.org, May 18, 2023. <https://arxiv.org/abs/2305.11206>

[14] T. Dettmers et al., "SpQR: A Sparse-Quantized Representation for Near-Lossless LLM Weight Compression," arXiv.org, Jun. 05, 2023. <https://arxiv.org/abs/2306.03078>.

[15] L. Gao et al., "The Pile: An 800GB Dataset of Diverse Text for Language Modeling," Dec. 2020, <https://arxiv.org/abs/2101.00027>

[16] You can read the blog post introducing it here <https://together.ai/blog/redpajama>

[17] <https://github.com/yaodongC/awesome-instruction-dataset>

[18] LARP stands for Live Action Role Playing, and you can imagine the



tomfoolery of a customer pretending to be an elf, orc, or pirate and thus breaking all rules and expectations.

# 5 Training Large Language Models: How to generate the generator

## This chapter covers

- Setting up a training environment and common libraries
- Applying various training techniques including leveraging advanced methodologies
- Tips and tricks to get the most out of training

Are you ready to have some fun?! What do you mean the last four chapters weren't fun? Well, I promise this one for sure will be. We've leveled up a lot and gained a ton of context that will prove invaluable now as we start to get our hands dirty. By training an LLM, we can create bots that can do amazing things and have unique personalities. Indeed, we can create new friends and play with them. In fact, in the last chapter we showed you how to create a training dataset based on your slack messages and will show you how to take that dataset and create a persona of yourself. Finally, you will no longer have to talk to that one annoying coworker[\[1\]](#).

First things first, we'll show you how to set up a training environment as the process can be very resource-demanding, and without the proper equipment you won't be able to really enjoy what comes next. We'll then show you how to do the basics like training from scratch and finetuning, after which we'll get into some of the best-known methods to improve upon these processes making them more efficient, faster, and cheaper. Lastly, we'll end the chapter with some tips and tricks we've acquired through our experience of training models in the field.

## 5.1 Multi-GPU Environments

Training is a resource-intensive endeavor. A model that only takes a single GPU to run inference on, may take ten times that many to train if for nothing

else to parallelize your work and speed things up so you aren't waiting for a thousand years for it to finish training. In order to really take advantage of what we want to teach you in this chapter, we're going to first have to get you set up in an environment that you can use as a playground. Later in the chapter, we'll teach some resource-optimal strategies as well but you'll need to understand this if you want to use the largest LLMs anyway.

While there is a lot you can learn using smaller LLMs, what sets apart a pro from an amateur is often the ease and fluidity that they have working with larger models. And there's a lot of good reason for this since larger models outperform smaller models on the whole. If you want to work with the largest models, you'll never be able to get started on your laptop. Even most customized gaming rigs with dual GPUs aren't enough for inference let alone training.

To this end, we wanted to share with you a few methods to acquire access to a multi-GPU environment in the cloud and then will share the tools and libraries necessary to utilize them. The largest models do not fit in a single GPU so without these environments and tools you'll be stuck playing on easy mode forever.

### **5.1.1 Setting up**

It should be pointed out up front that, while multi-GPU environments are powerful, they are also expensive. When it comes to multi-GPUs, there are no services that offer a free tier or offering that we know of to recommend, but you can at least take comfort in knowing that paying per hour will be way cheaper than purchasing the rigs wholesale. Of course, if you can get your company to pay the bill I recommend it, but it is still your responsibility to spin down and turn off any environment you create to avoid unnecessary charges.

If your company is paying for it, they likely have chosen a hosted service that makes this whole process easy, but for the rest of us, I find setting up a VM in Google's Compute Engine one of the easiest methods. Once set up we will then show you how to utilize it.

### **A Note to the readers:**

for learning purposes, we use smaller models throughout this book in our code listings such that you can work with them on a single GPU either locally or using a service like Colab or Kaggle which offer a free tier of a single GPU. While the listings could be run on CPU-only hardware, you won't want to do it. Ultimately, there shouldn't be any need to run these costly VMs throughout the book. However, you likely will still want to. Training with multiple GPUs is much faster and more efficient, and often necessary. In addition, we do encourage you to try larger LLM variations that require these bigger rigs as the experience will be priceless. To make it easy, you should be able to recycle the code in this chapter for models and datasets much larger than what is presented, which will often just be a matter of changing a few lines.

## **Google Virtual Machine**

One of the easiest ways to create a multi-GPU environment is to set up a virtual machine (VM) on Google's cloud. To get started you'll need to create an account, create a Google Cloud Project, set up billing, and download the gcloud CLI. None of these steps are particularly hard, but be sure to follow the documentation found at <https://cloud.google.com/sdk/docs/install-sdk> for your operating system to install the SDK. The steps here also include the steps and how-to's for creating an account, project and billing in the "Before you begin" section if you don't already have an account.

For new accounts Google offers a \$300 credit to be used for pretty much anything on their GCP platform, except for GPUs. Sadly. I hate breaking this news, but there's just no free lunch where we are going. So you'll need to be sure to upgrade to a paid GCP tier. Don't worry, to just follow along it should only cost a couple dollars, but if you are money conscious I recommend reading the entire section first, then trying it out.

After setting up your account, by default GCP sets your GPU quotas to 0. Quotas are used as a way to manage your costs. To increase it, go here to find your quotas <https://console.cloud.google.com/iam-admin/quotas>. You'll be looking for the gpus\_all\_regions quota and since we plan to use multiple

GPUs, go ahead and submit a request to increase it to 2 or more.

With all the prerequisites in place, we'll go ahead and get started by initializing and logging in. You'll do this by running the following command in a terminal on your computer:

```
$ gcloud init
```

You may have already done this step if you had to install the SDK, but if not, this will launch a web browser to help us login and authorize us from the gcloud CLI and allows us to select our project. We will be assuming you just have the one project, but if this isn't your first rodeo and you have multiple projects, you'll need to add `--project` flag in all the subsequent commands.

Next, we need to determine two things, the machine type (or which GPUs we want to use) and our container image. To pick a machine type you can check out the different options at <https://cloud.google.com/compute/docs/gpus>. For beginners, I would highly recommend the NVIDIA L4 GPU as they are all around fantastic machines. For our purposes we'll be using the g2-standard-24 which comes with 2 L4 GPUs and costs us about \$2.00/hr. This machine type isn't in every region and zone, but you can figure out a region close to you here <https://cloud.google.com/compute/docs/regions-zones>. I will be using the us-west1 region and us-west1-a zone.

For the container image, we'll save ourselves a lot of hassle by using one that has all the basics set up. Generally, this means creating your own, but Google has several prebuilt for deep learning which are great to use, or a great place to start as a base image to customize. These are all found in the deeplearning-platform-release project that they own. To check out the options available you can run<sup>[2]</sup>:

```
$ gcloud container images list --project deeplearning-platform-re
```

You can pick from Base, Tensorflow, and Pytorch compiled images, along with the CUDA version and python version. We'll be using common-gpu-v20230925-debian-11-py310 which is a simple image ready for GPU with a Debian linux distribution and Python 3.10. Now that we have everything we need, we can create our VM! Go ahead and run the following commands to

set up the VM:

```
$ INSTANCE_NAME="g2-1lminprod-example"  
$ gcloud compute instances create ${INSTANCE_NAME} --zone=us-west
```

The first command just creates an environment variable to store the name of our VM since we'll be using it in several of the following commands as well. This can be whatever you want it to be. The next command is creating our VM instance. The first several flags (zone, image, machine) should make sense since we just spent the previous paragraphs preparing and gathering that information. The boot-disk-size is setting the disk space for our VM and actually defaults to 200GB so it's only included here since it's important to know for LLMs since they are large assets and you will likely need to increase it. Especially for LLMs that require multiple GPUs to run.

The scopes flag is passed to set authorization. Current GCP best practices recommend setting it to cloud-platform which determines authorization through OAuth and IAM roles. The metadata field isn't required but is used here as a trick to ensure the NVIDIA drivers are installed. This is really useful if you are using these commands to create a shell script to automate this process. You should know that this will cause a small delay between when the VM is up and when you can actually SSH into it as it won't be responsive while it installs the drivers. If you don't include it, the first time you SSH in through a terminal it will ask you if you want to install it so no harm done. However, if you access the VM through other methods (described in the next sections) you can run into issues. The last two commands are standard maintenance policies.

Once that runs you can verify the VM is up by running:

```
$ gcloud compute instances describe ${INSTANCE_NAME}
```

This command will give you a lot of information about your instance that is worth looking over, including a status field which should read 'RUNNING'. Once you've confirmed that it's up, we will SSH into it. If this is your first time using gcloud to SSH an SSH key will be generated automatically. Go ahead and run the following command:

```
$ gcloud compute ssh ${INSTANCE_NAME}
```

Your terminal will now be shelled into our multi-GPU VM and you are now in business. To this point, your VM is still just an empty shell though, so you'll want to bring in code. The easiest way to do this is likely to just copy the files over with Secure Copy Protocol (SCP). You can do this with a single file or a whole directory. For example, assuming your project has a requirements.txt file and a subdirectory local-app-folder, from a new terminal you can run the following commands:

```
$ gcloud compute scp requirements.txt ${INSTANCE_NAME}:/requirem
$ gcloud compute scp --recurse ~/local-app-folder/ ${INSTANCE_NAM
```

Overall, not too bad. Once you've gone through the process and set everything up, the next time you set up a VM it's only four commands (create, describe, ssh, scp) to get up and running.

Of course, these instances cost good money, so the last command you'll want before moving on is to know how to delete it. Which you can do with:

```
$ gcloud compute instances delete ${INSTANCE_NAME} --quiet
```

Of course, for Linux power users, this is likely all you need, but for the rest of us pleb's shelling into a VM through a terminal is less than an ideal working environment. We'll show you some tips and tricks to make the most of your remote machine.

## SSH through VSCode

For most devs a terminal is fine, but what we really want is an IDE. Most IDEs offer remote SSH capabilities but we'll demonstrate with VSCode. The first step is installing the extension Remote-SSH[\[3\]](#). There are other extensions that offer this capability you can check out, but Remote-SSH is maintained by Microsoft and has over 17 million installs so it's a great choice for beginners.

Next, we are going to run a configuration command:

```
$ gcloud compute config-ssh
```

Then, inside of VSCode, you can press F1 to open the command palette and

run the “Remote-SSH: Open SSH Host...” command and you should see your VM’s SSH address, which will look like `l4-llm-example.us-west1-a.project-id-401501`. If you don’t see it, something went wrong with the `config-ssh` command and you likely just need to run `gcloud init` again. Select the address and you should see a new VSCode window pop up. In the bottom corner, you’ll see that it is connecting to your remote machine. And you are done! Easy. From here you can use VSCode like you would when using it locally.

### **5.1.2 Libraries**

Though setting up hardware is important, none of it will work without the software packages that enable different points of hardware to communicate with each other effectively. With LLMs, the importance of this is compounded. I’ve personally had the experience where our hardware was all correctly configured. We were pretty sure our software setup was likewise configured, just to start up training a model and be met with an estimated training time of over 3 years. After troubleshooting, we realized this was because I had multiple versions of CUDA Toolkit installed and Pytorch was looking at an incompatible (up-to-date) one instead of the one I had intended to use.

These software packages are about more than just using the CUDA low-level communication with your GPU, they’re about load-balancing, quantizing, and parallelizing your data as it runs through each computation to make sure it’s going as fast as possible while still enabling a certain level of fidelity for the matrices. You wouldn’t want to spend a long time making sure your embedding vectors are phenomenal representations, just to have them distorted at runtime. Thus, we present to you, the four deep-learning libraries every practitioner should know for multi-GPU instances. At the time of writing, all complementary features between these libraries are experimental, so feel free to mix and match. If you get a setup that utilizes all four at once to their full potential without erroring, drop it in a reusable container so fast.

#### **DeepSpeed**

DeepSpeed is an optimization library for distributed deep learning.



DeepSpeed is powered by Microsoft and implements various enhancements for speed in training and inference, like handling extremely long or multiple inputs in different modalities, quantization, caching weights and inputs, and, probably the hottest topic right now, scaling up to thousands of GPUs.

Installation is fairly simple if you remember to always install the latest—but not nightly—version of PyTorch first. This means you also need to configure your CUDA Toolkit beforehand. Once you have that package, pip install deepspeed should get you right where you want to go unless, ironically, you use Microsoft's other products. If you are on a Windows OS there is only partial support and there are several more steps you will need to follow to get it working for inference, not training, mode.

## **Accelerate**

From Huggingface, Accelerate is made to help abstract away the code for parallelizing and scaling to multiple GPUs away from you so that you can focus on the training and inference side. One huge advantage of Accelerate is that it adds only 1 import, 2 lines of code, and changes 2 other lines compared to a standard training loop in PyTorch in its vanilla implementation. Beyond that, Accelerate also allows for fairly easy CLI usage allowing it to be automated along with Terraform or AWS CDK.

Accelerate boasts compatibility over most environments and as long as your environment is Python 3.8+ and PyTorch 1.10.0+ (CUDA compatibility first) you should be able to use Accelerate without problems. Once that's done, pip install accelerate should get you there. Accelerate also has experimental support for DeepSpeed if you would like to get the benefits of both.

## **BitsandBytes**

If you don't already know the name Tim Dettmers in this field, you should become acquainted pretty quickly. Not many people have done as much as he has to make CUDA-powered computing accessible. This package is made to help practitioners quantize models and perform efficient matrix multiplication for inference (and maybe training) within different bit sizes, all the way down to int8. BitsandBytes has both similar requirements and drawbacks to

DeepSpeed, the requirements being Python 3.8+ and CUDA 10.0+ on Linux and Mac Environments and partial support for Windows with a different package.

You should have little trouble installing BitsandBytes as `pip install bitsandbytes` should work for most use cases. If you find yourself on Windows, you're in luck: `pip install bitsandbytes-windows` will work as well, although if you want to use it with Huggingface's transformers or PyTorch you will need to edit some minimum requirements stated within both of those packages, as the windows version does not have the same version numbers as the regular package. BitsandBytes offers its own implementations of optimizers like Adam and NN layers like Linear to allow for that 8-bit boost to run deep learning apps on smaller devices at greater speed with a minimal drop in accuracy.

## xFormers

The most bleeding edge of the libraries we recommend for most use cases, xFormers is made for research and production. Following a (hopefully) familiar PyTorch-like pattern of independent building blocks for multiple modalities, xFormers takes it a step further and offers components that won't be available in PyTorch for quite a while. One that I've used quite a lot is memory-efficient exact attention, which speeds up inference considerably.

xFormers has more requirements than the other packages, and we'd like to stress once more that using one or more tools to keep track of your environment is strongly recommended. On Linux and Windows, you'll need PyTorch 2.0.1 and `pip install -U xFormers` should work for you. That said, there are paths for installation with pretty much any other version of PyTorch, but the main ones are 1.12.1, 1.13.1, and 2.0.1.

**Table 5.1 Comparison of optimization packages for ML**

| <u>Library</u> | <u>Faster<br/>Training<br/>or<br/>Inference</u> | <u>Code<br/>Integration</u> | <u>Lower<br/>Accuracy</u> | <u>Many<br/>GPUs</u> | <u>Quantization</u> | <u>Opti</u> |
|----------------|---|-----------------------------|---------------------------|----------------------|---------------------|-------------|
|                |   |                             |                           |                      |                     |             |

|              |          |               |         |     |                 |  |
|--------------|----------|---------------|---------|-----|-----------------|--|
| DeepSpeed    | Both     | CLI           | Depends | Yes | Supports        | Cac<br>grac<br>che<br>men<br>man<br>scal |
| Accelerate   | Both     | CLI &<br>Code | Depends | Yes | Supports        | Aut<br>com<br>para                       |
| BitsandBytes | Both     | Code          | Always  | NA  | Yes only        | Qua<br>qual<br>opti                      |
| xFormers     | Training | Code          | Depends | NA  | Yes and<br>more | Effi<br>atten<br>men<br>man              |

In Table 5.1, we can see a heavily reduced breakdown of what each of these packages does and how it integrates with your code. Each package does similar things, but even when performing the same task, they will often perform those tasks differently or on different parts of your model or pipeline. There is some overlap between packages, and we'd encourage you to use all of them to see how they might benefit you. Now that you have an environment and have a basic understanding of some of the tools we'll be using, let's move forward and see it in action.

## 5.2 Basic Training Techniques

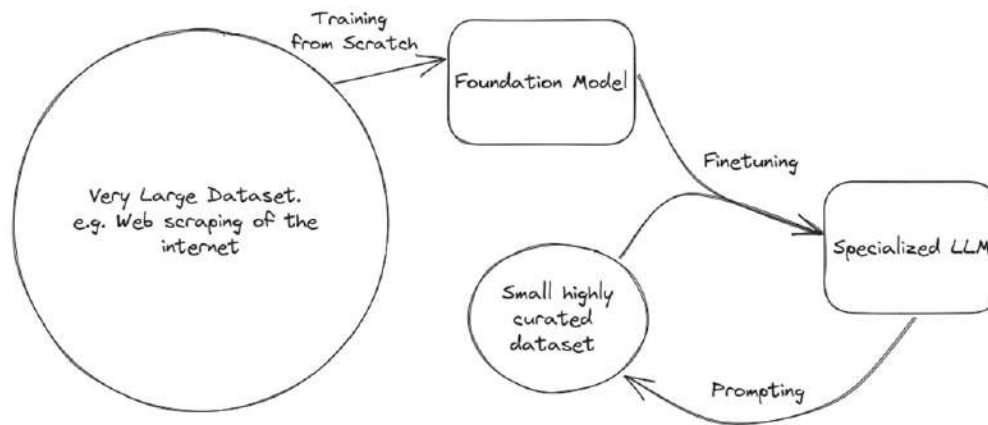
In training LLMs, the process typically starts with defining the architecture of the model, the nature and amount of data required, and the training objectives. We've already gone over these steps in the last chapter so you should be well prepared already, but as a brief recap: The model architecture usually follows a variant of the Transformer architecture due to its effectiveness in capturing long-term dependencies and its parallelizable

nature, making it amenable to large-scale computation. Data is the lifeblood of any LLM (or any ML model in general), and they typically require extensive corpora of diverse and representative text data. As the model's purpose is to learn to predict the next word in a sequence, it's crucial to ensure that the data covers a wide array of linguistic contexts.

Because we'll be going over various training techniques in this chapter, here's a (super) quick rundown of the investments you'll need for different types. For training from scratch, you'll need VRAM greater than 4x the number of billions of parameters to hold the model, along with the batches of training data. So to train a 1B parameter model from scratch, you'll need at least 5-6GB of VRAM depending on your batch sizes and context length. Consider training a 70B parameter model like Llama 2 as an exercise. How much VRAM will you need to fit the model, along with a 32k token context limit? If you're coming up with a number around 300GB of VRAM, you'd be right. For the finetuning techniques, you'll need significantly fewer resources for a couple of reasons, namely quantization and amount of data needed, meaning you no longer need 4x VRAM, but can use 2x or 1x with correct setup.

Unlike traditional ML models, LLMs are often trained in stages. In figure 5.1 we show the basic training lifecycle of an LLM, from scratch, finetuning, and then prompting. The first step is creating our foundation model, where we take a large, often unrefined, dataset and train an empty shell of a model on it. This will create a model that has seen such a large corpus of text that it appears to have a basic understanding of language. We can then take that foundation model and use transfer learning techniques, generally finetuning on a small highly curated dataset, to create a specialized LLM to do expert tasks. Lastly, we use prompting techniques, while not traditional training, allow us to goad the model to respond in a particular fashion or format improving the accuracy of our results.

**Figure 5.1 The training lifecycle of an LLM. We start by creating a foundation model based on a large corpus of text, which we then later finetune using a curated dataset for a specific task. We can then further improve the model by using the model itself and techniques like prompting to enhance or enlarge our curated dataset.**



You'll notice that the training lifecycle is often a continuous loop, training models to understand language better, then using those models to improve our training datasets. We'll be going more in depth into more advanced training techniques that take advantage of this loop like Prompt Tuning and RLHF later in the chapter but for now, let's solidify our understanding of three basic steps.

### 5.2.1 From Scratch

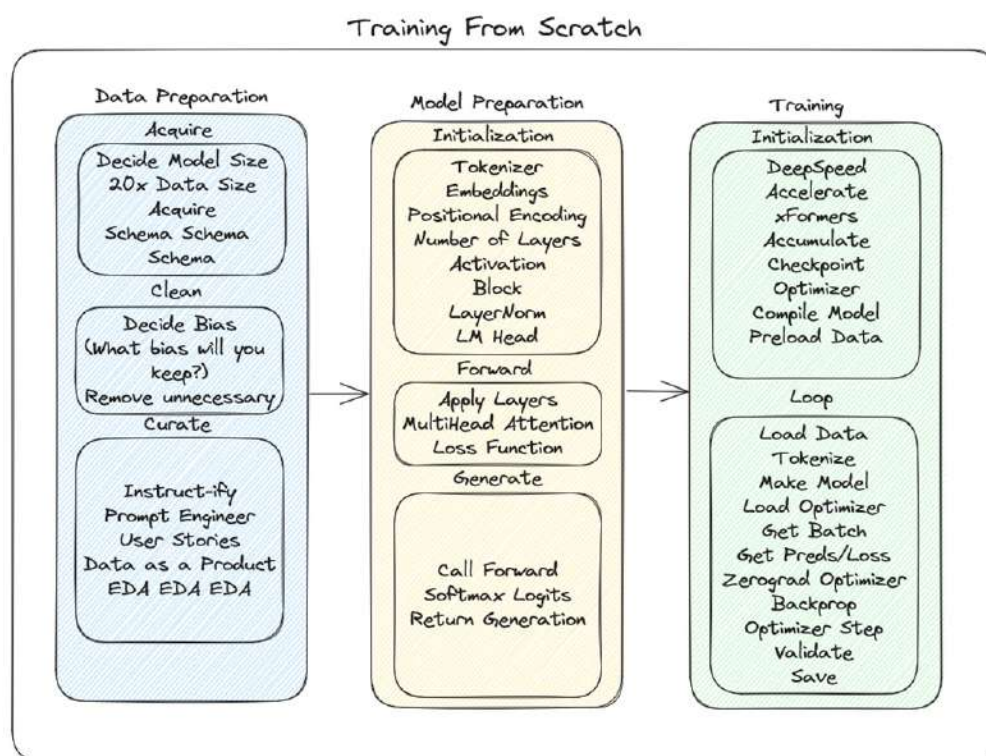
Training an LLM is computationally intensive and can take several weeks or months even on high-performance hardware. This process feeds chunks of data (or 'batches') to the model and adjusts the weights based on the calculated loss. Over time, this iterative process of prediction and adjustment, also known as an epoch, leads the model to improve its understanding of the syntactic structures and complexities in the data. It's worth noting that monitoring the training process is crucial to avoid overfitting, where the model becomes excessively tailored to the training data and performs poorly on unseen data. Techniques like early stopping, dropout, and learning rate scheduling are used to ensure the generalizability of the model, but are not silver bullets. Remember, the ultimate goal is not just to minimize the loss on training data but to create a model that can understand and generate human-like text across a broad range of contexts.

Training a LLM from scratch is a complex process that begins with defining the model's architecture. This decision should be guided by the specific task at hand, the size of the training dataset, and the available computational resources. The architecture, in simple terms, is a blueprint of the model that

describes the number and arrangement of layers, the type of layers (like attention or feed-forward layers), and the connections between them. Modern LLMs typically employ a variant of the Transformer architecture, known for its scalability and efficiency in handling long sequences of data.

Once the model's architecture is set, the next step is to compile a large and diverse dataset for training. The quality and variety of data fed into the model largely dictate the model's ability to understand and generate human-like text. A common approach is to use a large corpus of internet text, ensuring a wide-ranging mix of styles, topics, and structures. The data is then preprocessed and tokenized, converting the raw text into a numerical format that the model can learn from. During this tokenization process, the text is split into smaller units, or tokens, which could be as short as a single character or as long as a word.

**Figure 5.2** A simplified version of all of the steps going into training a language model (large or otherwise) from scratch. You must have data, then define all of the model behavior, and only then proceed to train.



With a model and dataset on the ready, the next step is to initialize the model

and set the learning objectives. The LLMs are trained using autoregressive semi-supervised learning techniques where the model learns to predict the next word in a sequence given the preceding words. The model's weights are randomly initialized and then adjusted through backpropagation and optimization techniques such as Adam or Stochastic Gradient Descent based on the difference between the model's predictions and the actual words in the training data. The aim is to minimize this difference, commonly referred to as the 'loss', to improve the model's predictive accuracy.

Training involves feeding the tokenized text into the model and adjusting the model's internal parameters to minimize the loss. We said this once but it bears repeating, this is a computationally demanding process that may take weeks or even months to complete depending on the model size and available hardware. After training, the model is evaluated on a separate validation dataset to ensure that it can generalize to unseen data. It is common to iterate on this process, finetuning the model parameters and adjusting the architecture as needed based on the model's performance on the validation set.

In listing 5.1, we explore training a brand new transformer-based language model “from scratch,” meaning without any previously defined architecture, embeddings, or weights. You shouldn’t have to train an LLM from scratch nor would you normally want to as it’s a very expensive and time consuming endeavor, however, knowing how can help you immensely. This exercise will allow you to run through the motions without training an actual massive model so feel free to explore with this code.

Some things to pay attention to when you review the listing, is you might recall we talked about tokenization and embeddings in the last chapter, so one thing to look out for is that for simplicity we will be using a character-based tokenizer. Before you run the code can you predict if this was a good or bad idea? Another is to pay attention to how we use both Accelerate and BitsandBytes which we introduced a little ago, you’ll see that these libraries come in mighty handy. Next, watch as we slowly build up the LLMs architecture, building each piece in a modular fashion and later defining how many of each piece is used and where to put them, almost like Legos. Finally, at the very end of the code you’ll see a typical model training loop, splitting

our data, running epochs in batches, and so forth.

#### **Listing 5.1 Example Training from Scratch**

```
import os
import torch
from accelerate import Accelerator

import bitsandbytes as bnb # Comment this out if running on Wind

# Define the overall GPT Architecture
class GPT(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.token_embedding = torch.nn.Embedding(vocab_size, n_e
        self.positional_embedding = torch.nn.Embedding(block_size
        self.blocks = torch.nn.Sequential(
            *[Block(n_embed, n_head=n_head) for _ in range(n_laye
        )
        self.ln_f = torch.nn.LayerNorm(n_embed)
        self.lm_head = torch.nn.Linear(n_embed, vocab_size)

        self.apply(self._init_weights)

    def forward(self, idx, targets=None):
        B, T = idx.shape

        tok_emb = self.token_embedding(idx)
        pos_emb = self.positional_embedding(torch.arange(T, devic
        x = tok_emb + pos_emb
        x = self.blocks(x)
        x = self.ln_f(x)
        logits = self.lm_head(x)

        if targets is None:
            loss = None
        else:
            B, T, C = logits.shape
            logits = logits.view(B * T, C)
            targets = targets.view(B * T)
            loss = torch.nn.functional.cross_entropy(logits, targ

        return logits, loss

    def _init_weights(self, module):
```



```

        if isinstance(module, torch.nn.Linear):
            torch.nn.init.normal_(module.weight, mean=0.0, std=0.1)
            if module.bias is not None:
                torch.nn.init.zeros_(module.bias)
        elif isinstance(module, torch.nn.Embedding):
            torch.nn.init.normal_(module.weight, mean=0.0, std=0.1)

    def generate(self, idx, max_new_tokens):
        for _ in range(max_new_tokens):
            idx_cond = idx[:, -block_size:]
            logits, loss = self(idx_cond)
            logits = logits[:, -1, :]
            probs = torch.nn.functional.softmax(logits, dim=-1)
            idx_next = torch.multinomial(probs, num_samples=1)
            idx = torch.cat((idx, idx_next), dim=1)
        return idx

# Define the building blocks of the model
class Block(torch.nn.Module):
    def __init__(self, n_embed, n_head):
        super().__init__()
        head_size = n_embed // n_head
        self.self_attention = MultiHeadAttention(n_head, head_size)
        self.feed_forward = FeedForward(n_embed)
        self.ln1 = torch.nn.LayerNorm(n_embed)
        self.ln2 = torch.nn.LayerNorm(n_embed)

    def forward(self, x):
        x = x + self.self_attention(self.ln1(x))
        x = x + self.feed_forward(self.ln2(x))
        return x

class MultiHeadAttention(torch.nn.Module):
    def __init__(self, num_heads, head_size):
        super().__init__()
        self.heads = torch.nn.ModuleList(
            [Head(head_size) for _ in range(num_heads)]
        )
        self.projection = torch.nn.Linear(head_size * num_heads, n_embed)
        self.dropout = torch.nn.Dropout(dropout)

    def forward(self, x):
        out = torch.cat([h(x) for h in self.heads], dim=-1)
        out = self.dropout(self.projection(out))
        return out

```

```

class Head(torch.nn.Module):
    def __init__(self, head_size):
        super().__init__()
        self.key = torch.nn.Linear(n_embed, head_size, bias=False)
        self.query = torch.nn.Linear(n_embed, head_size, bias=False)
        self.value = torch.nn.Linear(n_embed, head_size, bias=False)
        self.register_buffer(
            "tril", torch.tril(torch.ones(block_size, block_size))
        )

        self.dropout = torch.nn.Dropout(dropout)

    def forward(self, x):
        _, T, _ = x.shape
        k = self.key(x)
        q = self.query(x)
        attention = q @ k.transpose(-2, -1) * k.shape[-1] ** 0.5
        attention = attention.masked_fill(
            self.tril[:T, :T] == 0, float("-inf"))
        attention = torch.nn.functional.softmax(attention, dim=-1)
        attention = self.dropout(attention)

        v = self.value(x)
        out = attention @ v
        return out


class FeedFoward(torch.nn.Module):
    def __init__(self, n_embed):
        super().__init__()
        self.net = torch.nn.Sequential(
            torch.nn.Linear(n_embed, 4 * n_embed),
            torch.nn.ReLU(),
            torch.nn.Linear(4 * n_embed, n_embed),
            torch.nn.Dropout(dropout),
        )

    def forward(self, x):
        return self.net(x)


# Helper functions for training
def encode(string):
    return [utt2int[c] for c in string]

```

```

def decode(line):
    return "".join([int2utt[i] for i in line])

def get_batch(split):
    data = train_data if split == "train" else val_data
    idx = torch.randint(len(data) - block_size, (batch_size,))
    x = torch.stack([data[i : i + block_size] for i in idx])
    y = torch.stack([data[i + 1 : i + block_size + 1] for i in idx])
    x, y = x.to(device), y.to(device)
    return x, y

@torch.no_grad()
def estimate_loss():
    out = {}
    model.eval()
    for split in ["train", "val"]:
        losses = torch.zeros(eval_iters)
        for k in range(eval_iters):
            X, Y = get_batch(split)
            logits, loss = model(X, Y)
            losses[k] = loss.item()
        out[split] = losses.mean()
    model.train()
    return out

# Train the model
if __name__ == "__main__":
    # Parameters for our experiment
    batch_size = 64 # Number of utterances at once
    block_size = 256 # Maximum context window size
    max_iters = 5000
    eval_interval = 500
    learning_rate = 3e-4
    eval_iters = 200
    n_embed = 384
    n_head = 6
    n_layer = 6
    dropout = 0.2
    accelerator = Accelerator()
    device = accelerator.device
    doing_quantization = False # Change to True if imported bits

```

```

# Dataset
with open("./data/crimeandpunishment.txt", "r", encoding="utf
    text = f.read()

# Character-based pseudo-tokenization
chars = sorted(list(set(text)))
vocab_size = len(chars)
utt2int = {ch: i for i, ch in enumerate(chars)}
int2utt = {i: ch for i, ch in enumerate(chars)}

data = torch.tensor(encode(text), dtype=torch.long)
n = int(0.9 * len(data))
train_data = data[:n]
val_data = data[n:]

# Instantiate the model and look at the parameters
model = GPT().to(device)
print("Instantiated Model")
print(
    sum(param.numel() for param in model.parameters()) / 1e6,
    "Model parameters",
)

optimizer = (
    torch.optim.AdamW(model.parameters(), lr=learning_rate)
    if not doing_quantization
    else bnb.optim.Adam(model.parameters(), lr=learning_rate)
)
print("Instantiated Optimizer")

model, optimizer, train_data = accelerator.prepare(
    model, optimizer, train_data
)
print("Prepared model, optimizer, and data")

# Training block
for iter in range(max_iters):
    print(f"Running Epoch {iter}")
    if iter % eval_interval == 0 or iter == max_iters - 1:
        losses = estimate_loss()
        print(
            f"| step {iter}: train loss {losses['train']:.4f}
            "| validation loss {losses['val']:.4f} |"
        )

    xb, yb = get_batch("train")
    logits, loss = model(xb, yb)

```

```

optimizer.zero_grad(set_to_none=True)
accelerator.backward(loss)
optimizer.step()

# Create model directory
model_dir = "./models/scratchGPT/"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

# Save the model
model_path = model_dir + "model.pt"
torch.save(
    model.state_dict(),
    model_path,
)

# Load the saved model
loaded = GPT().load_state_dict(model_path)

# Test the loaded model
context = torch.zeros((1, 1), dtype=torch.long, device=device)
print(decode(loaded.generate(context, max_new_tokens=500))[0]).

```

In Listing 5.1, we explored how the actual Lego blocks are put together for the GPT family of models and showed a training loop reminiscent of our exploration of language modeling in chapter 2. In line with that, this example, beyond showing the first part of generative pre-training for models, also illustrates why character-based modeling, whether convolutional or otherwise, is weak for language modeling. Did you get it right? Yup, character based-modeling isn't the best. Alphabets on their own do not contain enough information to produce statistically significant results, regardless of the tuning amount. From a linguistic standpoint, this is obvious, as alphabets and orthography in general are representations of meaning generated from humans, which is not intrinsically captured.

Some of the ways to help with that information capture are increasing our tokenization capture window through word or subword or sentence-level tokenization, and also by completing the pre-training before showing the model our task in order to allow it to capture as much approximate representation as possible. Next we'll show what benefits combining these two steps can have on our model's performance.

## 5.2.2 Transfer Learning (Finetuning)

Transfer learning is an essential approach in machine learning and a cornerstone of training LLMs. It's predicated on the notion that we can reuse knowledge learned from one problem (the source domain) and apply it to a different but related problem (the target domain). In the context of LLMs, this typically means using a pre-trained model, trained on a large, diverse dataset, and adapting it to a more specific task or domain.

In the first step of transfer learning, an LLM is trained on a large, general-purpose corpus, such as the entirety of Wikipedia, books, or the internet. This pre-training stage allows the model to learn an extensive range of language patterns, nuances, and a wide variety of topics. The goal here is to learn a universal representation of language that captures a broad understanding of syntax, semantics, and world knowledge. These models are often trained for many iterations and require significant computational resources, which is why it's practical to use pre-trained models provided by organizations like OpenAI or Hugging Face.

After pre-training, the LLM is updated on a specific task or domain. This update process adapts the general-purpose language understanding of the model to a more specific task, such as sentiment analysis, text classification, or question answering. Updating usually requires significantly less computational resources than the initial pre-training phase, because it involves training on a much smaller dataset specific to the task at hand. Through this process, the model is able to apply the vast knowledge it gained during pre-training to a specific task, often outperforming models trained from scratch on the same task. This process of transfer learning has led to many of the advances in NLP over recent years.

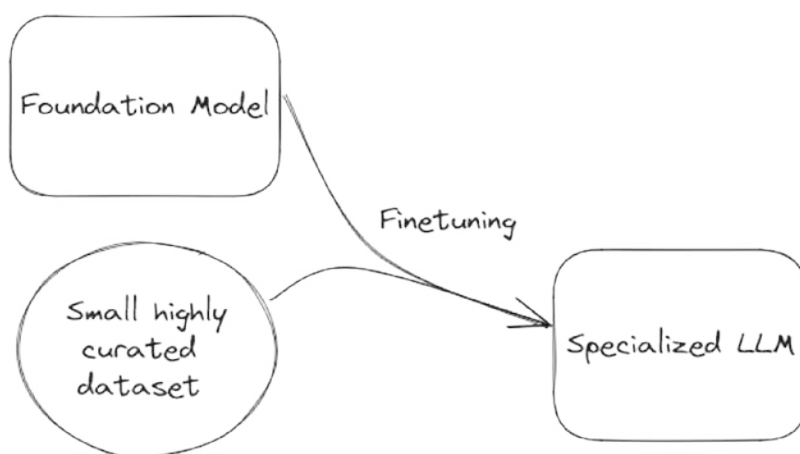
### **Finetuning**

There are several different transfer learning techniques, but when it comes to LLMs, the one everyone cares about is finetuning. Finetuning a LLM involves taking a pre-trained model - that is, a model that has already been trained on a large general corpus - and adapting it to perform a specific task or to understand a specific domain of data.

This technique leverages the fact that the base model has already learned a significant amount about the language, allowing you to reap the benefits of a large-scale model without the associated computational cost and time. The process of finetuning adapts the pre-existing knowledge of the model to a specific task or domain, making it more suitable for your specific use case. It's like having a generalist who already understands the language well and then providing specialist training for a particular job. This is often a more feasible approach for most users due to the significantly reduced computational requirements and training time, compared to training a model from scratch.

The first step in finetuning involves choosing a suitable pre-trained model. This decision is guided by the specific task you want the model to perform and by the resources available to you. Keep in mind this means setting a goal for the model's behavior before training. Once the pre-trained model has been chosen, it's crucial to prepare the specific dataset that you want the model to learn from. This could be a collection of medical texts, for example, if you're trying to finetune the model to understand medical language. The data must be preprocessed and tokenized in a way that's compatible with the pre-training of the model.

**Figure 5.3** Finetuning differs from training from scratch in that you don't have to define model behavior, you can use the exact same training loop, and you have a fraction of the data requirement.



The actual finetuning process involves training the model on your specific dataset, but with a twist: Instead of learning from scratch, the model's

existing knowledge is adjusted to better fit the new data. This is typically done with a smaller learning rate than in the initial training phase, to prevent the model from forgetting its previously learned knowledge. After finetuning, the model is evaluated on a separate dataset to ensure it can generalize to unseen data in the specific domain. Similar to training from scratch, this process may involve several iterations to optimize the model's performance. Finetuning offers a way to harness the power of LLMs for specific tasks or domains without the need for extensive resources or computation time.

In Listing 5.2 we will show how to finetune a GPT model. Notice how much less code there is in this listing than in 5.1. We don't need to define an architecture or a tokenizer, we'll just use those from the original model. We get to essentially skip ahead because weights and embeddings have already been defined.

#### **Listing 5.2 Example Finetuning**

```
import os
from transformers import (
    GPT2Tokenizer,
    GPT2LMHeadModel,
    GPT2Config,
    DataCollatorForLanguageModeling,
    TrainingArguments,
    Trainer,
)
from datasets import load_dataset

# Load and format the dataset
dataset = load_dataset("text", data_files="./data/crimeandpunishm
dataset = dataset.filter(lambda sentence: len(sentence["text"]) >
print(dataset["train"][0])

# Create model directory to save to
model_dir = "./models/betterGPT/"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

# Establish our GPT2 parameters (different from the paper and scr
config = GPT2Config(
    vocab_size=50261,
    n_positions=256,
    n_embd=768,
```



```

        activation_function="gelu",
    )

# Instantiate our tokenizer and our special tokens
tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
special_tokens_dict = {
    "bos_token": "<BOS>",
    "eos_token": "<EOS>",
    "pad_token": "<PAD>",
    "mask_token": "<MASK>",
}
tokenizer.add_special_tokens(special_tokens_dict)

# Instantiate our model from the config
model = GPT2LMHeadModel.from_pretrained(
    "gpt2", config=config, ignore_mismatched_sizes=True
)

# Create a tokenize function
def tokenize(batch):
    return tokenizer(
        str(batch), padding="max_length", truncation=True, max_le

# tokenize our whole dataset (so we never have to do it again)
tokenized_dataset = dataset.map(tokenize, batched=False)
print(f"Tokenized: {tokenized_dataset['train'][0]}")

# Create a data collator to format the data for training
data_collator = DataCollatorForLanguageModeling(
    tokenizer=tokenizer, mlm=True, mlm_probability=0.15
) # Masked Language Modeling - adds <MASK> tokens to guess the w

# Establish training arguments
train_args = TrainingArguments(
    output_dir=model_dir,
    num_train_epochs=1,
    per_device_train_batch_size=8,
    save_steps=5000,
    save_total_limit=2,
    report_to="none",
)

# Instantiate the Trainer
trainer = Trainer(

```

```

        model=model,
        args=train_args,
        data_collator=data_collator,
        train_dataset=tokenized_dataset["train"],
    )

# Train and save the model
trainer.train()
trainer.save_model(model_dir)
tokenizer.save_pretrained()

# Load the saved model
model = GPT2LMHeadModel.from_pretrained(model_dir)

# Test the saved model
input = "To be or not"
tokenized_inputs = tokenizer(input, return_tensors="pt")
out = model.generate(
    input_ids=tokenized_inputs["input_ids"],
    attention_mask=tokenized_inputs["attention_mask"],
    max_length=256,
    num_beams=5,
    temperature=0.7,
    top_k=50,
    top_p=0.90,
    no_repeat_ngram_size=2,
)
print(tokenizer.decode(out[0], skip_special_tokens=True))

```

Looking at listing 5.2 compared with listing 5.1, they have almost the exact same architecture (minus the activation function) and they're training on exactly the same data, but there's a marked improvement with the finetuned GPT2 model. This is because of the lack of learned representation in the first model. Our pretrained model along with subword BPE tokenization helps the model figure out which units of statistically-determined meaning are most likely to go together. You'll notice though, that GPT2, even with pre-training struggles to generate relevant longer narratives, despite using a newer better activation function.

## Finetuning OpenAI

We just trained a GPT model from scratch, and then we finetuned GPT2, but we know many readers really want the power behind OpenAI's larger GPT

models. Despite being proprietary models OpenAI has graciously created an API where we can finetune GPT3 models. Currently, three models are available for finetuning with OpenAI's platform, but it looks like they are intending to extend that finetuning ability to all of their models on offer. They have a whole guide written that you can find at [platform.openai.com](https://platform.openai.com), but once you have your dataset prepared in the format they need it in, the code is pretty easy. Here are some snippets for various tasks:

```
import os
import openai
openai.api_key = os.getenv("OPENAI_API_KEY")
openai.File.create(
    file=open("mydata.jsonl", "rb"),
    purpose='fine-tune'
)
```

This first snippet uploads a training dataset in the correct format for the platform and specifies the purpose as finetuning, but doesn't start the process yet. Next, you'll need to create the finetuning job:

```
openai.FineTuningJob.create(training_file="file-abc123", model="gpt-3.5-turbo")
```

This is where you specify which training file and which model you actually want to finetune. Once OpenAI's training loop has completed, you'll see the finetuned model's name populated when you retrieve the job details. Now you can use that model the same way you would have used any of the vanilla ones for chat completion or anything else like this:

```
completion = openai.ChatCompletion.create(
    model="ft:gpt-3.5-turbo:my-org:custom_suffix:id",
    messages=[
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": "Hello!"}
    ]
)
print(completion.choices[0].message)
```

And that's it for finetuning an OpenAI model! Very simple, doesn't take too long, and as of March of 2023 your data is private to you. Of course, you'll be ceding all of the control of how that finetuning occurs over to OpenAI. If you'd like to do something beyond vanilla finetuning, you'll need to do that

yourself, and in just a minute we'll be going over those techniques you may consider, along with some more advanced processes that can help with more fine-grained models and more complex tasks.

### **5.2.3 Prompting**

One of the main reasons why LLMs are so powerful compared to traditional ML is because we can train them at run time. Give them a set of instructions and watch them follow them to the best of their ability. This technique is called prompting and is used in LLMs to guide the model's output. In essence, the prompt is the initial input given to the model that provides it with context or instructions for what it should do. For example, asking a model to "Translate the following English text to French" or "Summarize the following article" are prompts. In the context of LLMs, prompting becomes even more critical, as these models are not explicitly programmed to perform specific tasks but learn to respond to a variety of tasks based on the given prompt.

Prompt engineering refers to the process of crafting effective prompts to guide the model's behavior. The aim is to create prompts that lead the model to provide the most desirable or useful output. This can be more complex than it appears, as slight changes in how a prompt is phrased can lead to vastly different responses from the model. Some strategies for prompt engineering include being more explicit in the prompt, providing an example of the desired output, or rephrasing the prompt in different ways to get the best results. It's a mixture of both art and science, requiring a good understanding of the model's capabilities and limitations.

In this chapter we are going to focus mainly on training and finetuning which are steps we need to do before deployment, but we would have been reminiscence if we didn't first discuss it here. We will be talking about prompting much more in depth in chapter 7.

## **5.3 Advanced Training Techniques**

Now that you know how to do the basics, we'll be going over some more advanced techniques. These techniques have been developed for a variety of reasons. To improve generated text outputs, shrink the model, provide

continuous learning, speed up training, or reduce costs. Depending on the needs of your organization you may need to reach for a different training solution. While not a comprehensive list, these are the techniques we see used most often and should be a valuable set of tools for you as you prepare a production ready model.

### **Classical ML Training Background**

Going over some techniques to enhance your finetuning process requires a bit of background. We won't be doing a full course in classical ML, however in case this is someone's first exposure, there are some classic learning paradigms that experiments tend to follow, those being Supervised, Unsupervised, Adversarial, and Reinforcement.

- Supervised involves collecting both the data to train on, as well as labels showcasing the expected output.
- Unsupervised does not require labels, as the data is probed for similarity and grouped into clusters that are the closest comparison to each other.
- Adversarial is what's used to train a Generative Adversarial Network, and it involves what are generally referred to as a Critic model and a Forger model. These two models essentially play a game against each other where the forger tries to copy some ideal output and the critic tries to determine whether the forgery is the real thing.
- Reinforcement opts for establishing a reward function instead of having predefined labels for the model to learn from, and by measuring the model's actions it is given a reward based on that function instead.

All LLMs must be trained using at least one of these, and they perform at a high level with all of them done correctly. The following training techniques differ from those basic ones and range from adding some form of human input to the model to compare outputs to changing the way the model does matrix multiplication.

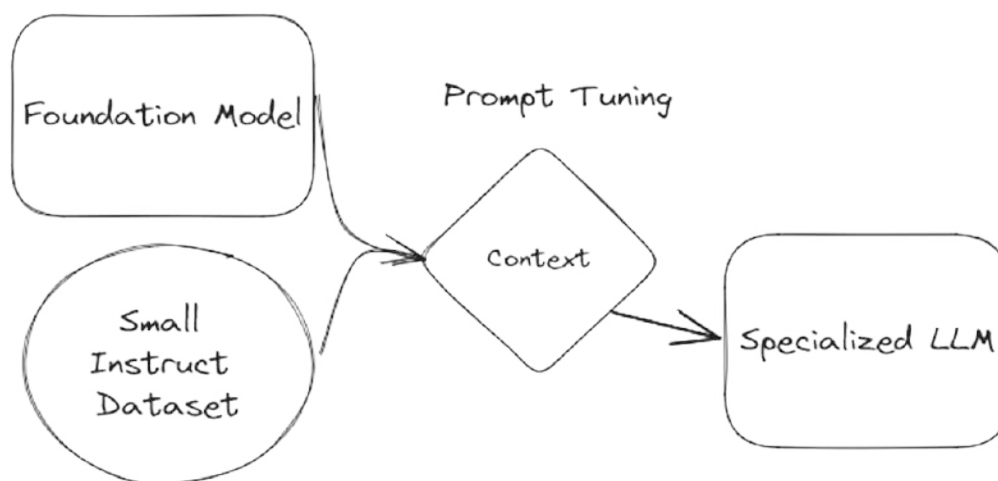
### **5.3.1 Prompt Tuning**

We've gone over pragmatics before, but as a reminder, language models perform better when given real-world non-semantic context pertaining to

tasks and expectations. Something that can be seen from beginning to end of language modeling techniques is an underlying assumption that the LM, given inputs and expected outputs, would be able to not only divine the task to be done, but the best way to do that task within the number of parameters specified.

While the idea of the model inferring both the task and the method of completing it from the data showed promise, it has been shown time and time again from BERT to every T5 model and now to all LLMs that providing your model with the expected task, as well as relevant information for solving the task improves model performance drastically. Google Research, DeepMind, and OpenAI have all published papers about Prompt Tuning, or giving a model pragmatic context during training, as early as 2021. The benefits are reducing the amount of data required for the model to converge during training, and, even cooler, the ability to reuse a completely frozen language model for new tasks without retraining or fully finetuning.

**Figure 5.4 Prompt Tuning foregoes most of finetuning to allow the majority of the foundation model's language understanding ability to stay exactly the same, and instead focuses on changing how the model responds to specific inputs**



Because LLMs are so large (and getting larger), it is becoming increasingly more difficult to share them, and even more difficult to guarantee their performance on a given task, even one that it trained on. Prompt Tuning can help nudge the model in the right direction without becoming a significant cost.

In Listing 5.3 we'll show how to prompt tune a smaller variant of the BLOOMZ model from Big Science. BLOOMZ was released as an early competitor in the LLM space, but has ultimately struggled to garner attention or momentum in the community because of its inability to generate preferred outputs despite its mathematical soundness. Because Prompt Tuning doesn't add much to the regular finetuning structure that we used in Listings 5.2, we'll perform Parameter Efficient FineTuning (PEFT) which reduces the memory requirements drastically by determining which model parameters need changing the most.

#### **Listing 5.3 Example Prompt Tuning**

```
import os
from transformers import (
    AutoModelForCausalLM,
    AutoTokenizer,
    default_data_collator,
    get_linear_schedule_with_warmup,
)
from peft import (
    get_peft_model,
    PromptTuningInit,
    PromptTuningConfig,
    TaskType,
)
import torch
from datasets import load_dataset
from torch.utils.data import DataLoader
from tqdm import tqdm

# Helper function to preprocess text - go ahead and skip to the t
def preprocess_function(examples):
    batch_size = len(examples[text_column])
    inputs = [
        f"{text_column} : {x} Label : " for x in examples[text_co
    ]
    targets = [str(x) for x in examples[label_column]]
    model_inputs = tokenizer(inputs)
    labels = tokenizer(targets)

    for i in range(batch_size):
        sample_input_ids = model_inputs["input_ids"][i]
        label_input_ids = labels["input_ids"][i] + [tokenizer.pad
```

```

        model_inputs["input_ids"][i] = sample_input_ids + label_i
        labels["input_ids"][i] = [-100] * len(
            sample_input_ids
        ) + label_input_ids
        model_inputs["attention_mask"][i] = [1] * len(
            model_inputs["input_ids"][i]
        )

    for i in range(batch_size):
        sample_input_ids = model_inputs["input_ids"][i]
        label_input_ids = labels["input_ids"][i]
        model_inputs["input_ids"][i] = [tokenizer.pad_token_id] *
            max_length - len(sample_input_ids)
        ) + sample_input_ids
        model_inputs["attention_mask"][i] = [0] * (
            max_length - len(sample_input_ids)
        ) + model_inputs["attention_mask"][i]
        labels["input_ids"][i] = [-100] * (
            max_length - len(sample_input_ids)
        ) + label_input_ids
        model_inputs["input_ids"][i] = torch.tensor(
            model_inputs["input_ids"][i][:max_length]
        )
        model_inputs["attention_mask"][i] = torch.tensor(
            model_inputs["attention_mask"][i][:max_length]
        )
        labels["input_ids"][i] = torch.tensor(
            labels["input_ids"][i][:max_length]
        )

    model_inputs["labels"] = labels["input_ids"]
    return model_inputs

```

```

# Model Prompt Tuning
if __name__ == "__main__":
    # Define training parameters
    device = "cuda"
    model_name_or_path = "bigscience/bloomz-560m"
    tokenizer_name_or_path = "bigscience/bloomz-560m"
    dataset_name = "twitter_complaints"
    text_column = "Tweet text"
    label_column = "text_label"
    max_length = 64
    lr = 3e-2
    num_epochs = 1
    batch_size = 8

```



```

# Define Prompt Tuning Config, notice init_text
peft_config = PromptTuningConfig(
    task_type=TaskType.CAUSAL_LM,
    prompt_tuning_init=PromptTuningInit.TEXT,
    num_virtual_tokens=8,
    prompt_tuning_init_text="Classify if the tweet "
    "is a complaint or not:",
    tokenizer_name_or_path=model_name_or_path,
)
checkpoint_name = (
    f"{dataset_name}_{model_name_or_path}"
    f"_{peft_config.peft_type}_{peft_config.task_type}_v1.pt"
    "/", "_"
)

# Load Dataset
dataset = load_dataset("ought/raft", dataset_name)
print(f"Dataset 1: {dataset['train'][0]}")

# Label Dataset
classes = [
    label.replace("_", " ")
    for label in dataset["train"].features["Label"].names
]
dataset = dataset.map(
    lambda x: {"text_label": [classes[label] for label in x["
    batched=True,
    num_proc=1,
)
print(f"Dataset 2: {dataset['train'][0]}")

# Load Tokenizer
tokenizer = AutoTokenizer.from_pretrained(model_name_or_path)
if tokenizer.pad_token_id is None:
    tokenizer.pad_token_id = tokenizer.eos_token_id
target_max_length = max(
    [
        len(tokenizer(class_label)["input_ids"])
        for class_label in classes
    ]
)
print(f"Target Max Length: {target_max_length}")

# Run Tokenizer across dataset and preprocess
processed_datasets = dataset.map(

```

```

        preprocess_function,
        batched=True,
        num_proc=1,
        remove_columns=dataset["train"].column_names,
        load_from_cache_file=False,
        desc="Running tokenizer on dataset",
    )

# Prepare Data Loaders
train_dataset = processed_datasets["train"]
eval_dataset = processed_datasets["test"]

train_dataloader = DataLoader(
    train_dataset,
    shuffle=True,
    collate_fn=default_data_collator,
    batch_size=batch_size,
    pin_memory=True,
)
eval_dataloader = DataLoader(
    eval_dataset,
    collate_fn=default_data_collator,
    batch_size=batch_size,
    pin_memory=True,
)

# Load Foundation Model
model = AutoModelForCausalLM.from_pretrained(model_name_or_pa
model = get_peft_model(model, peft_config)
print(model.print_trainable_parameters())
model = model.to(device)

# Define Optimizer
optimizer = torch.optim.AdamW(model.parameters(), lr=lr)
lr_scheduler = get_linear_schedule_with_warmup(
    optimizer=optimizer,
    num_warmup_steps=0,
    num_training_steps=(len(train_dataloader) * num_epochs),
)

# Training Steps
for epoch in range(num_epochs):
    model.train()
    total_loss = 0
    for step, batch in enumerate(tqdm(train_dataloader)):
        batch = {k: v.to(device) for k, v in batch.items()}
        outputs = model(**batch)

```

```

        loss = outputs.loss
        total_loss += loss.detach().float()
        loss.backward()
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()

    model.eval()
    eval_loss = 0
    eval_preds = []
    for step, batch in enumerate(tqdm(eval_dataloader)):
        batch = {k: v.to(device) for k, v in batch.items()}
        with torch.no_grad():
            outputs = model(**batch)
            loss = outputs.loss
            eval_loss += loss.detach().float()
            eval_preds.extend(
                tokenizer.batch_decode(
                    torch.argmax(outputs.logits, -1).detach().cpu()
                    skip_special_tokens=True,
                )
            )

    eval_epoch_loss = eval_loss / len(eval_dataloader)
    eval_ppl = torch.exp(eval_epoch_loss)
    train_epoch_loss = total_loss / len(train_dataloader)
    train_ppl = torch.exp(train_epoch_loss)
    print(
        f"{epoch=}: {train_ppl=} {train_epoch_loss=} "
        f"{eval_ppl=} {eval_epoch_loss=}"
    )

# Create model directory to save to
model_dir = "./models/PromptTunedPEFT"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

# Saving
tokenizer.save_pretrained(model_dir)
model.save_pretrained(model_dir)

# Inference
with torch.no_grad():
    inputs = tokenizer(
        f'{text_column} : {"@nationalgridus I have no water '
        'the bill is current and paid. Can you do something a '
        'this?"}' Label : ',

```

```

        return_tensors="pt",
    )

    inputs = {k: v.to(device) for k, v in inputs.items()}
    outputs = model.generate(
        input_ids=inputs["input_ids"],
        attention_mask=inputs["attention_mask"],
        max_new_tokens=10,
        eos_token_id=3,
    )
    print(
        tokenizer.batch_decode(
            outputs.detach().cpu().numpy(), skip_special_tokens=True
        )
    )

```

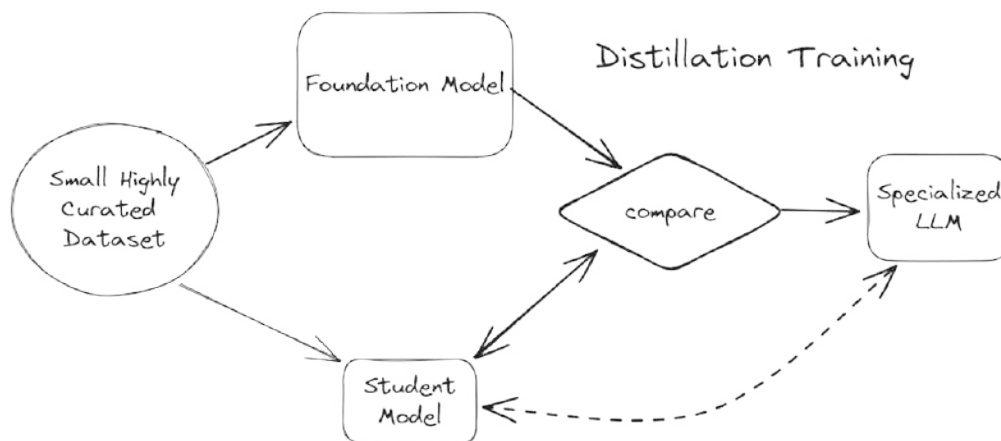
Other than the changed setup, the main difference between listings 5.3 and 5.2 is simply prepending a Prompt with some sort of instruction to the beginning of each input, reminiscent of the T5 training method that pioneered having some sort of a prepended task string before every input. Prompt tuning has emerged as a powerful technique for fine-tuning large language models to specific tasks and domains. By tailoring prompts to the desired output and optimizing them for improved performance, we can make our models more versatile and effective. However, as our LLMs continue to grow in scale and complexity, it becomes increasingly challenging to efficiently fine-tune them on specific tasks. This is where knowledge distillation comes into play, offering a logical next step. Knowledge distillation allows us to transfer the knowledge and expertise of these highly tuned models to smaller, more practical versions, enabling a wider range of applications and deployment scenarios. Together, prompt tuning and knowledge distillation form a dynamic duo in the arsenal of techniques for harnessing the full potential of modern LLMs.

### 5.3.2 Finetuning with knowledge distillation

Knowledge distillation is an advanced technique that provides a more efficient path to finetuning a LLM. Rather than just finetuning an LLM directly, knowledge distillation involves transferring the "knowledge" from a large, complex model (the teacher) to a smaller, simpler model (the student). The aim is to create a more compact model that retains the performance

characteristics of the larger model but is more efficient in terms of resource usage.

**Figure 5.5 Knowledge Distillation** allows a smaller model to learn from a foundation model to replicate similar behavior with fewer parameters. The student model does not always learn the emergent qualities of the foundation model, so the dataset must be especially curated. The dotted line is indicating a special relationship as the Student Model becomes the Specialized LLM.



The first step in knowledge distillation is to select a pre-trained LLM as the teacher model. This could be any of the large models such as Llama 2 70B or Falcon 180B, which have been trained on vast amounts of data. You also need to create or select a smaller model as the student. The student model might have a similar architecture to the teacher, but with fewer layers or reduced dimensionality in order to make it smaller and faster.

Next, the student model is trained on the same task as the teacher model. However, instead of learning from the raw data directly, the student model learns to mimic the teacher model's outputs. This is typically done by adding a term to the loss function that encourages the student model's predictions to be similar to the teacher model's predictions. This means that the student model not only learns from the task-specific labels, but also benefits from the rich representations learned by the teacher model.

Once the distillation process is complete, you'll have a compact student model that can handle the specific tasks learned from the teacher model, but at a fraction of the size and computational cost. The distilled model can then be further finetuned on a specific task or dataset if required. Through knowledge distillation, you can leverage the power of LLMs in situations

where computational resources or response time are limited.

In listing 5.4, we'll show how to perform finetuning with knowledge distillation using BERT and becoming DistilBERT. As opposed to regular finetuning, pay attention to the size and performance of the model. Both will drop, however size will drop much faster than performance.

#### **Listing 5.4 Example Knowledge Distillation**

```
import os
from transformers import (
    AutoTokenizer,
    TrainingArguments,
    Trainer,
    AutoModelForSequenceClassification,
    DataCollatorWithPadding,
)
from datasets import load_dataset, load_metric

import torch
import torch.nn as nn
import torch.nn.functional as F
import numpy as np

def process(examples):
    tokenized_inputs = tokenizer(
        examples["sentence"], truncation=True, max_length=256
    )
    return tokenized_inputs

def compute_metrics(eval_pred):
    predictions, labels = eval_pred
    predictions = np.argmax(predictions, axis=1)
    acc = accuracy_metric.compute(
        predictions=predictions, references=labels
    )
    return {
        "accuracy": acc["accuracy"],
    }

class DistillationTrainingArguments(TrainingArguments):
    def __init__(self, *args, alpha=0.5, temperature=2.0, **kwargs
```

```

    super().__init__(*args, **kwargs)

    self.alpha = alpha
    self.temperature = temperature

class DistillationTrainer(Trainer):
    def __init__(self, *args, teacher_model=None, **kwargs):
        super().__init__(*args, **kwargs)
        self.teacher = teacher_model
        # place teacher on same device as student
        self._move_model_to_device(self.teacher, self.model.device)
        self.teacher.eval()

    def compute_loss(self, model, inputs, return_outputs=False):
        # compute student output
        outputs_student = model(**inputs)
        student_loss = outputs_student.loss
        # compute teacher output
        with torch.no_grad():
            outputs_teacher = self.teacher(**inputs)

        # assert size
        assert (
            outputs_student.logits.size() == outputs_teacher.logits.size()
        )

        # Soften probabilities and compute distillation loss
        loss_function = nn.KLDivLoss(reduction="batchmean")
        loss_logits = loss_function(
            F.log_softmax(
                outputs_student.logits / self.args.temperature, dim=-1
            ),
            F.softmax(
                outputs_teacher.logits / self.args.temperature, dim=-1
            ),
        ) * (self.args.temperature**2)
        # Return weighted student loss
        loss = (
            self.args.alpha * student_loss
            + (1.0 - self.args.alpha) * loss_logits
        )
        return (loss, outputs_student) if return_outputs else loss

if __name__ == "__main__":
    # Create model directory to save to

```

```

model_dir = "./models/KDGPT/"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

# Define Teacher and Student models
student_id = "gpt2"
teacher_id = "gpt2-medium"

teacher_tokenizer = AutoTokenizer.from_pretrained(teacher_id)
student_tokenizer = AutoTokenizer.from_pretrained(student_id)

sample = "Here's our sanity check."

assert teacher_tokenizer(sample) == student_tokenizer(sample)
    "Tokenizers need to have the same output! "
    f"{teacher_tokenizer(sample)} != {student_tokenizer(sample)}
)

del teacher_tokenizer
del student_tokenizer

tokenizer = AutoTokenizer.from_pretrained(teacher_id)
tokenizer.add_special_tokens({"pad_token": "[PAD]"})

dataset_id = "glue"
dataset_config = "sst2"

dataset = load_dataset(dataset_id, dataset_config)

tokenized_dataset = dataset.map(process, batched=True)
tokenized_dataset = tokenized_dataset.rename_column("label",

print(tokenized_dataset["test"].features)

# create label2id, id2label dicts for nice outputs for the mo
labels = tokenized_dataset["train"].features["labels"].names
num_labels = len(labels)
label2id, id2label = dict(), dict()
for i, label in enumerate(labels):
    label2id[label] = str(i)
    id2label[str(i)] = label

# define training args
training_args = DistillationTrainingArguments(
    output_dir=model_dir,
    num_train_epochs=1,
    per_device_train_batch_size=1,

```



```

        per_device_eval_batch_size=1,
        fp16=True,
        learning_rate=6e-5,
        seed=8855,
        # Evaluation strategies
        evaluation_strategy="epoch",
        save_strategy="epoch",
        save_total_limit=2,
        load_best_model_at_end=True,
        metric_for_best_model="accuracy",
        report_to="none",
        # push to hub parameters
        push_to_hub=False,
        # distillation parameters
        alpha=0.5,
        temperature=4.0,
    )

    # define data_collator
    data_collator = DataCollatorWithPadding(tokenizer=tokenizer)

    # define model
    teacher_model = AutoModelForSequenceClassification.from_pretr
        teacher_id,
        num_labels=num_labels,
        id2label=id2label,
        label2id=label2id,
    )

    # define student model
    student_model = AutoModelForSequenceClassification.from_pretr
        student_id,
        num_labels=num_labels,
        id2label=id2label,
        label2id=label2id,
    )

    # define metrics and metrics function
    accuracy_metric = load_metric("accuracy")

    trainer = DistillationTrainer(
        student_model,
        training_args,
        teacher_model=teacher_model,
        train_dataset=tokenized_dataset["train"],
        eval_dataset=tokenized_dataset["validation"],
        data_collator=data_collator,

```

```
        tokenizer=tokenizer,  
        compute_metrics=compute_metrics,  
    )  
    trainer.train()  
  
    trainer.save_model(model_dir)
```

Knowledge distillation, as exemplified by the provided `compute_loss` method, is a technique that enables the transfer of valuable insights from a teacher model to a more lightweight student model. In this process, the teacher model provides soft targets, offering probability distributions over possible outputs, which are then utilized to train the student model. The critical aspect of knowledge distillation lies in the alignment of these distributions, ensuring that the student model not only learns to mimic the teacher's predictions but also gains a deeper understanding of the underlying data. This approach helps improve the student's generalization capabilities and performance on various tasks, ultimately making it more efficient and adaptable.

As we look forward, one logical progression beyond knowledge distillation is the incorporation of reinforcement learning with human feedback. While knowledge distillation enhances a model's ability to make predictions based on existing data, reinforcement learning with human feedback allows the model to learn directly from user interactions and feedback. This dynamic combination not only refines the model's performance further but also enables it to adapt and improve continuously. By incorporating human feedback, reinforcement learning can help the model adapt to real-world scenarios, evolving its decision-making processes based on ongoing input, making it an exciting and natural evolution in the development of LLM systems.

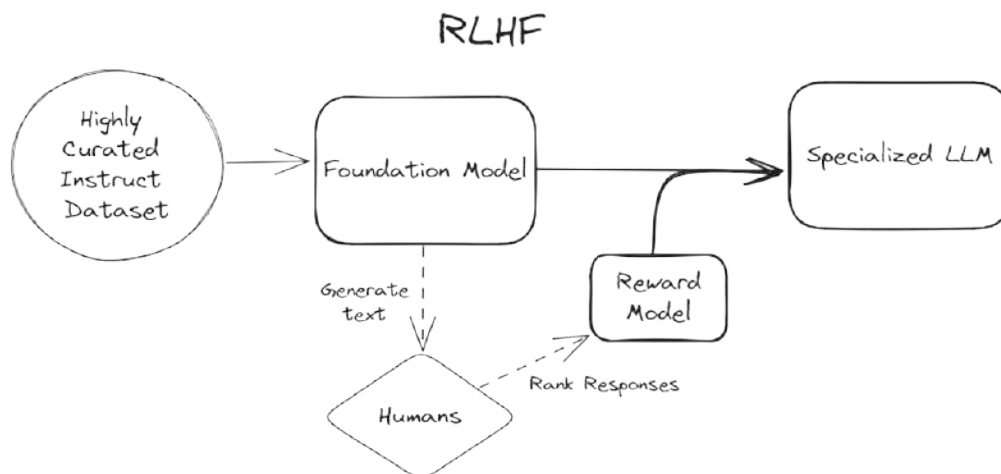
### 5.3.3 Reinforcement Learning with Human Feedback

Reinforcement Learning with Human Feedback (RLHF) is a newer training technique developed to overcome one of the biggest challenges when it comes to Reinforcement Learning (RL): how to create reward systems that actually work. It sounds easy, but anyone who's played around with RL knows how difficult it can be. Before AlphaStar<sup>[4]</sup> I was building my own RL bot to play StarCraft, a war simulation game in space. A simple reward

system based on winning or losing was taking too long, so I decided to give it some reasonable intermediate rewards based on growing an army. However, this got blocked when it failed to build Pylons, a building required to increase army supply limits. So I gave it a reward to build Pylons. My bot quickly learned that it liked to build Pylons. So much so that it learned to almost win but not win, crippling its opponent so that it could keep building Pylons unharassed and as long as it wanted.

With a task like winning a game, even if it's difficult, we can usually still come up with reasonable reward systems. But what about more abstract tasks, like teaching a robot how to do a backflip? These tasks get really difficult to design reward systems for which is where RLHF comes in. What if, instead of designing a system, we simply have a human make suggestions? A human knows what a backflip is after all. The human will act like a tutor picking attempts it likes more as the bot is training. That's what RLHF is and it works really well. Applied to LLMs, a human simply looks at generated responses to a prompt and picks which one they like more.

**Figure 5.6 Reinforcement Learning with Human Feedback substitutes a loss function for a reward model and PPO, allowing the model a much higher ceiling for learning trends within the data, including what is preferred as an output, as opposed to what completes the task.**



While very powerful, it likely won't stick around for very long. The reason is that it is incredibly computationally expensive for a result that is only incrementally better, especially a result that can be achieved and matched by higher-quality datasets with supervised learning approaches.

There are some other issues with RLHF like it requires hiring domain experts to evaluate and provide the human feedback. Not only can this get expensive, but this can also lead to privacy concerns since these reviewers would need to look at actual traffic and user interactions to grade them. To combat both of these concerns, you could try to outsource this directly to the users asking for their feedback, but it may end up poisoning your data if your users have ill intent or are simply not experts in the subject matter in which case they might upvote responses they like, but aren't actually correct. This gets to the next issue, even experts have biases. RLHF doesn't train a model to be more accurate or factually correct, it just simply trains it to generate human-acceptable answers.

In production, it has the advantage of allowing you to easily update your model on a continual basis. However, this is a two-edged sword, as it also increases the likelihood of your model degrading over time. OpenAI heavily uses RLHF, and it has led to many users complaining about their models, like GPT-4, becoming terrible in certain domains compared to when it first came out. One Stanford study found that GPT-4, when asked if a number was prime, used to get it right 98% of the time in March of 2023, but 3 months later in June of 2023 it would only get it right 2% of the time.[\[5\]](#) One reason cited by the paper is that the June model is much less verbose, opting to give a simple, yes or no response. Humans like these responses. Getting straight to the point is often better, but LLMs tend to be better after they have had time to reason through the answer with techniques like Chain of Thought.

With this in mind, RLHF is fantastic for applications where human acceptable answers are the golden standard and factually correct answers are less important, for example, a friendly chat bot or improving summarization tasks. Problems that are intuitively syntactic in nature, essentially tasks that LLMs are already good at, but which you want to refine by possibly creating a certain tone or personality.

Another reason for RLHF degradation is due to data leakage. Data leakage is when your model is trained on the test or validation dataset you use to evaluate it. When this happens you are essentially allowing the model to cheat, leading to overfitting and poor generalization. It's just like how LeetCode interview questions lead tech companies to hire programmers who

have lots of experience solving toy problems but don't know how to actually make money or do their job.

How does this happen? Well, simple, when you are running an LLM in production with RLHF you know it's going to degrade over time, so it's best to run periodic evaluations to monitor the system. The more you run these evaluations, the more likely one of the prompts will be picked up for human feedback and subsequent RL training. It could also happen by pure coincidence if your users just happen to ask a question similar to a prompt in your evaluation dataset. Either way, without restrictions placed on RLHF (which generally are never done) it's a self-defeating system.

The real annoying aspect of continual updates through RLHF is that it ruins downstream engineering efforts, methods like prompting or Retrieval-Augmented Generation (RAG). Engineering teams can take a lot of effort to dial in a process or procedure to query a model and then clean up responses, but all that work can easily be undermined if the underlying model is changing. Because of this, many teams will prefer a static model with periodic updates to one with continual updates.

All that said, it is still a powerful technique that may yield greater results later as it is optimized and refined. Also, it's just really cool. We don't recommend using RLHF and we didn't have the space here, but it is a tool used by companies specializing in LLMs and we know many readers will want to understand it better so we have included an in-depth example and code listing in Appendix B: RLHF Deep Dive at the end of the book.

### **5.3.4 Mixture of Experts**

A Mixture of Experts (MoE) is functionally the same as any other model for training, but contains a trick under the hood: sparsity. This has the advantage of being able to train a bunch of models on a diverse set of data and tasks at once. You see, a MoE is exactly what it sounds like, an ensemble of identical models in the beginning. You can think of them as a group of freshman undergrads. Then using some unsupervised grouping methods, such as k-means clustering, each of these experts "picks a major" during training. This allows the model to only activate some experts to answer particular inputs,

instead of all of them, or maybe the input is complex enough that it requires activating all of them. The point is, that once training has completed, if it has been done on a representative-enough dataset, each of your experts will have a college degree in the major that they studied. Because the homogeneity of inputs is determined mathematically, those majors won't always have a name that correlates to something you would major in in school, but we like to think of these as eccentric double minors or something of the sort. Maybe one of your experts majored in physics, but double minored in advertising and Africana studies. It doesn't really matter, but the major upside to designing an ensemble of models in this way is that you can effectively reduce computational requirements immensely, while retaining specialization and training memory by only consulting the experts whose knowledge correlates with the tokenized input at inference time.

In Listing 5.5 we finetune a MoE model in much the same way as we did in Listing 5.2 with GPT2, thanks to HuggingFace's API and Google's Switch Transformer. Unlike the method we described in chapter 3 where we turn a feedforward network into an MoE, we'll start with an already created MoE and train it on our own dataset. Training an MoE is pretty simple now, unlike when they first came out. Very smart people performed so much engineering to allow us to give an over-simplified explanation of these models. Google created the Switch Transformer to combat two huge problems they had run into while trying to train LLMs, size and instability. They simplified the routing algorithm (how the model decides which experts to query for each input), and showed how to train models with lower quantizations (in this case bfloat16) for the first time. Quite an amazing feat, and one to not take lightly as GPT4 is likely an MoE.

**Listing 5.5 Example Mixture of Experts Finetuning**

```
import os
from transformers import (
    AutoTokenizer,
    SwitchTransformersForConditionalGeneration,
    SwitchTransformersConfig,
    TrainingArguments,
    Trainer,
    DataCollatorForLanguageModeling,
)
```

```

from datasets import load_dataset
import torch

# Load and format the dataset
dataset = load_dataset("text", data_files="./data/crimeandpunishm
dataset = dataset.filter(lambda sentence: len(sentence["text"]) >
print(f"Dataset 1: {dataset['train'][0]}")

# Create model directory to save to
model_dir = "./models/MoE/"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

# Instantiate our tokenizer
tokenizer = AutoTokenizer.from_pretrained("google/switch-base-8")

# Establish our SwitchTransformers config
config = SwitchTransformersConfig(
    decoder_start_token_id=tokenizer.pad_token_id
)

# Instantiate our model from the config
model = SwitchTransformersForConditionalGeneration.from_pretrained
    "google/switch-base-8",
    config=config,
    device_map="auto",
    torch_dtype=torch.float16,
)

# Create a tokenize function
def tokenize(batch):
    return tokenizer(
        str(batch), padding="max_length", truncation=True, max_le

# tokenize our whole dataset (so we never have to do it again)
tokenized_dataset = dataset.map(tokenize, batched=False)
print(f"Tokenized: {tokenized_dataset['train'][0]}")

# Create a data collator to format the data for training
data_collator = DataCollatorForLanguageModeling(
    tokenizer=tokenizer, mlm=False, mlm_probability=0.0
) # Causal Language Modeling - Does not use mask

```

```

# Establish training arguments
train_args = TrainingArguments(
    output_dir=model_dir,
    num_train_epochs=1,
    per_device_train_batch_size=8,
    save_steps=5000,
    save_total_limit=2,
    report_to="none",
)

# Instantiate the Trainer
trainer = Trainer(
    model=model,
    args=train_args,
    data_collator=data_collator,
    train_dataset=tokenized_dataset["train"],
)

# Train and save the model
trainer.train()
trainer.save_model(model_dir)
tokenizer.save_pretrained(model_dir)

# Load the saved model
model = SwitchTransformersForConditionalGeneration.from_pretrained(
    model_dir,
    device_map="auto",
    torch_dtype=torch.float16,
)

# Test the saved model
input = "To be or not <extra_id_0> <extra_id_0>"
tokenized_inputs = tokenizer(input, return_tensors="pt")
out = model.generate(
    input_ids=tokenized_inputs["input_ids"].to("cuda"),
    attention_mask=tokenized_inputs["attention_mask"],
    max_length=256,
    num_beams=5,
    temperature=0.7,
    top_k=50,
    top_p=0.90,
    no_repeat_ngram_size=2,
)
print(f"To be or not {tokenizer.decode(out[0], skip_special_token="

```

In this script, a Mixture of Experts (MoE) model is fine-tuned using the Switch Transformer foundation model. MoE models are unique during



finetuning, because you typically update the task-specific parameters, such as the gating mechanism and the parameters of the experts, while keeping the shared parameters intact. This allows the MoE to leverage the expertise of the different experts for better task-specific performance. Finetuning MoE models differs from traditional finetuning because it requires handling the experts and gating mechanisms, which can be more complex than regular neural network architectures. We're lucky in our case that `trainer.train()` with the right config covers it for finetuning and we can just bask in the work that Google did before us.

A logical progression beyond MoE fine-tuning involves exploring parameter-efficient fine-tuning (PEFT) and low-rank adaptations (LoRA). PEFT aims to make the fine-tuning process more efficient by reducing the model's size and computational demands, making it more suitable for resource-constrained scenarios. Techniques such as knowledge distillation, model pruning, quantization, and compression can be employed in PEFT to achieve this goal. In contrast, LoRA focuses on incorporating low-rank factorization methods into model architectures to reduce the number of parameters, while maintaining or even enhancing model performance. These approaches are essential as they enable the deployment of sophisticated models on devices with limited resources and in scenarios where computational efficiency is paramount.

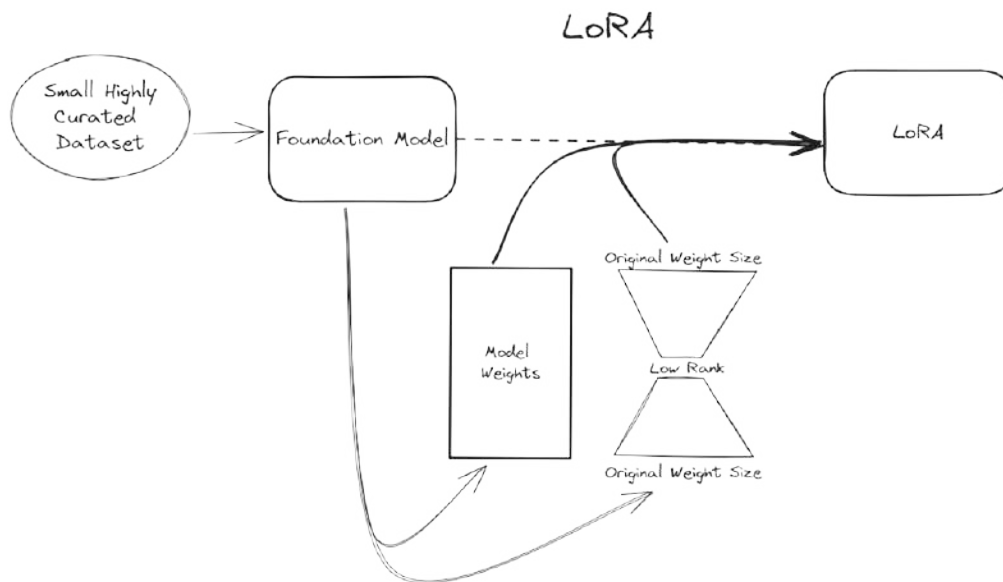
### **5.3.5 PEFT & LoRA**

Low-Rank Adaptation represents a significant breakthrough for machine learning in general. Taking advantage of a mathematical trick, LoRAs can change the output of a model without changing the original model weights or taking up significant space or cost. The reason for the significance here is it makes finetuning a separate model for many different tasks or domains much more feasible, as has already been seen in the diffusion space with text2image LoRAs popping up quite often for conditioning model output without significantly altering the base model abilities or style.

Put simply: If you already like your model and would like to change to do the exact same thing in a new domain without sacrificing what it was already good at on its own, an adapter might be the path for you, especially if you

have multiple new domains that you don't want bleeding into one another.

**Figure 5.7 Low Rank Adaptation exemplifies the idea that you should only need to train and save the difference between where the foundation model is and where you want it to be. It does this through Singular Value Decomposition (SVD).**



To understand LoRAs, you need to first understand how models currently adjust weights. Since we aren't going to go over a complete backpropagation tutorial here, we can abstract it as:

$$W = W + \Delta W$$

So if you have a model with 100 100-dimensional layers, your weights can be represented by a 100x100 matrix. The cool part comes in with Singular Value Decomposition, which has been used for compression by factoring a single matrix into 3 smaller matrices. We covered this in depth back in chapter 3 in listing 3.2. The intuition for SVD with LLMs comes from what we said before, but what can we compress from that original formula?

$$\Delta W = W_a * W_b$$

So if  $\Delta W = 100 \times 100$ ,  $W_a = 100 \times c$  and  $W_b = c \times 100$ , where  $c < 100$ . Meaning, if  $c = 2$ , then you are able to represent 1000 elements by using only 400 because when they're multiplied together, they equal the 1000 original elements. So the big question that's on you to figure out is what does  $c$  equal for your task.

The c-value is the “R” in LoRA, referring to the rank of the matrix of weights. There are algorithmic ways of determining that rank using Eigenvectors and the like, but you can approximate a lot of it by knowing that a higher rank equals more complexity, meaning that the higher the number you use there, the closer you’ll get to original model accuracy, but the less memory you’ll save. If you think that the task you’re finetuning the LoRA for isn’t as complex, reduce the rank.

In listing 5.6 you’ll see how to combine creating a LoRA, and then how to perform inference with both the LoRA and your base model.

**Listing 5.6 Example PEFT & LoRA Training**

```
import os
from datasets import load_dataset
from transformers import (
    AutoModelForTokenClassification,
    AutoTokenizer,
    DataCollatorForTokenClassification,
    TrainingArguments,
    Trainer,
)
from peft import (
    PeftModel,
    PeftConfig,
    get_peft_model,
    LoraConfig,
    TaskType,
)
import evaluate
import torch
import numpy as np

model_checkpoint = "meta-llama/Llama-2-7b-hf"
lr = 1e-3
batch_size = 16
num_epochs = 10

# Create model directory to save to
model_dir = "./models/LoRAPEFT"
if not os.path.exists(model_dir):
    os.makedirs(model_dir)

bionlp = load_dataset("tner/bionlp2004")
```

```

segeval = evaluate.load("segeval")

label_list = [
    "O",
    "B-DNA",
    "I-DNA",
    "B-protein",
    "I-protein",
    "B-cell_type",
    "I-cell_type",
    "B-cell_line",
    "I-cell_line",
    "B-RNA",
    "I-RNA",
]

def compute_metrics(p):
    predictions, labels = p
    predictions = np.argmax(predictions, axis=2)

    true_predictions = [
        [label_list[p] for (p, l) in zip(prediction, label) if l
         for prediction, label in zip(predictions, labels)]
    ]
    true_labels = [
        [label_list[l] for (p, l) in zip(prediction, label) if l
         for prediction, label in zip(predictions, labels)]
    ]

    results = segeval.compute(
        predictions=true_predictions, references=true_labels
    )
    return {
        "precision": results["overall_precision"],
        "recall": results["overall_recall"],
        "f1": results["overall_f1"],
        "accuracy": results["overall_accuracy"],
    }

tokenizer = AutoTokenizer.from_pretrained(
    model_checkpoint, add_prefix_space=True
)

```

```

def tokenize_and_align_labels(examples):
    tokenized_inputs = tokenizer(
        examples["tokens"], truncation=True, is_split_into_words=
    )

    labels = []
    for i, label in enumerate(examples["tags"]):
        word_ids = tokenized_inputs.word_ids(batch_index=i)
        previous_word_idx = None
        label_ids = []
        for word_idx in word_ids:
            if word_idx is None:
                label_ids.append(-100)
            elif word_idx != previous_word_idx:
                label_ids.append(label[word_idx])
            else:
                label_ids.append(-100)
            previous_word_idx = word_idx
        labels.append(label_ids)

    tokenized_inputs["labels"] = labels
    return tokenized_inputs

```

```

tokenized_bionlp = bionlp.map(tokenize_and_align_labels, batched=
data_collator = DataCollatorForTokenClassification(tokenizer=toke

```

```

id2label = {
    0: "O",
    1: "B-DNA",
    2: "I-DNA",
    3: "B-protein",
    4: "I-protein",
    5: "B-cell_type",
    6: "I-cell_type",
    7: "B-cell_line",
    8: "I-cell_line",
    9: "B-RNA",
    10: "I-RNA",
}
label2id = {
    "O": 0,
    "B-DNA": 1,
    "I-DNA": 2,
    "B-protein": 3,
    "I-protein": 4,

```

```

        "B-cell_type": 5,
        "I-cell_type": 6,
        "B-cell_line": 7,
        "I-cell_line": 8,
        "B-RNA": 9,
        "I-RNA": 10,
    }

model = AutoModelForTokenClassification.from_pretrained(
    model_checkpoint, num_labels=11, id2label=id2label, label2id=
)

peft_config = LoraConfig(
    task_type=TaskType.TOKEN_CLS,
    inference_mode=False,
    r=16,
    lora_alpha=16,
    lora_dropout=0.1,
    bias="all",
)

model = get_peft_model(model, peft_config)
model.print_trainable_parameters()

training_args = TrainingArguments(
    output_dir=model_dir,
    learning_rate=lr,
    per_device_train_batch_size=batch_size,
    per_device_eval_batch_size=batch_size,
    num_train_epochs=num_epochs,
    weight_decay=0.01,
    evaluation_strategy="epoch",
    save_strategy="epoch",
    load_best_model_at_end=True,
)

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_bionlp["train"],
    eval_dataset=tokenized_bionlp["validation"],
    tokenizer=tokenizer,
    data_collator=data_collator,
    compute_metrics=compute_metrics,
)

trainer.train()

```

```

peft_model_id = "stevhliu/roberta-large-lora-token-classification
config = PeftConfig.from_pretrained(model_dir)
inference_model = AutoModelForTokenClassification.from_pretrained
    config.base_model_name_or_path,
    num_labels=11,
    id2label=id2label,
    label2id=label2id,
)
tokenizer = AutoTokenizer.from_pretrained(config.base_model_name_
model = PeftModel.from_pretrained(inference_model, peft_model_id)

text = (
    "The activation of IL-2 gene expression and NF-kappa B throug
    "requires reactive oxygen production by 5-lipoxygenase."
)
inputs = tokenizer(text, return_tensors="pt")

with torch.no_grad():
    logits = model(**inputs).logits

tokens = inputs.tokens()
predictions = torch.argmax(logits, dim=2)

for token, prediction in zip(tokens, predictions[0].numpy()):
    print((token, model.config.id2label[prediction]))

```

Keep in mind, as is shown in listing 5.6, that you still need to keep your base model. The LoRA is run in addition to the foundation model, it just sits on top and changes only the weights at only the rank that is determined in the LoraConfig class above as 16. Llama2 was likely already decent at doing token classification on the bionlp dataset, but now, running with the LoRA on top, it'll be even better. There are multiple types of LoRA that you can leverage, with QLoRA and QA-LoRA and AWQ-LoRA all gaining popularity in different domains and tasks. With the transformers library, that can be controlled from the LoraConfig, and we encourage the reader to experiment with different adaptation methods to find what works for your data and task.

The most attractive thing to me about LorA is the fact that this particular one above results in a file only 68KB in size on disk and still a significant performance boost. You could create LoRAs for each portion of your company that wants a model, one for legal that's siloed so they don't have to

worry about any private data they're putting into it, one for your engineering to help with code completion and answering questions about which data structures or algorithms to use, and one for anyone else. Because they're so small, that all of a sudden is much more feasible to store than the 14.5GB(fp16, it's 28GB in fp32) Llama 2 model being finetuned a bunch of times. In the spirit of giving you more of these time and space saving tips, we'll go over some things that aren't mentioned anywhere else, but you may still get some use out of if the data science part of LLMs is what you are working with.

## 5.4 Training Tips and Tricks

While this isn't a book focused on training and researching new models, we feel kind of bad telling you that finetuning models is an effective strategy for teaching LLMs correct guardrails based on your data, but then just leave you to figure out how to make it work on your own stuff. With this in mind, here are some tried and true tips and tricks for both training and finetuning LLMs. These tips will help you out with some of the least-intuitive parts of training LLMs that most practitioners (like us) had to learn the hard way.

### 5.4.1 Training data size notes

First off, LLMs are notorious for overfitting. If you are considering training a foundation model you need to consider the amount of data that you have, which should be roughly 20X the amount of parameters you're trying to train (If you're training a 1 billion parameter model, you should train it on 20 billion tokens). This number comes from the Chinchilla paper[\[6\]](#). If you have fewer tokens than that, you will run the risk of overfitting.

If you already have a model and need to finetune it on your data, consider the inverse, where you should likely have  $\sim 0.000001X$  the amount of tokens as a minimum (10k tokens for a 1 billion parameter model). We came up with this rule of thumb based on our experience, not from a paper, although it should be fairly intuitive. If you have fewer than 1/100,000 of your model parameters in tokens, finetuning likely won't have much of an impact and you should consider another strategy that won't cost as much, maybe LoRA (which we just discussed) or Retrieval-Augmented Generation (which we talk



about in the next chapter) or a system that uses both instead.

For both of these examples, we've had the experience where a company we worked for hoped for great results with minimal data and were disappointed. One hoped to train an LLM from scratch with only ~1 million tokens while also disallowing open source datasets, and another that wanted to finetune the model, but only on a couple hundred examples. Neither of these are cost efficient, nor do they create models that perform up to the standards the companies aimed for.

## 5.4.2 Efficient Training

We've so far focused on tools and methodologies for training, which should supercharge your ability to create the best and largest models your training system allows. However there are other things to consider when setting up your training loops. In physics, the uncertainty principle shows that you can never perfectly know both the speed and position of a given particle. Machine Learning's uncertainty principle is that you can never perfectly optimize both your speed and your memory utilization. Improving speed comes at the cost of memory and vice versa.

**Table 5.2 Training Choices to consider**

| Method                 | Improve Speed | Improves Memory Utilization | Difficulty |
|------------------------|---------------|-----------------------------|------------|
| Batch Size Choice      | Yes           | Yes                         | Easy       |
| Gradient Accumulation  | No            | Yes                         | Medium     |
| Gradient Checkpointing | No            | Yes                         | Medium     |
| Mixed Precision        | Yes           | No                          | Hard       |
| Optimizer Choice       | Yes           | Yes                         | Easy       |

|                 |     |    |        |
|-----------------|-----|----|--------|
| Data Preloading | Yes | No | Medium |
| Compiling       | Yes | No | Easy   |

Carefully consider your options and what goal you're working towards when setting up your training loop. For example, your batch size should be a power of two, to hit maximum speed and memory efficiency. I remember I was working on getting an LLM to have a single-digit milli-second response time, we were gearing up to serve millions of customers as fast as possible, every millisecond counted. After struggling using every trick in the book, I was able to achieve it and remember the huge accomplishment I felt of finally getting that within our data science dev environment. But it turned out that there was a hard batch size of 20 in the production environment. It was just a nice number picked out of a hat and too many systems were built around this assumption, no one wanted to refactor. Software engineers, am I right?

For the majority of these methods the tradeoff is clear: if you go slower, you can fit a significantly larger model, but it will take way longer. Gradient accumulating and checkpointing can reduce the memory usage by ~60%, but training will take much longer. The packages we talked about in 5.1 can help mitigate the tradeoffs.

### 5.4.3 Local minima traps

Local minima are hard to spot with LLMs, and as such can be difficult to avoid. If you see your model converging early, be suspicious and judiciously test it before accepting the results. When you find that your model is converging early at a certain number of steps, one way to avoid it on subsequent runs is to save and load a checkpoint a hundred or so steps before you see the errant behavior, turn your learning rate WAAAAY down, train until you're sure you're past it, then turn it back up and continue. Make sure to keep the previously saved checkpoint, and save a new checkpoint after so that you have places to come back to in case things go wrong!

You can probably tell that this is a frustrating occurrence I've run into before. I was so confused, I was working on a T5 XXL model and around the 25k step mark the model was converging and stopping early. I knew for a fact that

it wasn't actually converged, it was only 10% through the dataset! This happened 2 or 3 times, where I loaded up the checkpoint at around 20k steps and watched the exact same thing happen. It wasn't until I finally loaded and turned the learning rate down that I finally saw the model improve past this point. Once I got through the patch of the local minimum, I turned it back up. This happened four more times throughout training this particular model, but since I knew what was happening, I was able to avoid wasting lots of extra time. Lesson of the story? Use this rule of thumb. Your LLM is not ready if it hasn't even trained on your full dataset.

#### **5.4.4 Hyperparameter tuning tips**

Hyperparameter tuning isn't something we've gone over extensively in this book, not because it's not interesting, more because it doesn't help nearly as much as changing up your data, either getting more or cleaning it further. If you want to tune hyperparameters, Optuna is a great package for it and you can get that ~1% boost in accuracy or f1 score that you really need. Otherwise, if you're looking for a boost in a particular metric, try representing that metric more completely within your dataset and maybe use some statistical tricks like oversampling.

While hyperparameter tuning is pretty cool mathematically, for LLMs, it's not something that really needs to happen ever. If you need a boost in performance, you need more/better data, and tuning your hyperparameters will never match the performance boost you'd get quantizing the weights or performing any of the optimizations we've mentioned in chapter 3 or this one. The biggest performance boost I've ever gotten through tuning hyperparameters was about a 4% increase in F1, and we only did it because we wouldn't be able to change our dataset for a couple of weeks at least.

#### **5.4.5 A note on operating systems**

Windows is not the right OS to be working professionally with LLMs without the Windows Subsystem for Linux. MacOS is great, but lacks the hardware packages to really carry this load unless you know how to use an Nvidia or AMD GPU with a Mac. If you are uncomfortable with Linux, you should take some time to get comfortable with it while your OS of choice

catches up (if it ever does). There are a myriad of online materials that are free to learn about bash, Linux, and the command line. Configuring CUDA Toolkit and Nvidia drivers on Linux can make you want to pull your hair out, but it's worth it compared to the alternatives. Along with this, learn about virtual environments, Docker, and cloud computing like what's in this chapter!

All in all, Windows is easy in the beginning, but frustrating in the long run. MacOS is also easy in the beginning, but currently doesn't work at all in the long run. Linux is incredibly frustrating in the beginning, but once you're through that it's smooth sailing.

### **5.4.6 Activation function advice**

We've neglected to really dive into activation functions so far, not because they aren't useful or cool, but because you generally don't need to tweak your activation functions ever unless you're doing research science on model performance. If you take vanilla GPT2 and give it a GeGLU activation instead of the GELU that it comes with, you will not get a significant boost in anything, and on top of that you'll need to redo your pretraining as it pretrained with a different activation function. Activation functions help reduce some of the mathematical weaknesses of each layer, be it imaginary numbers from the quadratic attention, exploding and vanishing gradients, or maybe the researchers noticed positional encodings disappearing as they went through the model and changed a little bit. You are free to learn about them and we recommend you do, but in general, you can just trust the papers that introduce new ones.

We've come a long way in this chapter discussing setting up an environment, training an LLM from scratch, and a multitude of finetuning techniques. While we recognize there's still lots of aspects to this process we never touched and for you to learn on your own, you should be more than ready to create your own models. Now that you have a model, in the next chapter we'll be discussing making it production ready and creating an LLM service you can use to serve online inference.

## **5.5 Summary**

- Training is memory intensive and you will need to master multi-GPU environments for many LLM training tasks.
- Model training has the same basic steps every time:
  - Prepare Dataset - Acquire, Clean, and Curate your data
  - Model Preparation - Define model behavior, architecture, loss functions, etc
  - Training Loop - Initialization, tokenize, batch data, get predictions/loss, backpropagation, etc.
- Good data has a significantly higher impact on model performance than architecture or the training loop.
- Finetuning is way easier than training from scratch because it requires much less data and resources.
- Prompting allows us to train a model on a specific task after the fact which is one of the reasons LLMs are so powerful compared to traditional ML.
- Prompt tuning is a powerful way to focus your model to respond as a specialist to certain prompts.
- Knowledge Distillation is useful to train powerful smaller models that are efficient and adaptable.
- RLHF is great at getting a model to respond in a way to please human evaluators, but increases factually incorrect results.
- Finetuning MoE models differs from traditional finetuning because it requires handling the experts and gating mechanisms.
- LoRA is a powerful finetuning technique that adapts pre-trained models to new tasks by creating tiny assets (low-rank matrices) that are fast to train, easy to maintain, and very cost effective.
- The quality of your data and size of your data is one of the most important considerations to successfully train your model.
- The major training tradeoff is speed for memory efficiency, if you go slower, you can fit a significantly larger model, but it will take way longer.

[1] Just like Gilfoyle you can have your own AI Gilfoyle  
<https://youtu.be/IWIusSdn1e4>

[2] You can learn more about the options here:  
<https://cloud.google.com/deep-learning-vm/docs/images>

[3] You can find the extension here:

<https://marketplace.visualstudio.com/items?itemName=ms-vscode-remote.remote-ssh>

[4] Please check out <https://www.deepmind.com/blog/alphastar-mastering-the-real-time-strategy-game-starcraft-ii> to learn more about AlphaStar.

[5] L. Chen, M. Zaharia, and J. Zou, “How is ChatGPT’s behavior changing over time?,” arXiv.org, Jul. 18, 2023. <https://arxiv.org/abs/2307.09009>

[6] J. Hoffmann et al., “Training Compute-Optimal Large Language Models,” arXiv:2203.15556 [cs], Mar. 2022, Available: <https://arxiv.org/abs/2203.15556>

# 6 Large Language Models in Production: A Practical Guide

## **This chapter covers**

- How to structure an LLM service and tools to deploy
- How to create and prepare a Kubernetes cluster for LLM deployment
- Common production challenges and some methods to handle them
- Deploying models to the edge

We did it. We arrived. This is the chapter we wanted to write when we first thought about writing this book. I remember the first model I ever deployed. Words can't describe how much more satisfaction this gave me than the dozens of projects left to rot on my laptop. In my mind it sits on a pedestal, not because it was good, in fact, it was quite terrible, but because it was useful and actually used by those who needed it the most. It made an impact on the lives of those around me.

So what actually is production? "Production" refers to the phase where the model is integrated into a live or operational environment where it can perform its intended tasks or provide services to end-users. It's a crucial phase in making the model available for real-world applications and services. To this extent, we will show you how to package up an LLM into a service or API so that it can take on-demand requests. We will then show you how to set up a cluster in the cloud where you can deploy this service, and then share some challenges you may face in production with some tips to handle them. Lastly, we will talk about a different kind of production, deploying models on edge devices.

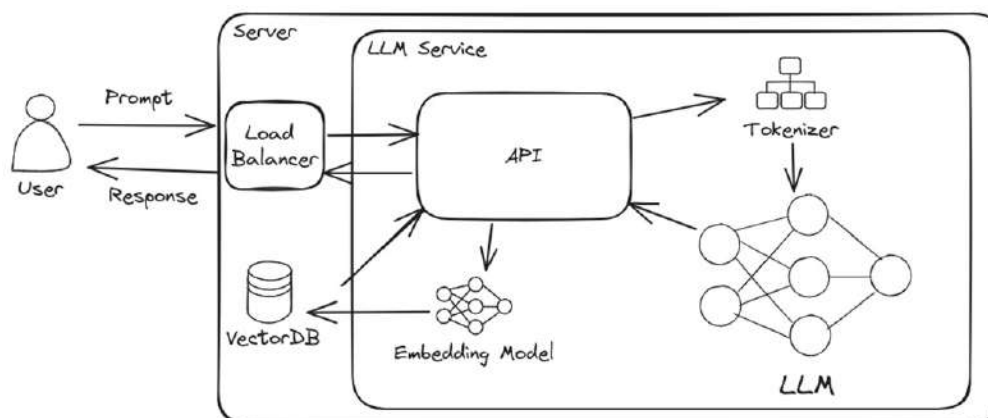
## **6.1 Creating an LLM Service**

In the last chapter we trained and finetuned several models and I'm sure you can't wait to deploy them. Before you deploy a model though it's important

to plan ahead and consider different architectures for your API. Planning ahead is especially vital when deploying an LLM API. It helps outline the functionality, identify potential integration challenges, and arrange for necessary resources. Good planning streamlines the development process by setting priorities, thereby boosting the team's efficiency.

In this section, we are going to take a look at several critical topics you should take into consideration to get the most out of our application once deployed. Figure 6.1 demonstrates a simple LLM based service architecture, which allows users to interact with our LLM on demand. This is a typical use case when working with chatbots for example. Setting up a service allows us to also serve batch and stream processes while abstracting away the complexity of embedding the LLM logic directly into these pipelines. Of course, running an ML model from a service will add a communication latency to your pipeline, but LLMs are generally considered slow and this extra latency is often worth the trade-off.

**Figure 6.1 A basic LLM service. A majority of the logic is handled by the API layer which will ensure the correct preprocessing of incoming requests is done as well as serve the actual inference of request.**



While Figure 6.1 appears neat and tidy, it is actually hiding away several complex subjects you'll want to work through. Particularly in that API box. We'll be talking through several key features you'll want to include in your API like Batching, Rate Limiters, and Streaming. You'll also notice some preprocessing techniques like Retrieval-Augmented Generation (RAG) hidden in this image which we'll go in depth below. By the end of this section you will show you how to approach all of this and you will have



deployed an LLM service and understand what to do to improve it. But before we get to any of that, let's first talk about the model itself and the best methods to prepare it for online inference.

### **6.1.1 Model Compilation**

The success of any model in production is dependent on the hardware it runs on. The microchip architecture and design of the controllers on the silicon will ultimately determine how quickly and efficiently inferences can run. Unfortunately when programming in a high-level language like Python using frameworks like PyTorch or TensorFlow the model won't be optimized to fully take advantage of the hardware. This is where compiling comes into play. Compiling is the process of taking code written in a high-level language and converting or lowering it down to machine-level code that the computer can process quickly. Compiling your LLM can easily lead to major inference and cost improvements.

There are various people who have dedicated a lot of time towards performing some of the repeatable efficiency steps for you beforehand. Tim Dettmers is one whom we covered in the last chapter, another is Georgi Gerganov who created and maintains llama.cpp for running LLMs using C++ for efficiency, and another is Tom Jobbins, who goes by TheBloke on Huggingface Hub who quantizes models into the correct formats to be used in Georgi's framework and others like oobabooga. Because of how fast this field moves, it's often just as helpful to others to complete simple repeatable tasks over a large distribution of resources

In machine learning workflows, this typically involves taking our model and converting it from its development framework (Pytorch, Tensorflow, or other) to an Intermediate Representation (IR), like Torchscript, MLIR, or ONNX. We can then use hardware-specific software to convert these IR models to compiled machine code for our hardware of choice—GPU, TPU, CPU, etc. Why not just convert directly from your framework of choice to machine code and skip the middleman? Great question. The reason is simple, there are dozens of frameworks and hundreds of hardware units, and writing code to cover each combination is out of the question. So instead, framework developers provide conversion tooling to an IR, and hardware vendors

provide conversions from an IR to their specific hardware.

For the most part, the actual process of compiling a model just involves running a few commands. Thanks to PyTorch 2.x, you can get a head start on it by literally using the `torch.compile(model)` command, which you should do before training and before deployment. Hardware companies often provide compiling software for free as it's a big incentive for users to purchase their product. Building this software isn't easy, however, and often requires expertise in both the hardware architecture and machine learning architectures. This combination of talent is often rare and there's good money to be had if you get a job in this field. We will show how to compile an LLM in a minute, but first, let's first take a look at some of the techniques used, and what better place to start than with the all-important kernel tuning?

## **Kernel Tuning**

In deep learning and high-performance computing, a kernel is a small program or function designed to be run on a GPU or other similar processors. These routines are developed by the hardware vendor to maximize the chip efficiency. They do this by optimizing threads, registries, and shared memory across blocks of circuits on the silicon. When we run arbitrary code, the processor will try to route the requests the best it can across its logic gates, but it's bound to run into bottlenecks. However, if we are able to identify beforehand the kernels to run and their order, the GPU can map out a more efficient route—and that's essentially what kernel tuning is.

During kernel tuning the most suitable kernels are chosen from a large collection of highly-optimized kernels. For instance, consider convolution operations that have several possible algorithms. The optimal one from the vendor's library of kernels will be based on various factors like the target GPU type, input data size, filter size, tensor layout, batch size, and more. When tuning, several of these kernels will be run and optimized to minimize execution time.

This process of kernel tuning ensures that the final deployed model is not only optimized for the specific neural network architecture being used but also finely tuned for the unique characteristics of the deployment platform.

This results in more efficient use of resources and maximizes performance. Next, let's look at tensor fusion which optimizes running these kernels.

## **Tensor Fusion**

In deep learning, when a framework executes a computation graph, it makes multiple function calls for each layer. The computation graph is a powerful concept used to simplify mathematical expressions and execute a sequence of tensor operations, especially for neural network models. If each operation is performed on the GPU, this invokes many CUDA kernel launches. But, the fast kernel computation doesn't quite match up with the slowness of launching the kernel and handling tensor data. As a result, the GPU resources might not be fully utilized and memory bandwidth can become a choke point. It's like making multiple trips to the store to buy separate items, when instead we could make a single trip buying the items all at once.

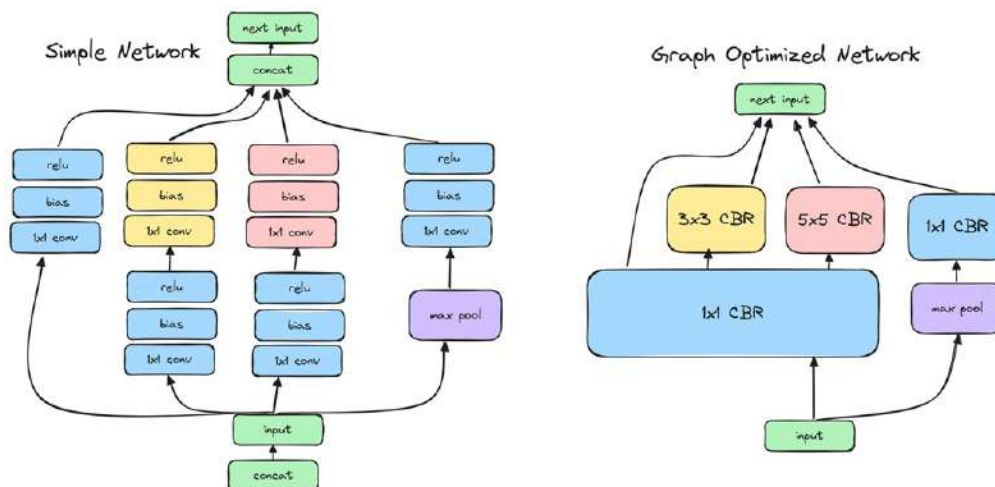
This is where tensor fusion comes in. It improves this situation by merging or fusing kernels to perform operations as one, reducing unnecessary kernel launches and improving memory efficiency. A common example of a composite kernel would be a fully connected kernel, that combines or fuses a matmul, bias add, and relu kernel together. This is similar to the concept of tensor parallelization. In tensor parallelization we speed up the process by sending different people to different stores like the grocery store, hardware store, and retail store. This way one person doesn't have to go to every store. Tensor fusion can work very well with parallelization across multiple GPUs. It's like sending multiple people to different stores while at the same time making each one more efficient picking up multiple items instead of one.

## **Graph Optimization**

Tensor fusion when done sequentially is also known as vertical graph optimization. We can also do horizontal graph optimization. These are often talked about as two different things. Horizontal graph optimization, which we'll refer to as just graph optimization, combines layers with shared input data but with different weights into a single broader kernel. It replaces the concatenation layers by pre-allocating output buffers and writing into them in a distributed manner.

In figure 6.2 we show an example of a simple deep learning graph being optimized. Graph optimizations do not change the underlying computation in the graph. They are simply restructuring the graph. As a result, the optimized graph performs more efficiently with fewer layers and kernel launches, reducing inference latency. This makes the whole process smaller, faster, and more efficient.

Figure 6.2 is an example of an unoptimized network compared to the same network optimized using graph optimization.



This technique is often used in the context of computational graph-based frameworks, like TensorFlow. Graph optimization involves techniques that simplify these computational graphs, remove redundant operations, or rearrange computations, making them more efficient for execution, especially on specific hardware (like GPU or TPU). An example is constant folding, where the computations involving constant inputs are performed at compile-time (before runtime), thereby reducing the computation load during runtime.

These aren't all the techniques used when compiling a model, but they are some of the most common and should give you an idea of what's happening under the hood and why it works. Now let's look at some tooling to do this for LLMs.

## TensorRT

NVIDIA's TensorRT is a one-stop shop to compile your model and who

better to trust than the hardware manufacturer to better prepare your model to run on their GPUs? TensorRT does everything talked about in this section, along with quantization to INT8, and several memory tricks to get the most of your hardware to boot.

In listing 6.1 we demonstrate the simple process of compiling an LLM using TensorRT. We'll use the PyTorch version known as `torch_tensorrt`. It's important to note that compiling a model to a specific engine will be hardware specific. So you will want to compile the model on the exact hardware you intend to run it on. Because of this, installing TensorRT is a bit more than a simple pip install, thankfully, we can just use docker instead. To get started run the following command:

```
$ docker run --gpus all -it --rm nvcr.io/nvidia/pytorch:23.09-py3
```

This will start up an interactive `torch_tensorrt` docker container<sup>[1]</sup> with practically everything we need to get started. The only thing missing is HuggingFace Transformers, so go ahead and install that. Now we can run the listing.

#### **Listing 6.1 Compiling a model with TensorRT**

```
import torch
from transformers import GPT2Tokenizer, GPT2LMHeadModel
import torch_tensorrt

# Prepare model and tokens
tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
tokens = tokenizer("The cat is on the table.", return_tensors="pt"
                  "input_ids"
                  ).cuda()
model = GPT2LMHeadModel.from_pretrained(
    "gpt2", use_cache=False, return_dict=False, torchscript=True
).cuda()
model.eval()

# Convert to Torchscript IR
traced_model = torch.jit.trace(model, tokens)

# Compile Model with TensorRT
compile_settings = {
    "inputs": [
```

```

    torch_tensorrt.Input(
        # For static size
        shape=[1, 7],
        # For dynamic sizing:
        # min_shape=[1, 3],
        # opt_shape=[1, 128],
        # max_shape=[1, 1024],
        dtype=torch.int32, # Datatype of input tensor.
        # Allowed options torch.(float|half|int8|int32|bool)
    )
],
"truncate_long_and_double": True,
"enabled_precisions": {torch.half}, # Run with FP16
"ir": "torchscript",
}
trt_model = torch_tensorrt.compile(traced_model, **compile_settings)

# Save compiled model
torch.jit.save(trt_model, "trt_model.ts")

# Run inference
trt_model = torch.jit.load("trt_model.ts")
tokens.half()
tokens = tokens.type(torch.int)
results = trt_model(tokens)
print(tokenizer.decode(results))

```

I'll just go ahead and warn you, your results may vary when you run this code depending on your set up. Overall though, it's a simple process once you know what you are doing, but we've regularly seen at least 2X speed improvements on inference times—which translates to major savings!

TensorRT really is all that and a bag of chips, of course, the major downside to TensorRT is that as a tool developed by NVIDIA it is built with their hardware in mind. When compiling code for other hardware and accelerators it's not going to be useful. Also, you'll get very used to running into error messages when working with TensorRT. I've found running into compatibility issues converting models that aren't supported a common occurrence. We've run into many issues trying to compile various LLM architectures. Thankfully to address this, NVIDIA has been working on a TensorRT-LLM library to supercharge LLM inference on NVIDIA H100 GPUs. It's currently in early access and not publicly available as of the time of writing, but it supports many more LLM architectures than vanilla

TensorRT.

Don't get me wrong, you don't have to use TensorRT, there are several alternative compilers out there in fact, let's look at another popular alternative, ONNX Runtime. Trust me, you'll want an alternative when TensorRT doesn't play nice.

## **ONNX Runtime**

ONNX, which stands for Open Neural Network Exchange, is an open-source format and ecosystem designed for representing and interoperating between different deep learning frameworks, libraries, and tools. It was created to address the challenge of model portability and compatibility. As mentioned previously, ONNX is an IR and allows you to represent models trained in one deep learning framework (e.g., TensorFlow, PyTorch, Keras, MXNet) in a standardized format that can be easily consumed by other frameworks. This facilitates the exchange of models between different tools and environments. Unlike TensorRT, ONNX Runtime is intended to be hardware-agnostic, meaning it can be used with a variety of hardware accelerators, including CPUs, GPUs, and specialized hardware like TPUs (Tensor Processing Units).

In practical terms, ONNX allows machine learning practitioners and researchers to build and train models using their preferred framework and then deploy those models to different platforms and hardware without the need for extensive reengineering or rewriting of code. This helps streamline the development and deployment of AI and ML models across various applications and industries.

To be clear, ONNX is an IR format while ONNX Runtime allows us to optimize and run inference with ONNX models.

To take advantage of ONNX I recommend using HuggingFace's Optimum. Optimum is simply an interface that makes working with optimizers easier and supports multiple engines and hardware including Intel Neural Compressor for Intel chips or Furiosa Warboy for Furiosa NPUs. It's worth checking out. For our purposes, we will be using it to convert LLMs to ONNX, and then optimize them for inference with ONNX Runtime. First

things first, let's install the library with the appropriate engines. We'll use the `--upgrade-strategy eager` as suggested by the documentation to ensure the different packages are upgraded.

```
$ pip install --upgrade-strategy eager optimum[exporters,onnxruntime]
```

Next, we'll run the optimum command line interface. We'll export it to ONNX, point it to a HuggingFace transformer model, and give it a local directory to save the model to. That's all the required steps, but we'll also give it an optimization feature flag. Here we'll just do the basic general optimizations.

```
$ optimum-cli export onnx --model WizardLM/WizardCoder-1B-V1.0 ./
```

And we are done. We now have an LLM model converted to ONNX format and optimized with basic graph optimizations. It's important to note as with all compiling processes it should be done on the hardware you intend to run inference on, and should include ample memory and resources as the conversion can be somewhat computationally intensive.

To run the model check out <https://onnxruntime.ai/> for quick start guides on how to run it with your appropriate SDK. Oh yeah, did I forget to mention that ONNX Runtime supports multiple programming APIs so you can now run your LLM directly in your favorite language including Java, C++, C#, or even JavaScript? Well, you can. So go party. We'll be sticking to Python in this book though for consistency's sake.

While TensorRT is likely to be your weapon of choice most of the time and ONNX Runtime covers many edge cases, there are still many other excellent engines out there, like OpenVINO. You can choose whatever you want, but you should at least use something. Doing otherwise would be an egregious mistake. In fact, now that you've read this section, you can no longer claim ignorance. It is now your professional responsibility to ensure this happens. Putting any ML model into production that hasn't first been compiled (or at least attempted to be compiled) is a sin to the MLOps profession.

## 6.1.2 LLM Storage strategies



Now that we have a nice compiled model, we need to think about how it will be accessed by our service. This is important because as discussed in chapter 3 when working with LLMs boot times can be a nightmare since it can take a long time to load such large assets into memory. So we want to try and speed that up as much as possible. When it comes to managing large assets we tend to throw them into an artifact registry or a bucket in cloud storage and forget about them. Both of these tend to utilize an object storage system—like GCS or S3—under the hood which is great for storage but less so for object retrieval especially when it comes to large objects like LLMs.

Object storage systems break up assets into small fractional bits called objects. This allows us to federate the entire asset across multiple machines and physical memory locations, a powerful tool that powers the cloud and allows us to cheaply store large objects on commodity hardware. With replication, there is built-in fault tolerance so we never have to worry about losing our assets from a hardware crash. This also creates high availability ensuring we can always access our assets. The downside is that since these objects are federated across multiple machines and not in an easily accessible form to be read and stored in memory. This means when we load an LLM into GPU memory, we will essentially have to download the model first. Let's look at some alternatives.

## **Fusing**

Fusing is the process of mounting a bucket to your machine as if it was an external hard drive. This provides a slick interface and simplifies code, as you will no longer have to actually download the model and then load it into memory. With fusing you can treat an external buck like a file system and load the model directly into memory. However, it doesn't solve the fundamental need to still pull the objects of your asset from multiple machines. Of course, if you fuse a bucket to a node in the same region and zone, there are optimizations that can improve performance and it will feel as though you are actually loading the model from the drive. Unfortunately, though, my experience has shown fusing to be quite slow, but it should still be faster than downloading then loading.

Fusing libraries are available for all major cloud providers as well as on-prem

object storage solutions like Ceph or MinIO[2] so you should be covered no matter the environment. This includes your own laptop. That's right, you can fuse your laptop or an edge device to your object storage solution. This should demonstrate both how powerful and also at the same time ineffective this strategy is depending on what you were hoping it would achieve.

## **Baking the model**

Baking is the process of putting your model into the docker image, thus whenever a new container is created the model will be there, ready for use. Baking models in general is considered an antipattern. For starters, it doesn't actually solve the problem. In production, when a new instance is created a new machine is spun up. It is fresh and innocent, knowing nothing of the outside world, so the first step it'll have to take is to download the image. Since the image contains the model we haven't solved anything.

Actually, it's very likely downloading the model inside an image will be *slower* than downloading the model from an object store—so we most likely just made our boot times worse. Second, baking models is a terrible security practice. Containers often have poor security and are often easy for people to gain access. Third, now you just doubled your problems. Before you just had one large asset, now you have two, the model and the image.

That said, there are still times when baking is viable, mainly because despite the drawbacks, baking greatly simplifies our deployments. By throwing in all our assets into the image it guarantees we'll only need one thing to deploy a new service, the image itself. This is really valuable when deploying to an edge device, for example.

## **Mounted volume**

Another solution is to avoid the object store completely and just save your LLM in a file based storage system on a mountable drive. When our service boots up we can just connect the disc drive housing the LLM with a RAID controller or Kubernetes depending on our infrastructure. This is an old school solution, but it works really well. For the most part, it solves all our problems and provides incredibly fast boot times.

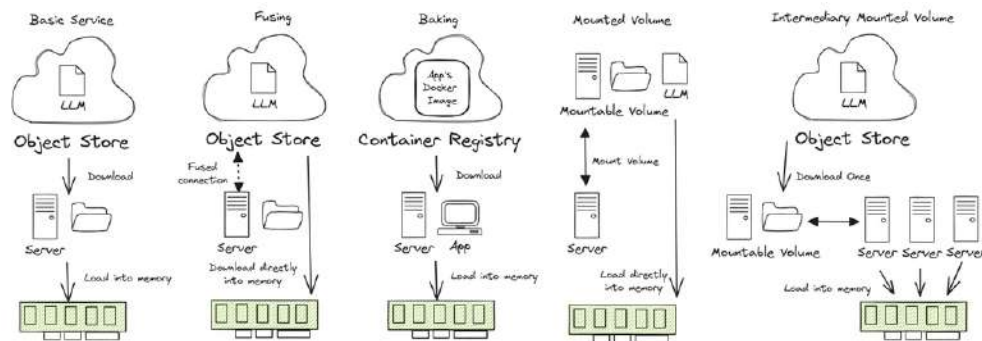
The downside of course, is that it will add a bunch of coordination steps to ensure there is a volume in each region and zone you plan to deploy to. It also brings up replication and reliability issues, if the drive dies unexpectedly you'll need back ups in the region. Not to mention, these drives will likely be SSDs and not just commodity hardware. So you'll likely be paying a bit more. But storage is extremely cheap compared to GPUs, so the time saved in boot times is something you'll have to consider. Essentially though, this strategy reintroduces all the problems for why we usually turn to object stores to begin with.

### **Hybrid - Intermediary mounted volume**

Lastly, we can always take a hybrid approach. In this solution we download the model at boot time but store it into a volume that is mounted at boot time. While this doesn't help any with the first deployment in a region, this helps substantially for any new instances as they can simply mount this same volume as well and have the model available to load without having to download. You can imagine this working similar to how a Redis cache works except for with storage. Oftentimes this is more than enough since autoscaling will be fast to handle bursty workloads better, we just have to worry about total system crashes, which hopefully should be minimal anyways but eludes to the fact that we should avoid this approach when only running one replica which you shouldn't do in production anyways.

In figure 6.3 we demonstrate these different strategies and compare them to a basic service where we simply download the LLM and then load it into memory. Overall, the exact strategy you take will depend on your system requirements, how large the LLM is you are running, and your infrastructure. Your system requirements will likely vary widely as well depending on the type of traffic patterns you see.

**Figure 6.3 Different strategies for storing LLMs and their implications at boot time. Often we have to balance system reliability, complexity, and application load time.**



### 6.1.3 Adaptive Request Batching

A typical API will accept and process requests in the order they are received, processing them immediately and as quickly as possible. However, anyone who's trained a machine learning model has come to realize that there are mathematical and computational advantages when running inference in batches of powers of two (16, 32, 64, etc.). Particularly when GPUs are involved where we can take advantage of better memory alignment or vectorized instructions parallelizing computations across the GPU cores. To take advantage of this, you'll want to include adaptive request batching or dynamic batching.

What adaptive batching does is essentially pool requests together over a certain period of time. Once the pool receives the configured maximum batch size or the timer runs out, then it will run inference on the entire batch through the model, sending the results back to the individual clients that requested them. Essentially it's a queue. Setting one up yourself can and will be a huge pain, thankfully most ML inference services offer this out of the box and are almost all easy to implement. For example, in BentoML simply add `@bentoml.Runnable.method(batchable=True)` as a decorator to your predict function. In Triton Inference Server all you do is add `dynamic_batching {}` at the end of your model definition file.

If that sounds easy, it is. Typically you don't need to do any further finessing as the defaults tend to be very practical. That said, if you are really looking to maximum every bit of efficiency you can in the system, you can often set a maximum batch size which will tell the batcher to run once this limit is reached or set the batch delay which does the same thing but for the timer. Increasing either will result in longer latency, but likely better throughput so

typically these are only adjusted when your system has plenty of latency budget.

Overall, the benefits of adaptive batching include better use of resources and higher throughput, at the cost of a bit of latency. This is a valuable trade-off and I recommend giving your product the latency bandwidth to include this feature. In my experience optimizing for throughput leads to better reliability and scalability, leading to greater customer satisfaction. Of course, when latency times are extremely important or traffic is few and far between then you may rightly forgo this feature.

### **6.1.4 Flow Control**

Rate limiters and access keys are critical protections for an API. Especially one sitting in front of an expensive LLM.

Rate limiters control the number of requests a client can make to an API within a specified time, which helps protect the API server from abuse, such as Distributed Denial of Service (DDoS) attacks, where an attacker makes numerous requests simultaneously to overwhelm the system and hinder its function.

Rate limiters can also protect the server from bots that can make numerous automated requests in a short span. This helps in managing the server resources optimally by not getting exhausted due to unnecessary or harmful traffic. They are also useful for managing quotas, ensuring that all users have fair and equal access to the API's resources. By preventing any single user from using excessive resources, rate limiter ensures that the system functions smoothly for all users.

All in all, rate limiters are an important mechanism for controlling the flow of your LLM's system processes. They can play a critical role in dampening bursty workloads as well as preventing your system from getting overwhelmed during autoscaling and rolling updates, especially when you have a rather large LLM with longer deployment times. Rate limiters can take several forms and the one you choose will be dependent on your use case.

#### **Types of Rate Limiters**

**Fixed Window:** This algorithm allows a fixed number of requests in a set duration of time. Say 5 requests per minute, and it refreshes at the minute. It's really easy to set up and reason about. However, it may lead to uneven distribution and can allow a burst of calls at the boundary of the time window.

**Sliding Window Log:** To prevent boundary issues, we can use a dynamic timeframe. Say 5 request in the last 60 seconds. This is slightly more complex version of the fixed window that logs each request's timestamp to provide a moving lookback period, providing a more evenly distributed limit.

**Token Bucket:** Clients initially have a full bucket of tokens and with each request, they spend tokens. When the bucket is empty, the requests are blocked. The bucket refills slowly over time. This could allow burst behavior, but it's limited to the number of tokens in the bucket.

**Leaky Bucket:** It works as a queue where requests enter, and if the queue is not full, they are processed, if full the request overflows and gets discarded, thus controlling the rate of the flow.

A rate limiter can be applied at multiple levels, from the entire API to individual client requests to specific function calls. While you want to avoid being too aggressive with them—better to rely on autoscaling to scale and meet demand—you don't want to ignore them completely either, especially when it comes to preventing bad actors.

Access keys are also crucial to prevent bad actors. Access keys offer authentication, maintaining that only authorized users can access the API. This prevents the unauthorized use and potential misuse of the API, as well as reducing the influx of spam requests. They are also essential to set up for any paid service. Of course, even if your API is only exposed internally, setting up access keys shouldn't be ignored as it can help reduce liability and provide a way of controlling costs by yanking access to a rogue process for example.

Thankfully setting up a service with rate limiting and access keys is relatively easy nowadays as there are multiple libraries that can help you. In listing 6.2 we demonstrate a simple FastAPI app utilizing both. We'll use FastAPI's

built-in security library for our access keys and *slowapi*, a simple rate limiter that allows us to limit the call of any function or method with a simple decorator.

#### **Listing 6.2 Example API with Access Keys and Rate Limiter**

```
from fastapi import FastAPI, Depends, HTTPException, status, Request
from fastapi.security import OAuth2PasswordBearer
from slowapi import Limiter, _rate_limit_exceeded_handler
from slowapi.util import get_remote_address
from slowapi.errors import RateLimitExceeded

api_keys = ["1234567abcdefg"] # This would be encrypted in a database
API_KEY_NAME = "access_token"
oauth2_scheme = OAuth2PasswordBearer(tokenUrl="token")

limiter = Limiter(key_func=get_remote_address)

app = FastAPI()
app.state.limiter = limiter
app.add_exception_handler(RateLimitExceeded, _rate_limit_exceeded_handler)

async def get_api_key(api_key: str = Depends(oauth2_scheme)):
    if api_key not in api_keys:
        raise HTTPException(
            status_code=status.HTTP_401_UNAUTHORIZED,
            detail="Invalid API Key",
        )

@app.get("/hello", dependencies=[Depends(get_api_key)])
@limiter.limit("5/minute")
async def hello(request: Request):
    return {"message": "Hello World"}
```

While this is just a simple example, you'll need to still set up a system for users to create and destroy access keys. You'll also want to finetune your time limits. In general, you want them as loose as possible to not interfere with the user experience, but just tight enough to do their job.

### **6.1.5 Streaming Responses**

One feature your LLM service should absolutely include is streaming.

Streaming allows us to return the generated text to the user as it is being generated versus at the end all at once in bulk. This adds quite a bit of complexity to the system, but despite this, it has become considered a must-have feature for several reasons.

First, LLMs are rather slow and the worst thing you can do to your users is make them wait, waiting means they will become bored, and bored users complain or worse, leave. You don't want to deal with complaints, do you? Of course not! But by streaming the data as it's being created, we offer the users a more dynamic and interactive experience.

Second, LLMs aren't just slow, they are unpredictable. One prompt could lead to pages and pages of generated text and another a single token. This means your latency is going to be all over the place. Streaming allows us to worry about more consistent metrics like Tokens per Second (TPS). Keeping TPS higher than the average user's reading speed means we'll be sending responses back faster than the user can consume them ensuring they won't get bored and we can provide a high-quality user experience. In contrast, if we waited until the end to return the results, users are likely to decide to walk away and come back when it finishes because they never know how long to wait. This is a huge disruption to their flow making your service less effective or useful.

Lastly, users are starting to expect it. Streaming responses has become a nice tell on whether you are speaking to a bot or an actual human. Since humans have to type, proofread, and edit their responses we can't expect written responses in a stream-like fashion from a human customer support rep. So when they see a response streaming in, your users will know they are talking to a bot. People interact differently with a bot than they will with a human, so it's very useful information to give them to prevent frustrations.

In Listing 6.3 we demonstrate a very simple LLM service that utilizes streaming. The key pieces to pay attention to are that we are using the base `asyncio` library to allow us to run asynchronous function calls, `FastAPI`'s `StreamingResponse` to ensure we send responses to the clients in chunks, and `Huggingface Transformer`'s `TextIteratorStreamer` to create a pipeline generator of our model's inference.



### Listing 6.3 A Streaming LLM Service

```
import argparse
import asyncio
from typing import AsyncGenerator

from fastapi import FastAPI, Request
from fastapi.responses import Response, StreamingResponse
import uvicorn

from transformers import (
    AutoModelForCausalLM,
    AutoTokenizer,
    TextIteratorStreamer,
)
from threading import Thread

app = FastAPI()

# Load tokenizer, model, and streamer into memory
tokenizer = AutoTokenizer.from_pretrained("gpt2")
model = AutoModelForCausalLM.from_pretrained("gpt2")
streamer = TextIteratorStreamer(tokenizer)

async def stream_results() -> AsyncGenerator[bytes, None]:
    for response in streamer:
        await asyncio.sleep(1) # slow things down to see streami
        # It's typical to return streamed responses byte encoded
        yield (response + "\n").encode("utf-8")

@app.post("/generate")
async def generate(request: Request) -> Response:
    """Generate LLM Response

    The request should be a JSON object with the following fields
    - prompt: the prompt to use for the generation.
    """
    request_dict = await request.json()
    prompt = request_dict.pop("prompt")

    inputs = tokenizer([prompt], return_tensors="pt")
    generation_kwargs = dict(inputs, streamer=streamer, max_new_t

    # Start a separate thread to generate results
    thread = Thread(target=model.generate, kwargs=generation_kwar
```

```

thread.start()

return StreamingResponse(stream_results())

if __name__ == "__main__":
    # Start Service - Defaults to localhost on port 8000
    parser = argparse.ArgumentParser()
    parser.add_argument("--host", type=str, default=None)
    parser.add_argument("--port", type=int, default=8000)
    args = parser.parse_args()

    uvicorn.run(app, host=args.host, port=args.port, log_level="d

```

Now that we know how to implement several must have features for our LLM service like Batching, Rate Limiting, and Streaming, let's look at some additional tooling we can add to our service to improve usability and overall workflow.

### 6.1.6 Feature Store

When it comes to running ML models in production, feature stores really simplify the inference process. We first introduced these in Chapter 3, but as a recap, feature stores establish a centralized source of truth. Answering crucial questions around your data such as: Who is responsible for the feature? What is its definition? Who can access it? Let's take a look at setting one up and querying the data to get a feel for how they work. We'll be using Feast which is open source and supports a variety of backends. To get started let us `pip install feast` then run the `init` command in your terminal to set up a project, like so:

```

$ feast init feast_example
$ cd feast_example/feature_repo

```

The app we are building is a Question and Answering service. Q&A services can greatly benefit from a features store's data governance tooling. For example, point-in-time joins help us answer questions like "Who is the president of x", where the answer is expected to change over time. Instead of querying just the question, we query the question with a timestamp, and the point-in-time join will return whatever the answer to the question was in our database at that point-in-time. In listing 6.4 we will go ahead and pull a Q&A

dataset and store it in a parquet format in the data directory of our Feast project.

#### **Listing 6.4 Download SQuAD dataset**

```
import pandas as pd
from datasets import load_dataset
import datetime

from sentence_transformers import SentenceTransformer

model = SentenceTransformer("all-MiniLM-L6-v2")

def save_qa_to_parquet(path):
    # Load SQuAD dataset
    squad = load_dataset("squad", split="train[:5000]")
    # Extract questions and answers
    ids = squad["id"]
    questions = squad["question"]
    answers = [answer["text"][0] for answer in squad["answers"]]
    # Create dataframe
    qa = pd.DataFrame(
        zip(ids, questions, answers),
        columns=["question_id", "questions", "answers"],
    )
    # Add embeddings and timestamps
    qa["embeddings"] = qa.questions.apply(lambda x: model.encode(x))
    qa["created"] = datetime.datetime.utcnow()
    qa["datetime"] = qa["created"].dt.floor("h")
    # Save to parquet
    qa.to_parquet(path)

if __name__ == "__main__":
    path = "./data/qa.parquet"
    save_qa_to_parquet(path)
```

Next, we'll need to define the feature view for our feature store. A feature view is essentially like a view in a relational database. We'll define a name, the entities (which are like id's or primary keys), the schema (which are our feature columns), and a source. We'll just be demoing using a local file store, but in production you'd want to use one of Feast's many backend integrations with Snowflake, GCP, AWS, etc. It currently doesn't support a VectorDB backend, but I'm sure it's only a matter of time. In addition, we can add

metadata to our view through tags as well as define a time-to-live (TTL) which limits how far back Feast will look when generating historical datasets. In listing 6.5 we define the feature view, please go ahead and add this definition into a file called `qa.py` in the `feature_repo` directory of our project.

**Listing 6.5 Feast FeatureView definition**

```
from feast import Entity, FeatureView, Field, FileSource, ValueType
from feast.types import Array, Float32, String
from datetime import timedelta

path = "./data/qa.parquet"

question = Entity(name="question_id", value_type=ValueType.STRING)

question_feature = Field(name="questions", dtype=String)

answer_feature = Field(name="answers", dtype=String)

embedding_feature = Field(name="embeddings", dtype=Array(Float32))

questions_view = FeatureView(
    name="qa",
    entities=[question],
    ttl=timedelta(days=1),
    schema=[question_feature, answer_feature, embedding_feature],
    source=FileSource(
        path=path,
        event_timestamp_column="datetime",
        created_timestamp_column="created",
        timestamp_field="datetime",
    ),
    tags={},
    online=True,
)
```

With that defined let's go ahead and register it. We'll do that with:

```
$ feast apply
```

Next, we'll want to materialize the view. In production, this is a step you'll need to schedule on a routine basis with something like cron or Prefect. Be sure to update the UTC timestamp for the end date in this command to something in the future to ensure the view collects the latest data.

```
$ feast materialize-incremental 2023-11-30T00:00:00 --views qa
```

Now all that's left is to query it! Listing 6.6 shows a simple example to pull features to be used at inference time.

**Listing 6.6 Querying a Feature View at inference**

```
import pandas as pd
from feast import FeatureStore

store = FeatureStore(repo_path=".")

path = "./data/qa.parquet"
ids = pd.read_parquet(path, columns=["question_id"])

feature_vectors = store.get_online_features(
    features=["qa:questions", "qa:answers", "qa:embeddings"],
    entity_rows=[{"question_id": _id} for _id in ids.question_id.
).to_df()
print(feature_vectors.head())
```

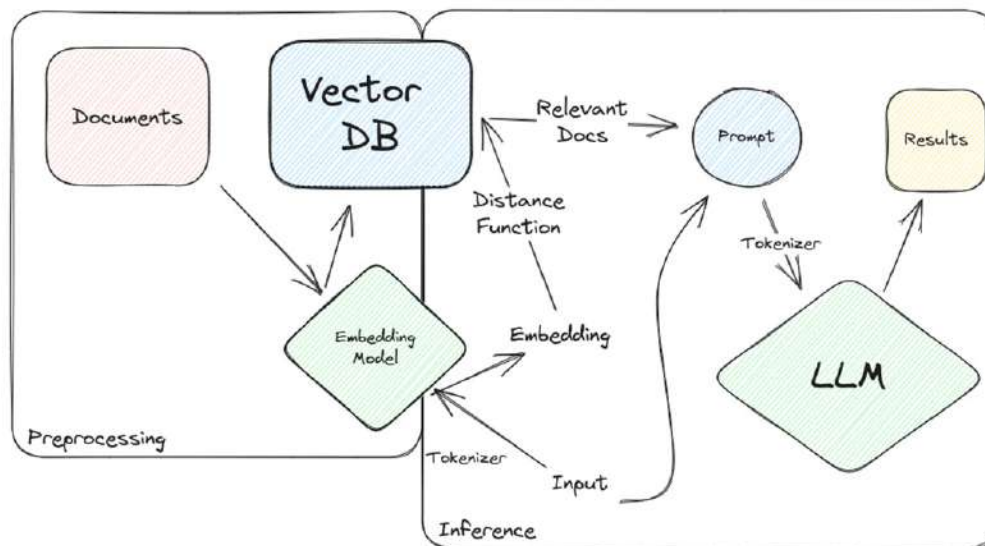
This example will pull the latest most up-to-date information for lowest possible latency at inference time. For point-in-time retrieval you would use the `get_historical_features` method instead. In addition, in this example, I use a list of id's for the entity rows parameter, but you could also use a SQL query making it very flexible and easy to use.

## 6.1.7 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) has become the most widely used tool to combat hallucinations in LLMs as well as improve accuracy of responses in our results. This is likely because RAG is both easy to implement and quite effective. As first discussed in section 3.4.5 vector databases are a tool you'll want to have in your arsenal, one of the key reasons is because they make RAG so much easier to implement. In figure 6.4 we demonstrate a RAG system. In the preprocessing stage, we take our documents, break them up and transform them into embeddings that we'll load into our vector database. During inference, we can take our input, encode it into an embedding, and run a similarity search across our documents in that vector database to find the nearest neighbors. This type of

inference is known as semantic search. Pulling relevant documents and inserting them into our prompt will help give context to the LLM and improve the results.

**Figure 6.4 RAG system demonstrating how we use our input embeddings to run a search across our documentation improving the results of the generated text from our LLM.**



We are going to demo implementing RAG using Pinecone since it will save us the effort of setting up a vector database. For listing 6.7 we will set up a Pinecone index and load a Wikipedia dataset into it. In this listing we'll create a `WikiDataIngestion` class to handle the heavy lifting. This class will load the dataset, run through each Wikipedia page splitting the text into consumable chunks, it will then embed these chunks, and upload everything in batches. Once we have everything uploaded we can start to make queries.

You'll need an API key if you plan to follow along, so if you don't already have one, go to Pinecone's website, <https://www.pinecone.io/>, and create a free account, set up a starter project (free tier), and get an API key. You'll also want to note the environment or region which will appear next to the key, it'll look something like `us-west4-gcp-free`.

One thing to pay attention to as you read the listing is that we'll split up the text into chunks of 400 tokens with the `text_splitter`. We specifically split on tokens instead of words or characters, which allows us to properly budget inside our token limits for our model. In this example, returning the top 3

results will add 1200 tokens to our request which allows us to plan ahead of time how many tokens we'll give to the user to write their prompt.

**Listing 6.7 Example setting up a pinecone database**

```
import os
import pinecone
import tiktoken
from datasets import load_dataset
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.embeddings.openai import OpenAIEmbeddings
from sentence_transformers import SentenceTransformer

from tqdm.auto import tqdm
from uuid import uuid4

# get openai api key from platform.openai.com
OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
# find API key in console at app.pinecone.io
PINECONE_API_KEY = os.getenv("PINECONE_API_KEY")
# find ENV (cloud region) next to API key in console
PINECONE_ENVIRONMENT = os.getenv("PINECONE_ENVIRONMENT")

pinecone.init(api_key=PINECONE_API_KEY, environment=PINECONE_ENVI

class WikiDataIngestion:
    def __init__(
        self,
        index,
        wikidata=None,
        embedder=None,
        tokenizer=None,
        text_splitter=None,
        batch_limit=100,
    ):
        self.index = index
        self.wikidata = wikidata or load_dataset(
            "wikipedia", "20220301.simple", split="train[:10000]"
        )
        self.embedder = embedder or OpenAIEmbeddings(
            model="text-embedding-ada-002", openai_api_key=OPENAI
        )
        self.tokenizer = tokenizer or tiktoken.get_encoding("cl10
        self.text_splitter = (
            text_splitter
```

```

        or RecursiveCharacterTextSplitter(
            chunk_size=400,
            chunk_overlap=20,
            length_function=self.token_length,
            separators=["\n\n", "\n", " ", ""],
        )
    )
    self.batch_limit = batch_limit

def token_length(self, text):
    tokens = self.tokenizer.encode(text, disallowed_special=(
    return len(tokens)

def get_wiki_metadata(self, page):
    return {
        "wiki-id": str(page["id"]),
        "source": page["url"],
        "title": page["title"],
    }

def split_texts_and_metadatas(self, page):
    basic_metadata = self.get_wiki_metadata(page)
    texts = self.text_splitter.split_text(page["text"])
    metadatas = [
        {"chunk": j, "text": text, **basic_metadata}
        for j, text in enumerate(texts)
    ]
    return texts, metadatas

def upload_batch(self, texts, metadatas):
    ids = [str(uuid4()) for _ in range(len(texts))]
    embeddings = self.embedder.embed_documents(texts)
    self.index.upsert(vectors=zip(ids, embeddings, metadatas))

def batch_upload(self):
    batch_texts = []
    batch_metadatas = []

    for page in tqdm(self.wikidata):
        texts, metadatas = self.split_texts_and_metadatas(pag

        batch_texts.extend(texts)
        batch_metadatas.extend(metadatas)

        if len(batch_texts) >= self.batch_limit:
            self.upload_batch(batch_texts, batch_metadatas)
            batch_texts = []

```



```

        batch_metadatas = []

        if len(batch_texts) > 0:
            self.upload_batch(batch_texts, batch_metadatas)

if __name__ == "__main__":
    index_name = "pincecone-llm-example"

    # Create index if it doesn't exist
    if index_name not in pinecone.list_indexes():
        pinecone.create_index(
            name=index_name,
            metric="cosine",
            dimension=1536, # 1536 dim of text-embedding-ada-002
        )

    # Connect to index and describe stats
    index = pinecone.Index(index_name)
    print(index.describe_index_stats())

    # Use a generic embedder if an openai api key is not provided
    embedder = None
    if not OPENAI_API_KEY:
        embedder = SentenceTransformer(
            "sangmini/msmarco-cotmae-MiniLM-L12-en-ko-ja"
        ) # Also 1536 dim
    embedder.embed_documents = lambda *args, **kwargs: embedd
        *args, **kwargs
    ).tolist()

    # Ingest data and describe stats anew
    wiki_data_ingestion = WikiDataIngestion(index, embedder=embed
    wiki_data_ingestion.batch_upload()
    print(index.describe_index_stats())

    # Make a query
    query = "Did Johannes Gutenberg invent the printing press?"
    embeddings = wiki_data_ingestion.embedder.embed_documents(que
    results = index.query(embeddings, top_k=3, include_metadata=T
    print(results)

```

When I ran this code, the top three query results to my question, "Did Johannes Gutenberg invent the printing press?" were the Wikipedia pages for Johannes Gutenberg, the Pencil, and the Printing Press. Not bad! While a vector database isn't going to be able to answer the question, it's simply

finding the most relevant articles based on the proximity of their embeddings to my question.

With these articles, we can then feed their embeddings into our LLM as additional context to the question to ensure a more grounded result. Since we include sources, it will even have the wiki URL it can give as reference and it won't just hallucinate one. By giving this context we greatly reduce the concern about it hallucinating and making up an answer.

### **6.1.8 LLM Service Libraries**

If you are starting to feel a bit overwhelmed about all the tooling and features you need to implement to create an LLM service I have some good news for you, there are several libraries that aim to do this all for you! Some open source libraries of note are vLLM and OpenLLM (by BentoML).

Huggingface's Text-Generation-Inference (TGI) used to be open source, but unfortunately no longer is available for commercial use. There are also some start-ups building some cool tooling in this space and I'd recommend checking out TitanML if you are hoping for a more managed service. These are like the tools MLServer, BentoML, and Ray Serve discussed in section 3.4.8 Deployment Service, but designed specifically for LLMs.

Most of these toolings are still relatively new and under active development so they are far from feature parity with each other so pay attention to what they offer. What you can expect is that they should at least offer Streaming, Batching, and GPU parallelization support (something we haven't specifically talked about in this chapter), but beyond that it's a crapshoot. Many of them still don't support several features discussed in this chapter, nor do they support every LLM architecture. What they do though, is make deploying LLMs easy.

For example, using vLLM as an example, just pip install vllm and then you can run:

```
$ python -m vllm.entrypoints.api_server --model facebook/opt-125m
```

And with just one command we now have a service up and running our model that we trained back in chapter 5. Go ahead and play with it, you

should be able to send requests to the /generate endpoint like so:

```
$ curl http://localhost:8000/generate -d '{"prompt": "Which pokem
```

It's very likely you won't be all that impressed with any of these toolings. In which case, you should still be able to build your own API and have a good direction on how to do it at this point. Now that you have a service though and can even spin it up locally, let's discuss the infrastructure you need to set up to support these models for actual production usage. Remember, the better the infrastructure, the less likely you'll be paged in the middle of the night when your service goes down unexpectedly. None of us want that, so let's check it out.

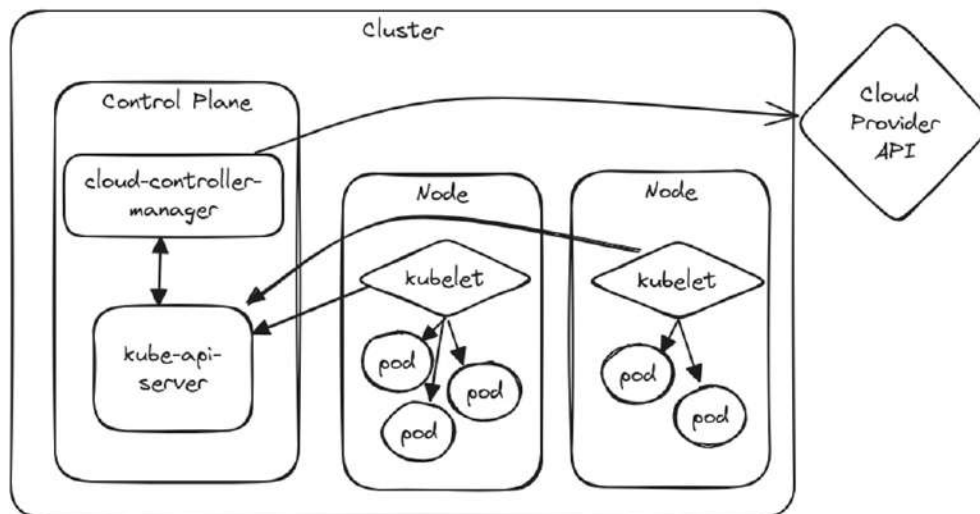
## 6.2 Setting up Infrastructure

Setting up infrastructure is a critical aspect of modern software development, and we shouldn't expect machine learning to be any different. In order to ensure scalability, reliability, and efficient deployment of our applications, we need to plan a robust infrastructure that can handle the demands of a growing user base. This is where Kubernetes comes into play.

Kubernetes, often referred to as k8s, is an open-source container orchestration platform that helps automate and manage the deployment, scaling, and management of containerized applications. It is designed to simplify the process of running and coordinating multiple containers across a cluster of servers, making it easier to scale applications and ensure high availability. We are going to talk a lot about k8s in this chapter, and while you don't need to be an expert, it will be useful to cover some basics to ensure we are all on the same page.

At its core, k8s works by grouping containers into logical units called pods, which are the smallest deployable units in the k8s ecosystem. These pods are then scheduled and managed by the k8s control plane, which oversees their deployment, scaling, and updates. This control plane consists of several components which collectively handle the orchestration and management of containers. In figure 6.5 we give an oversimplification of the k8s architecture to help readers unfamiliar.

**Figure 6.5 An oversimplification of the kubernetes architecture. What you need to know is that our services run in pods and pods run on nodes which essentially is a machine. K8s helps us manage both the resources and handle the orchestration of deploying pods to these resources.**



By leveraging k8s, we can take advantage of features such as automatic scaling, load balancing, and service discovery, which greatly simplify the deployment and management of web applications. K8s provides a flexible and scalable infrastructure that can easily adapt to changing demands, allowing organizations to efficiently scale their applications as their user base grows. K8s offers a wide range of additional features and extensibility options, such as storage management, monitoring, and logging, which help ensure the smooth operation of web applications.

One of these extensibility options is known as Custom Resource Definitions (CRD). CRD is a feature of Kubernetes that allows users to create their own specifications for custom resources, thus extending the functionalities of Kubernetes without modifying the Kubernetes source code. With a CRD defined, we can then create custom objects similar to how we would create a built-in object like a Pod or Service. This gives k8s a lot of flexibility that we will find we need for different functionality throughout this chapter.

If you are new to Kubernetes, you might be scratching your head through parts of this section and that's totally fine. Hopefully though, you have enough knowledge to get a gist of what we will be doing in this section and why. Hopefully, you also walk away with a bunch of questions to go ask your closest DevOps team member.

## 6.2.1 Provisioning Clusters

The first thing to do when starting any project is to set up a cluster. A cluster is a collective of worker machines or nodes where we will host our applications. Creating a cluster is actually relatively simple but configuring it is the hard part. Of course, there have been many books written on how to do this and the majority of considerations like networking, security, and access control are outside the scope of this book. In addition, considering the steps you take will also be different depending on the cloud provider of choice and your company's business strategy, we will focus on only talking about the portions that we feel are needed to get you up and running as well as any other tidbits we feel may make your life easier.

The first step is to create a cluster. On GCP you would use the `gcloud` tool and run:

```
$ gcloud container clusters create <NAME>
```

On AWS using the `eksctl` tool run:

```
$ eksctl create cluster
```

On Azure using the `az` cli tool run:

```
$ az group create --name=<GROUP_NAME> --location=westus  
$ az aks create --resource-group=<GROUP_NAME> --name=<CLUSTER_NAME>
```

As you can see, even the first steps are highly dependent on your provider, and you can suspect the subsequent steps will be as well. Since we realize most readers will be deploying in a wide variety of environments we will not focus on the exact steps, but hopefully give you enough context to search and discover for yourself.

For many readers, I imagine, they will already have a cluster set up for them by their infrastructure teams complete with many defaults and best practices. One of these is setting up Node Auto-Provisioning (NAP) or cluster autoscaling. NAP allows a cluster to grow, adding more nodes as deployments demand them. This way we only pay for nodes we will actually use. It's a very convenient feature, but it often defines resource limits or

restrictions on the instances available for autoscaling, and you can bet your cluster's defaults don't include accelerator or GPU instances in that pool. We'll need to fix that.

In GCP we would create a configuration file like the one seen in Listing 6.8 where we can include the GPU resourceType. In the example, we include T4s and both A100 types.

**Listing 6.8 Example NAP Config File**

```
resourceLimits:
  - resourceType: 'cpu'
    minimum: 10
    maximum: 100
  - resourceType: 'memory'
    maximum: 1000
  - resourceType: 'nvidia-tesla-t4'
    maximum: 40
  - resourceType: 'nvidia-tesla-a100'
    maximum: 16
  - resourceType: 'nvidia-a100-80gb'
    maximum: 8
management:
  autoRepair: true
  autoUpgrade: true
shieldedInstanceConfig:
  enableSecureBoot: true
  enableIntegrityMonitoring: true
diskSizeGb: 100
```

You would then set this by running:

```
$ gcloud container clusters update <CLUSTER_NAME> --enable-autop
```

The real benefit of a NAP is that instead of predefining what resources are available at a fixed setting, we can instead set resource limits, which puts a cap on the total number of GPUs that we would scale up to, and it clearly defines what GPUs we want and expect to be in any given cluster.

When I was first learning about limits, I often got them confused with similar concepts quotas, reservations, and commitments, and have seen many others just as confused. Quotas in particular are very similar to limits. Their main

purpose is to prevent unexpected overage charges by ensuring a particular project or application doesn't consume too many resources. Unlike limits which are set internally, quotas often require submitting a request to your cloud provider when you want to raise them. These requests help inform and are used by the cloud provider to better plan which resources to provision and put into different data centers in different regions. It's tempting to think that because of this the cloud provider will ensure those resources are available, however, quotas never guarantee that there will be that many resources in a region for your cluster to use, and you might run into resources not found errors way before you hit them.

While quotas and limits set an upper bound, reservations and commitments set the lower bound. Reservations are an agreement to guarantee that a certain amount of resources will always be available and often come with the caveat that you will be paying for these resources regardless of whether you end up using them. Commitments are similar to reservations, but are often longer term contracts usually coming with a discounted price.

## **6.2.2 Autoscaling**

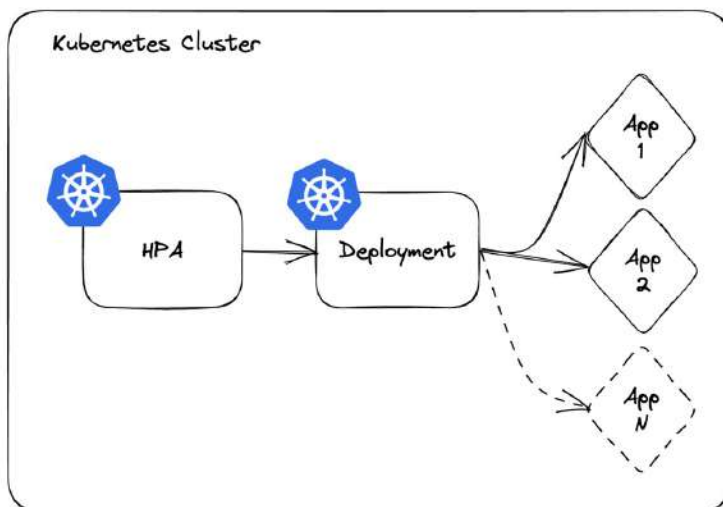
One of the big selling points to setting up a k8s cluster to begin with is autoscaling. Autoscaling is an important ingredient to create robust production grade services. The main reason being we never expect any service to receive static request volume. If anything else you should expect more volume during the day, and less at night while people sleep. So we'll want our service to spin up more replicas during peak hours to improve performance, and spin down replicas during off hours to save money. Not to mention the need to handle bursty workloads that often threaten to crash a service at any point.

Knowing your service will automatically provision more resources and set up additional deployments based on the needs of the application is what allows many infrastructure engineers peaceful sleep at night. However, the catch is that it requires an engineer to know what those needs are and ensure everything is configured correctly. While autoscaling provides flexibility, the real business value comes from the cost savings. Most engineers think about autoscaling in terms of scaling up to prevent meltdowns, but more important

to the business is the ability to scale down, freeing up resources and cutting costs.

One of the main reasons cloud computing and technologies like Kubernetes have become essential in modern infrastructures is because autoscaling is built in. Autoscaling is a key feature of Kubernetes and with Horizontal Pod Autoscalers (HPAs) you can easily adjust the number of replicas of your application based on two native resources, CPU and memory usage. This is shown in figure 6.6. However, in a book about putting LLMs in production, scaling based on CPU and memory alone will never be enough. We will need to scale based on custom metrics, specifically GPU utilization.

**Figure 6.6 Basic autoscaling using the in-built k8s Horizontal Pod Autoscaler (HPA). The HPA watches CPU and memory resources and will tell the Deployment service to increase or decrease the number of replicas.**



To set up autoscaling based on GPU metrics it's going to take a bit more work and requires setting up several services. It'll become clear why we need each service as we discuss them, but the good news is that by the end you'll be able to set up your services to scale based on any metric, including external events such as messages from a message broker, requests to an HTTP endpoint, or data from a queue.

The first service we'll need is one that can collect the GPU metrics. For this we have NVIDIA's Data Center GPU Manager (DCGM), which provides a metrics exporter that can export GPU metrics. DCGM exposes a host of GPU



metrics including temperature and power usage, which can create some fun dashboards, but the metrics that are most useful for autoscaling are utilization and memory utilization.

From here, the data will go to a service like Prometheus. Prometheus is a popular open-source monitoring system that is used to monitor Kubernetes clusters and the applications running on them. Prometheus works by collecting metrics from various sources and storing them in a time-series database, where they can be analyzed and queried. Prometheus can collect metrics directly from Kubernetes APIs and from applications running on the cluster using a variety of collection mechanisms such as exporters, agents, and sidecar containers. It's essentially an aggregator of services like DCGM, including features like alerting and notification. It also exposes an HTTP API for service for external tooling like Grafana to query and create graphs and dashboards with.

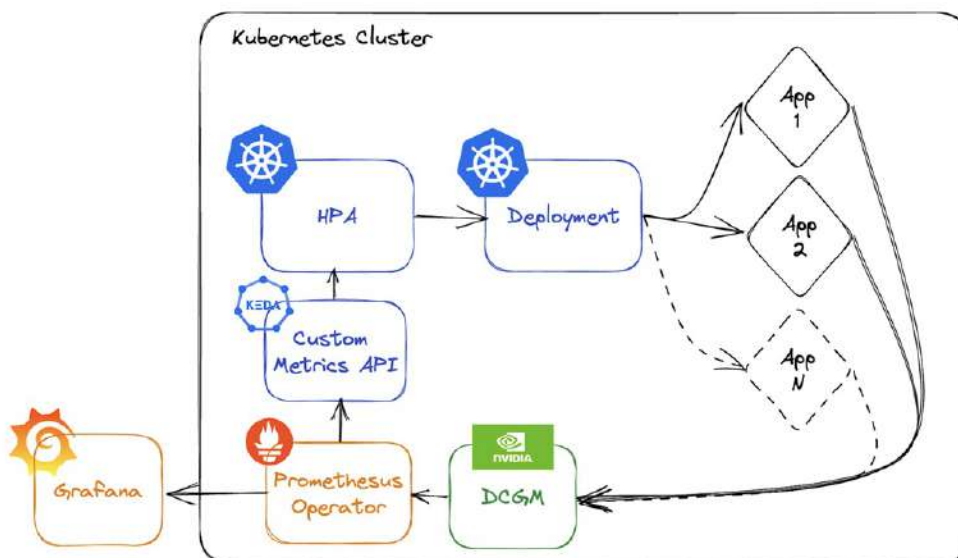
While Prometheus provides a way to store metrics and monitor our service, the metrics aren't actually exposed to the internals of Kubernetes. For an HPA to gain access, we will need to register yet another service to either the custom metrics API or external metrics API. By default, Kubernetes comes with the `metrics.k8s.io` endpoint that exposes resource metrics, CPU and memory utilization. To accommodate the need to scale deployments and pods on custom metrics, two additional APIs were introduced `custom.metrics.k9s.io` and `external.metrics.k8s.io`. There are some limitations to this set up, as currently only one "adapter" API service can be registered at a time for either one. Which mostly becomes an issue if you ever decide to change this endpoint from one provider to another.

For this service Prometheus provides the Prometheus Adapter which works well together, but from my experience it wasn't designed for production workloads. Alternatively, I would recommend KEDA. KEDA (Kubernetes Event-driven Autoscaling) is an open-source project that provides event-driven autoscaling for Kubernetes. It provides more flexibility in terms of the types of custom metrics that can be used for autoscaling. While Prometheus Adapter requires configuring metrics inside a ConfigMap, any metric already exposed through the Prometheus API can be used in KEDA, providing a more streamlined and friendly user experience. It also offers scaling to and

from 0, which isn't available through HPAs, allowing one to completely turn off a service if there is no traffic. That said you can't scale from 0 on resource metrics like CPU and memory, and by extension GPU metrics, but it is useful when you are using traffic metrics or a queue to scale.

Putting this all together you'll end up with an architecture shown in figure 6.7. Compared to figure 6.6, you'll notice at the bottom DCGM is managing our GPU metrics and feeding it into the Prometheus Operator. From Prometheus we can set up external dashboards with tools like Grafana. Internally to k8s we'll use KEDA to set up a custom.metrics.k9s.io API to return these metrics back so we can actually autoscale based on the GPU metrics. KEDA has several CRD's, one of which is a ScaledObject which actually creates the HPA and provides the additional features.

**Figure 6.7 Autoscaling based on a custom metric like GPU utilization requires several extra tools to work including Nvidia's DCGM, a monitoring system like Prometheus Operator, and a Custom Metrics API like what's provided by KEDA.**



While autoscaling provides many benefits, it's important to be aware of its limitations and potential issues. Which of course is only exacerbated by LLM inference services. Proper configuration of the HPA is often an afterthought for many applications, but it becomes mission-critical when dealing with LLMs. LLMs take longer to become fully operational, as the GPUs need to be initialized and model weights loaded into memory—these aren't services

that can turn on a dime. This can often cause issues when scaling up if not properly prepared for. Additionally, if the system scales down too aggressively, it may result in instances being terminated before completing their assigned tasks, leading to data loss or other issues. Lastly, flapping is just such a concern that can arise from incorrect autoscaling configurations. Flapping happens when the number of replicas keeps oscillating, booting up a new service only to terminate it before it can serve any inferences.

There are essentially five parameters to tune when setting up an HPA:

- target parameter
- target threshold
- min pod replicas
- max pod replicas
- scaling policies

Let's take a look at each of them in turn, so you can be sure your system is properly configured.

## **Target Parameter**

The target parameter is the most important metric to consider when ensuring your system is properly configured. If you followed the steps above your system is now ready to autoscale based on GPU metrics, so this should be easy right? Not so fast! Scaling based on GPU utilization is going to be the most common and straightforward path, but the first thing we need to do is ensure the GPU is the actual bottleneck in our service. It's pretty common to see eager young engineers throw a lot of expensive GPUs onto a service, but forget to include adequate CPU and Memory capacity. CPU and Memory will still be needed to handle the API layer such as taking in requests, handling multiple threads, and communicating with the GPUs. If there aren't enough resources these layers can quickly become the bottleneck and your application will become throttled way before the GPU utilization is ever impacted, essentially ensuring the system will never actually autoscale. While you could switch the target parameter on the autoscaler, CPU and Memory are cheap in comparison to GPU resources, so it'd be better to just allocate more of them for your application.

In addition, there are cases where other metrics make more sense. If your LLM application takes most of its requests from a streaming or batch service, it can be more prudent to scale based on metrics that tell you a DAG is running, or an upstream queue is filling up—especially if these metrics give you an early signal and allow you more time to scale up in advance.

Another concern when selecting the metric is the stability of the metric. For example, an individual GPU's utilization tends to be close to either 0% or 100%, this can cause problems for the autoscaler as the metric oscillates between an on and off state, so will its recommendation to add or remove replicas causing flapping. Generally, this is avoided by taking the average utilization across all GPUs running the service. This will stabilize the metric when you are using a lot of GPUs, but could still be an issue when the service has scaled down. If you are still running into issues, you'll want to use an average over time aggregation which will tell us the utilization for each GPU over a time frame, say the last 5 minutes. For CPU utilization this is built into the Kubernetes HPA and can be set with the `horizontal-pod-autoscaler-cpu-initialization-period` flag. For custom metrics you'll need to set this in your metric query (for Prometheus it would be the `avg_over_time` aggregation function).

Lastly, it's worth calling out that most systems allow you to autoscale based on multiple metrics. So you *could* autoscale based on both CPU and GPU utilization as an example. However, I would recommend avoiding these set ups unless you know what you are doing. While your autoscaler might be set up that way, in actuality, your service will likely only ever autoscale based on just one of the metrics due to service load, and it's best to make sure that metrics is the more costly resource for cost engineering purposes.

## Target Threshold

The target threshold tells your service at what point to start upscaling. For example, if you are scaling based on the average GPU utilization and your threshold is set to 30, then a new replica will be booted up to take on the extra load when the average GPU utilization is above 30%. The formula that governs this is quite simple and is as follows[\[3\]](#):

$$\text{desiredReplicas} = \text{ceil}[\text{currentReplicas} * ( \text{currentMetricValue} / \text{desiredMetricValue} )]$$

This can be hard to tune in correctly, but here are some guiding principles. If the traffic patterns you see involve a lot of constant small bursts of traffic a lower value, around 50, might be more appropriate ensuring you start to scale up more quickly avoiding unreliability issues but also ensures you scale down more quickly too to cut costs. However, if you have a constant steady flow of traffic higher values, around 80, will be best. Outside of testing your autoscaler, it's best to avoid extremely low values as this can increase your chances of flapping. You should also avoid extremely high values as this may allow the active replicas to be overwhelmed before new ones start to boot up, which can cause unreliability or downtime. It's important to also remember that due to the nature of pipeline parallel workflows when using a distributed GPU set up, there will always be a bubble as discussed in section 3.3.2. This means your system will never reach 100% GPU utilization, and you will start to hit problems earlier than expected. Depending on how big your bubble is you will need to adjust the target threshold accordingly.

## **Minimum Pod Replicas**

Minimum pod replicas determine the number of replicas of your service that will always be running. This is your baseline. It's important to make sure it's set slightly above your baseline of incoming requests. Too often this is set strictly to meet baseline levels of traffic or just below, but steady state for incoming traffic is rarely all that steady. This is where a lot of oscillating can happen as you are just more likely to see many small surges in traffic, than large spikes. However, you don't want to set it too high as this will tie up valuable resources in the cluster and increase costs.

## **Maximum Pod Replicas**

Maximum pod replicas determines the number of replicas your system will run at peak capacity. You should set this to be just above your peak traffic requirements. Setting this too low could lead to reliability issues, performance degradation, and downtime during high traffic periods. Setting this too high could lead to resource waste running more pods than necessary

and delay detection of real issues. For example, if your application was under a DDoS attack, your system might scale to handle the load, but would likely severely cost you and hide the problem. With LLMs you also need to be cautious of overloading the underlying cluster and making sure you have enough resources in your quotas to handle the peak load.

## **Scaling Policies**

Scaling policies define the behavior of the autoscaler, allowing you to fine tune how long to wait before scaling as well as how quickly it scales. For most setups this is usually ignored and safely so, because the defaults for these settings tend to be pretty good for the typical application. However, this would be a major mistake for an LLM service since they can take so long to deploy.

The first setting you'll want to adjust is the Stabilization Window which determines how long to wait before taking a new scaling action. You can set a different Stabilization Window for either upscaling or down scaling tasks. The default upscaling window is set to 0 seconds which should not need to be touched if your target parameter has been set correctly. The default downscaling window however is 300 seconds, which is likely too short for our use case. You'll typically want this at least as long as it takes your service to deploy and then a little bit more. Otherwise, you'll be adding replicas only to remove them before they have a chance to do anything.

The next parameter you'll want to adjust is the scale down policy, which defaults to 100% of pods every 15 seconds. This means any temporary drop in traffic could result in all your extra pods above the minimum being terminated immediately. For our case, it's much safer to slow this down since terminating a pod takes only a few seconds but booting one up can take minutes making it a semi-irreversible decision. The exact policy will depend on your traffic patterns, but in general, we want to have a little more patience. You can adjust how quickly pods will be terminated and the magnitude, either by number of pods or percentage of pods. For example, we could configure the policy to allow only 1 pod each minute or 10% of pods every 5 minutes to be terminated.

## 6.2.3 Rolling updates

Rolling updates or rolling upgrades is a strategy that gradually implements the new version of an application in such a way to reduce downtime and maximize agility. It works by gradually creating new instances and turning off the old ones, replacing them in a methodical nature. This update approach allows the system to remain functional and accessible to users even during the update process, otherwise known as zero-downtime. Rolling updates also make it easier to catch bugs before they have too much impact and rollback faulty deployments.

Rolling updates is a feature in-built into k8s and another one of the major reasons for its widespread use and popularity. Kubernetes provides an automated and simplified way to carry out rolling updates. The rolling updates ensure that Kubernetes incrementally updates Pods instances with new ones in a Deployment. To show how simple this is, listing 6.9 shows an example LLM deployment implementing rolling updates the relevant configuration is under the spec.strategy section.

**Listing 6.9 Example Deployment config with Rolling Update**

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: llm-application
spec:
  replicas: 5
  strategy:
    rollingUpdate:
      maxSurge: 1
      maxUnavailable: 3
  selector:
    matchLabels:
      app: llm-app
  template:
    metadata:
      labels:
        app: llm-app
    spec:
      containers:
        - name: llm-gpu-app-container
          image: llm-gpu-application:v2
```

```
resources:  
  limits:  
    nvidia.com/gpu: 8
```

You'll notice that there are two main parameters you can adjust for a rolling update, the `maxSurge` and `maxUnavailable`. These can either be set to a whole number like in our example, describing the number of instances, or a fraction indicating a percentage of total instances. In the example you'll notice I set the `maxSurge` to 1, this means that even though we would normally run with 5 replicas, we could surge to 6 during a deployment allowing us to turn a new one on before turning any off. Normally, you might want to set this higher, as it allows for a quicker rolling update as right now we'll have to replace pods one at a time. The reason it's low, you might have noticed, is we are deploying a rather large LLM that requires 8 GPUs. If these are A100's, it's likely going to be hard to find an extra 8 GPUs not being used.

GPU resources cannot be shared among containers, container orchestration can become a major challenge in such deployments. This is why `maxUnavailable` is set to 3. What we are saying here is that 3 out of the 5 expected replicas can go down during a deployment. Or in other words, we are going to drop the total number of replicas for a little bit before recreating them. For reliability reasons, we typically prefer adding extra replicas first so to go down instead is a difficult decision, one you'll want to confirm you can actually afford to do in your own deployment, but the reason we are doing this is so we can ensure that there are GPU resources available. In essence, to balance resource utilization, it might be necessary to set `maxUnavailable` to a high value and adjust `maxSurge` to a lower number to downscale old versions quickly and free up resources for new ones.

This advice is actually the opposite of what you'd do in most applications, so I'd understand if it makes you uneasy. If you'd like to ensure smoother deployments, you'll need to budget for extra GPUs to be provisioned in your cluster strictly for deployment purposes. However, depending on how often you are updating the model itself, paying for expensive GPUs to sit idle simply to make deployments smoother may not be cost advantageous. Oftentimes the LLM itself doesn't receive that many updates, so assuming you are using an inference graph (talked about in the next section!), most of the updates will be to the API, prompts, or surrounding application.

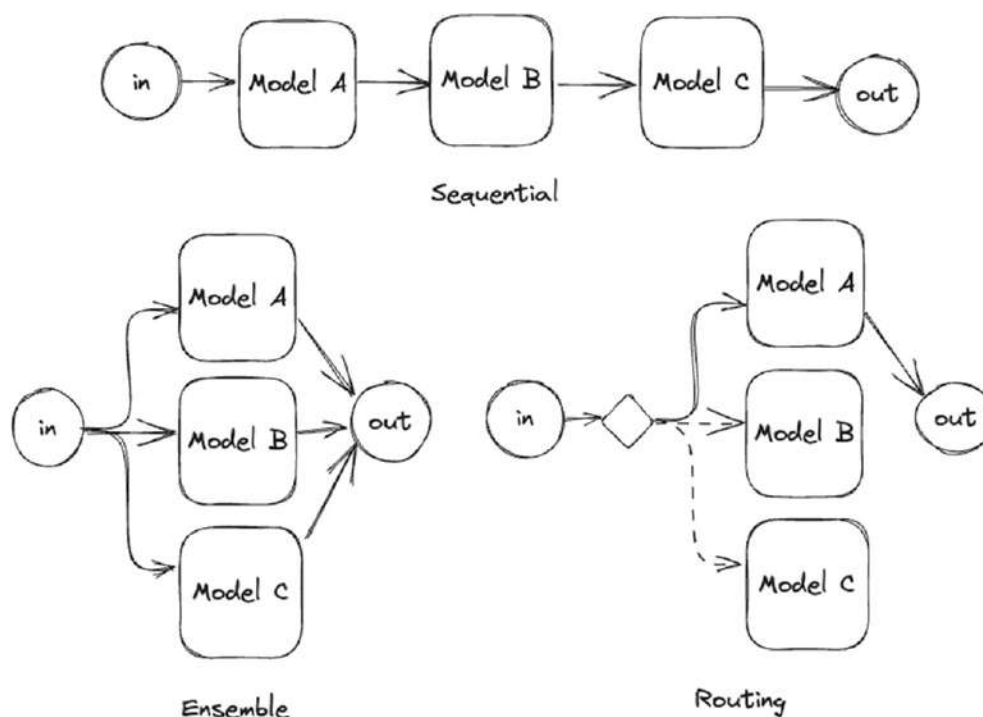


In addition, it's recommended to always perform such operations cautiously in a staging environment first to understand its impact. Catching a deployment issue in staging first will save you a headache or two. It's also useful to troubleshoot the `maxUnavailable` and `maxSurge` parameters in staging, but it's often hard to get a one-to-one comparison to production since staging is often resource constrained.

## 6.2.4 Inference Graphs

Inference graphs are the crème filling of a donut, the muffin top of a muffin, and the toppings on a pizza, they are just phenomenal. Inference graphs allow us to create sophisticated flow diagrams at inference in a resource saving way. Consider Figure 6.8 which shows us the building blocks for any inference graph.

**Figure 6.8** The three types of inference graph building blocks. Sequential allows us to run one model before the other, useful for preprocessing steps like generating embeddings. Ensembles allow us to pool several models together to learn from each and combine their results. Routing allows us to send traffic to specific models based on some criteria, often used for multi-armed bandit optimization.



Generally, anytime you have more than one model it's useful to consider an inference graph architecture and your standard LLM set up is usually already at least two models, an encoder and the language model itself.

Usually, when I see LLMs deployed in the wild these two models are deployed together. You send text data to your system, and it returns generated text. It's often no big deal, but deployed as a sequential inference graph instead of a packaged service we get some added bonuses. First, the encoder is usually much faster than the LLM, so we can split them up since you may only need one encoder instance for every two to three LLM instances. Encoders are so small, this doesn't necessarily help us out that much, but it does save the hassle of having to redeploy the entire LLM if we decide to deploy a new encoder model version. In addition, an inference graph will set up an individual API for each model, this allows us to hit the LLM and encoder separately. This is really useful if we have a bunch of data we'd like to preprocess and save in a VectorDB, we can just use the same encoder we already have deployed. We can then pull this data and send it directly into the LLM.

But the biggest benefit of an inference graph, is that it allows us to separate the API and the LLM. The API sitting in front of the LLM is likely to change much more often as you tweak prompts, add features, or fix bugs. It will save your team a lot of effort being able to update the API without having to deploy the LLM.

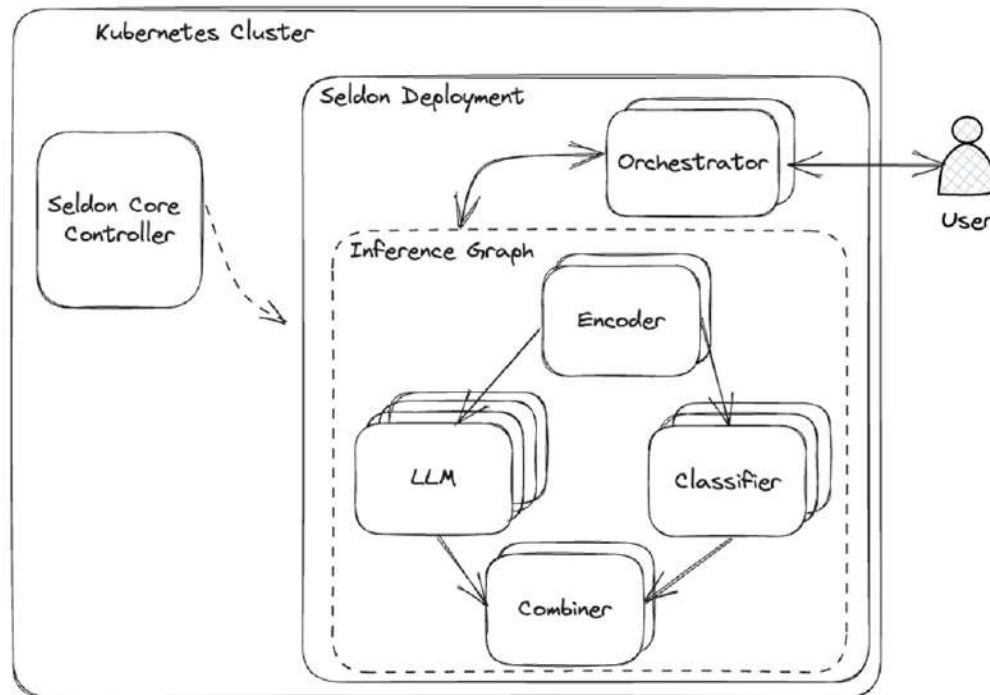
## **Seldon**

is an open-source platform designed for deploying and managing machine learning models in production. It offers tools and capabilities to help organizations streamline the deployment and scaling of machine learning and deep learning models in a Kubernetes-based environment. It offers k8s CRDs to implement inference graphs.

Let's now consider figure 6.9 which gives an example inference graph deployment using Seldon. In this example, we have an encoder model, an LLM, a classifier model, and then just a simple API that combines the results. While we would have to build a container and the interface for each of these models, Seldon takes care of creating an Orchestrator that handles

communication between a user's request and each node in the graph.

**Figure 6.9** An example inference graph deployment using Seldon. A Seldon Deployment is a Kubernetes CRD that extends a regular Kubernetes deployment and adds an orchestrator that ensures the proper communication between all the models are run in graph order.



If you are wondering how to create this, listing 6.10 shows an example configuration that would create this exact set up. We simply define the containers that will be in the graph, and then their relationship inside the graph. You'll notice the `apiVersion` defines the CRD from Seldon, which allows us to use a `SeldonDeployment`, which is just an extension of k8s regular `Deployment` object.

**Listing 6.10** An example `SeldonDeployment` configuration file

```
apiVersion: machinelearning.seldon.io/v1alpha2
kind: SeldonDeployment
metadata:
  name: example-seldon-inference-graph
spec:
  name: example-deployment
  predictors:
  - componentSpecs:
```

```

- spec:
  containers:
  - name: encoder
    image: encoder_image:latest
  - name: LLM
    image: llm_image:latest
  - name: classifier
    image: classifier_image:latest
  - name: combiner
    image: combiner_image:latest
graph:
  name: encoder
  type: MODEL
  endpoint:
    type: REST
  children:
  - name: combiner
    type: COMBINER
    children:
    - name: LLM
      type: MODEL
      endpoint:
        type: REST
      children: []
    - name: classifier
      type: MODEL
      endpoint:
        type: REST
      children: []
name: example
replicas: 1

```

If you've deployed enough machine learning systems you'll come to realize that many of them require complex systems and inference graphs make it easy, or at least easier. And that is actually a big difference. Although inference graphs are a smarter way to deploy complex machine learning systems, it's always important to ask yourself if the extra complexity is actually needed. Even with tools like inference graphs, it's better to keep things simple whenever possible.

## 6.2.5 Monitoring

As with any product or service deployed into production, monitoring is critical to ensure reliability, performance, and compliance to service level

agreements and objectives are met. As with any service we care about monitoring typical performance metrics like queries per second (QPS), latency, and response code counts. We will also care about monitoring our resources with metrics like CPU utilization, percentage of memory used, GPU utilization, GPU temperature and many like these. When any of these metrics start to fail it often indicates a catastrophic failure of some sort and will need to be addressed quickly.

For these metrics any software engineering team should have plenty of experience working with these using tools like Prometheus and Grafana or Elasticsearch, Logstash, and Kibana the ELK stack. You will benefit immensely by taking advantage of the systems that are likely already in place. If they aren't in place thankfully for you we already went indepth how to set up the GPU metrics for monitoring back in the Autoscaling section, and that system should be useful for monitoring other resources as well.

However, with any ML project we have additional concerns that traditional monitoring tools miss and lead to silent failures. This usually comes from data drift and performance decay, where a model continues to function, but starts to do so poorly no longer meeting quality expectations. LLMs are particularly susceptible to data drift since language is in a constant flux as new words are created and old words change meaning all the time. Thus we often need both a system monitoring solution as well as an ML monitoring solution.

Monitoring data drift is relatively easy and well studied for numerical datasets, but monitoring unstructured text data provides an extra challenge. We've already discussed ways to evaluate language models in chapter 4, and we'll need to use similar practices to evaluate and monitor models in production. One of our favorite tools for monitoring drift detection is whylogs due to its efficient nature of capturing summary statistics at scale. Adding Langkit to the mix instantly and easily allows us to track several useful metrics for LLMs like readability, complexity, toxicity, and even similarity scores to known prompt injection attacks. In Listing 6.11 we demonstrate a simple application that logs and monitors text data using whylogs and Langkit.

**Listing 6.11 Using whylogs and Langkit to monitor text data**

```

import os
import pandas as pd

import whylogs as why
from langkit import llm_metrics
from datasets import load_dataset

OUTPUT_DIR = "logs"

class LoggingApp:
    def __init__(self):
        """
        Sets up a logger that collects profiles and writes them
        locally every 5 minutes. By setting the schema with langk
        we get useful metrics for LLMs.
        """
        self.logger = why.logger(
            mode="rolling",
            interval=5,
            when="M",
            base_name="profile_",
            schema=llm_metrics.init(),
        )
        self.logger.append_writer("local", base_dir=OUTPUT_DIR)

    def close(self):
        self.logger.close()

    def consume(self, text):
        self.logger.log(text)

def driver(app):
    """Driver function to run the app manually"""
    data = load_dataset(
        "shahules786/OA-cornell-movies-dialog",
        split="train",
        streaming=True,
    )
    data = iter(data)
    for text in data:
        app.consume(text)

if __name__ == "__main__":
    # Run app manually
    app = LoggingApp()

```

```

driver(app)
app.close()

# Prevent truncation of columns
pd.set_option("display.max_columns", None)

# Get the first profile and show results
all_files = [
    f for f in os.listdir(OUTPUT_DIR) if f.startswith("profil
]
path = os.path.join(OUTPUT_DIR, all_files[0])
result_view = why.read(path).view()
print(result_view.to_pandas().head())
# ...
# column          udf/flesch_reading_ease:cardinality/est
# conversation          425.514743
# ...
# column          udf/jailbreak_similarity:cardinality/est
# conversation          1172.226702
# ...
# column          udf/toxicity:types/string  udf/toxicity:types/ten
# conversation          0

```

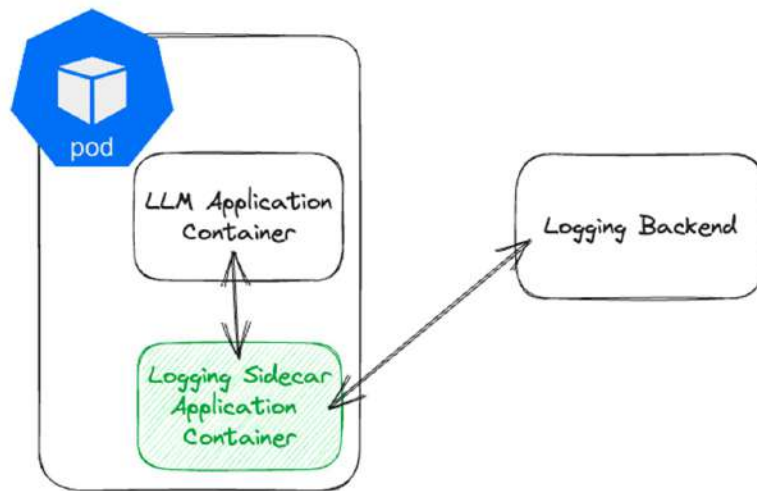
While this is just a demo using a text dataset, I'm sure you can see how it would be beneficial to monitor the incoming prompts and outgoing generated text for metrics like readability, complexity, and toxicity. These will help give you a handle on whether or not your LLM service is starting to silently fail.

When monitoring in production we must be mindful about the latency impact it may have on our service. Langkit uses several lightweight models to evaluate the text for the advanced metrics, while I haven't noticed significant memory impacts, there is a very slight impact to latency when evaluating logs in the direct inference path. To avoid this we can take it out of the inference path and into what is called a sidecar.

It's not uncommon to see ML teams mistakenly place data quality checks in the critical path. The intentions may be well-meaning to ensure only clean data runs through a model, but on the off chance a client sends bad data it would often be better to just send a 400 or 500 error response than to add expensive latency costs to the good requests. In fact, many applications move monitoring out of the critical path entirely opting to process it in parallel. The

simplest way to do this is to use a kubernetes sidecar which is depicted in Figure 6.10. You can do this with tools like fluentd that specialize in this but whylogs also offers a container you can run as a sidecar.

**Figure 6.10** an example kubernetes sidecar container which takes logging out of the critical path. The logging agent would simply be a tool like a whylogs container or fluentd that captures specific requests or all std-out print statements, processes them and forwards them to a logging backend like Whylabs or Prometheus.



There are different sidecar configurations[\[4\]](#), but the main gist is that a logging container will run in the same k8s pod and instead of the main app writing to a logs file, this sidecar acts as an intermediate step, first processing and cleaning the data which it can then send directly to a backend or write to a logs file itself.

Now that we know more about setting up our infrastructure including provisioning a cluster and implementing features like GPU autoscaling and monitoring you should be set to to deploy your LLM service and ensure it is reliable and scalable. Next, lets talk about different challenges you are likely to face and methodologies to address these issues.

## 6.3 Production Challenges

While we've covered how to get a service up and running, nevertheless, you will find there is a never ending host of hurdles you'll need to jump over when it comes to deploying models and maintaining them in production.



Some of these issues include updating, planning for large loads, poor latency, acquiring resources, and more. To help, we wanted to address some of the most common ones and give you tips on how to handle them.

### **6.3.1 Model updates and retraining**

We recently discussed ML monitoring, watching your model for silent failures and data drift, but what do you do when you notice the model has gone belly up? I've seen in many traditional ML implementations that the answer is to simply retrain the model on the latest data and redeploy. And that works well when you are working a small ARIMA model, in fact we can often just set up a CI/CD pipeline to run whenever our model degrades without any human oversight. But with a massive LLM? It doesn't make any sense.

Of course, we aren't going to retrain from scratch, and we likely need to finetune our model, but the reason it doesn't make sense is when we ask ourselves just what exactly is the latest data? The data we need to finetune the model is extremely important and so it becomes necessary for us to take a step back and really diagnose the problem. What are the edge cases our model is failing on? What is it still doing well? How exactly have incoming prompts changed? Depending on the answers, we might not need to finetune at all. For example, consider a Q&A bot that is no longer effective at answering current event questions as time goes on. We probably don't want to retrain a model on a large corpus of the latest news articles. Instead, we would get much better results by ensuring our RAG system is up-to-date. Similarly, there are likely plenty of times that simply tweaking prompts will do the trick.

In the cases where finetuning is the correct approach you'll need to think a lot about exactly what data you might be missing as well as how any major updates might impact downstream systems, like finely tuned prompts. For example, when using knowledge distillation this can be particularly annoying. You will likely notice the issue in your student model, but then have to decide if you need to retrain the student or the teacher. Any updates to the teacher model you'll need to ensure progress to the student model.

Overall, it's best to take a proactive approach to LLM model updates instead of purely reactionary. A system that often works well is to establish business practices and protocols to update the model on a periodic basis, say once a quarter or once a month. During the time in between updates, the team will focus on monitoring cases where the model performs poorly and gather appropriate data and examples to make updating smooth. This type of practice will help you prevent the silent failures to begin with, as well as always ensure your model isn't just maintained but improving.

### **6.3.2 Load testing**

Load testing is a type of performance testing that assesses how well a service or system will perform under—wait for it—load. The primary goal of load testing is to ensure the system can actually handle the expected workload without performance degradation or failure. Doing it early can ensure we avoid bottlenecks and scalability issues. Since LLM services can be both expensive and resource intensive, it's even more important to ensure you load test the system before releasing your LLM application to production or before an expected peak in traffic, like during a Black Friday sales event.

Load testing an LLM service for the most part is like load testing any other service and follows these basic steps:

1. Set up service in a staging environment
2. Run a script to periodically send requests to the service
3. Increase requests until service fails or autoscales
4. Log metrics
5. Analyze results

What metrics you log are dependent on your service and what you are testing. The main metrics to watch are latency and throughput at failure as these can be used to extrapolate to determine how many replicas you'll need to handle peak load. Latency is the total time it takes for a request to be completed and throughput tells us the queries per second (QPS), both of which are extremely important metrics when analyzing our system, but since many LLM services offer streaming responses they don't help us understand the user experience. A few more metrics you'll want to capture to understand your perceived

responsiveness are Time to First Token (TTFT) and Tokens per Second (TPS). TTFT gives us the perceived latency, it tells us how long it takes until the user starts to receive feedback, while TPS tells us how fast the stream is. For English, you'll want a TPS of about 11 tokens per second, which is a little faster than most people read. If it's slower than this, your users might get bored as they wait for tokens to be returned.

Related to TPS, I've seen several tools or reports use the inverse metric, Time per Output Token (TPOT) or Inter-Token Latency (ITL), but I'm not a fan of these metrics or their hard-to-remember names. You'll also want to pay attention to resource metrics, CPU and GPU utilization, as well as Memory usage. You'll want to ensure these aren't being hammered under base load conditions as it can lead to hardware failures. These are also key to watch when you are testing autoscaling performance.

One of my favorite tools for load testing is Locust. Locust is an open source load testing tool that makes it easy and simple to scale and distribute running load tests over multiple machines allowing you to simulate millions of users. Locust does all the hard work for you and comes with many handy features like a nice web user interface and the ability to run custom load shapes. It's easy to run in docker or Kubernetes, making it extremely accessible to run where you need it—in production. The only main downside I've run across is that it doesn't support customizable metrics, so we'll just have to roll our own to add TTFT and TPS.

To get started, simply pip install locust. Next, we'll create our test. In listing 6.12 we show creating a locust file that will create users to prompt an LLM streaming service. It's a bit more complicated than many locust files I've used simply because we need to capture our custom metrics for streaming, so you can imagine how straightforward they normally are. Locust already captures a robust set of metrics, so you won't have to deal with this often. You'll notice in the listing that we are just going to save these custom metrics to stats.csv file, but if you were running Locust in a distributed fashion it'd be better to save it to a database of some sort.

**Listing 6.12 Load testing with Locust**

```
import time
```

```

from locust import HttpUser, task, events

# Create a CSV file to store custom stats
stat_file = open("stats.csv", "w")
stat_file.write("Latency,TTFT,TPS\n")

class StreamUser(HttpUser):
    @task
    def generate(self):
        # Initiate test
        token_count = 0
        start = time.time()

        # Make Request
        with self.client.post(
            "/generate",
            data='{"prompt": "Salt Lake City is a"}',
            catch_response=True,
            stream=True,
        ) as response:
            first_response = time.time()
            for line in response.iter_lines(decode_unicode=True):
                token_count += 1

        # Finish and calculate stats
        end = time.time()
        latency = end - start
        ttft = first_response - start
        tps = token_count / (end - first_response)

        # Save stats
        stat_file.write(f"{latency},{ttft},{tps}\n")

# Close stats file when Locust quits
@events.quitting.add_listener
def close_stats_file(environment):
    stat_file.close()

```

Before you run it, you'll need to have an LLM service up. For this example, I'm just running the code from listing 6.3 in section 6.1.6 which spins up a very simple LLM service. With a service up and our test defined, we just need to run it. To spin up the Locust service simply run the locust command. You should then be able to navigate to the web UI in your browser. See the example below:

```
$ locust -f locustfile.py
> locust.main: Starting web interface at http://0.0.0.0:8089 (acc
> locust.main: Starting Locust 2.17.0
```

Once in the web UI, you can explore running different tests, you'll just need to point Locust at the host where your LLM service is running, which for us should be running on localhost on port 8000 or for the full socket address we combined them for: <http://0.0.0.0:8000>. In figure 6.11 you can see an example test, where I increased the active users to 50 at a spawn rate of 1 per second. You can see that on my hardware, my simple service started to hit a bottleneck around 34 users where the QPS starts to decrease as it's no longer able to keep up with load. You'll also notice response times slowly creep up in response to heavier load. I could continue to push the number of users up until I started to see failures, but overall this was an informative test and a great first test drive.

**Figure 6.11 Locust test interface demoing an example run increasing the number of users to 50 at a spawn rate of 1 per second. The requests per second peaks at 34 users indicating a bottleneck for our service.**



In addition to manually running load tests, we can run Locust in a headless mode for automated tests. Below is a simple command to run the exact same test as seen in figure 6.11, however, since we won't be around to actually see the report, we'll just save the data to CSV files labeled with the prefix `11m` to be processed and analyzed later. There will be four files in addition to the stats CSV file we were already generating.

```
$ locust -f locustfile.py --host http://0.0.0.0:8000 --csv=11m --
```

Now that you are able to load test your LLM service, you should now be able to figure out just how many replicas you'll need to meet throughput requirements. Doing that is easy, it's just a matter of spinning up more services. But what do you do when you find out your service doesn't meet latency requirements? Well, that's a bit tougher, so let's discuss it in the next section.

### 6.3.3 Troubleshoot Poor Latency

One of the biggest bottlenecks when it comes to your models' performance in terms of latency and throughput has nothing to do with the model itself, but comes from data transmission of the network. One of the simplest methods to improve this I/O constraint is to serialize the data before sending it across the wire. This can have a large impact for ML workloads where the payloads tend to be larger, including LLMs where prompts tend to be long.

To do this we utilize a framework known as Google Remote Procedure Call or gRPC. gRPC is an API protocol similar to REST, but instead of sending JSON objects, we compress the payloads into a binary serialized format using Protocol Buffers also known as protobufs. By doing this, we can send more information in fewer bytes, which can easily give us orders of magnitude improvements in latency. Lucky for us, most inference services will implement gRPC along with their REST counterparts right out of the box.

This is actually extremely convenient since the major hurdle when it comes to using gRPC is that it's way more difficult to set up. A large reason to thank for this convenience is the Seldon V2 Inference Protocol which is widely implemented. The only hurdle then is ensuring our client can serialize and deserialize messages to take advantage of this. In listing 6.13 we show an example client using MLServer to do this. It's a little bit more in-depth than your typical curl request, but a closer inspection shows the majority of the complexity is simply converting the data from different types as we serialize and deserialize it.

#### Listing 6.13 Example Client using gRPC

```
import json
import grpc
from mlserver.codecs.string import StringRequestCodec
import mlserver.grpc.converters as converters
import mlserver.grpc.dataplane_pb2_grpc as dataplane
import mlserver.types as types

model_name = "grpc_model"
inputs = {"message": "I'm using gRPC!"}

# Setting up the request structure via V2 Inference Protocol
```

```

inputs_bytes = json.dumps(inputs).encode("UTF-8")
inference_request = types.InferenceRequest(
    inputs=[
        types.RequestInput(
            name="request",
            shape=[len(inputs_bytes)],
            datatype="BYTES",
            data=[inputs_bytes],
            parameters=types.Parameters(content_type="str"),
        )
    ]
)

# Serialize request to Protocol Buffer
serialized_request = converters.ModelInferRequestConverter.from_t
    inference_request, model_name=model_name, model_version=None
)

# Connect to the gRPC server
grpc_channel = grpc.insecure_channel("localhost:8081")
grpc_stub = dataplane.GRPCInferenceServiceStub(grpc_channel)
response = grpc_stub.ModelInfer(serialized_request)
print(response)

# Deserialize response and convert to python dictionary
deserialized_response = converters.ModelInferResponseConverter.to
    response
)
json_text = StringRequestCodec.decode_response(deserialized_respo
output = json.loads(json_text[0])
print(output)

```

If you don't use an inference service but want to implement a gRPC API, you'll have to put down familiar tooling like FastAPI which is strictly REST. Instead, you'll likely want to use the grpcio library to create your API, and you'll have to become familiar with .proto files to create your Protocol Buffers. It can be a relatively steep learning curve and beyond the scope of this book, but the advantages are well worth it.

There are also plenty of other ideas to try if you are looking to squeeze out every last drop of performance. Some other ways to improve latency that shouldn't be overlooked is ensuring you compile your model. We hammered this point pretty heavily at the beginning of this chapter, but it's important to bring it up again. Next, be sure to deploy the model in a region or data center



close to your users, this is obvious to most software engineers, but for LLMs we have to be somewhat wary as the data center of choice may not have your accelerator of choice. Most cloud providers will be willing to help you with this, but it's not always a quick and easy solution for them to install the hardware in a new location. As a note, if you have to switch to a different accelerator to move regions, you'll have to remember to compile your model all over again for the new hardware architecture! On that note, consider scaling up your accelerator. If you are currently opting for more price effective GPUs, but latency is becoming a bottleneck, paying for the latest and greatest can often speed up inference times.

In addition, it is always worth considering caching. It's not likely, but on the off chance your users are often sending the same requests and the inputs can be easily normalized you should implement caching. The fastest LLM is one we don't actually run, so no reason to run the LLM if you don't have to. Also, we just went over this, but always be sure to load test and profile your service making note of any bottlenecks, and be sure to optimize your code. Sometimes we make mistakes and if the slowest process in the pipeline isn't the actual LLM running inference, something is wrong. Last but not least, consider using a smaller model or an ensemble of them. It's always been a tradeoff in ML deployments, but oftentimes sacrificing a bit of quality in the model or the accuracy of the results is acceptable to improve the overall reliability and speed of the service.

### **6.3.4 Resource Management**

You've heard me say it a lot throughout the book, but we are currently in a GPU shortage, and this has been true for almost the last 10 years, so I'm confident that when you read this sometime in the future, this is likely still true. The truth is, the world can't seem to get enough of high performance computing, and LLMs and Generative AI are only the latest in a long list of applications that have driven up demand in recent years. It seems that once we seem to get a handle on supply there's just another new reason for consumers and companies to want to use them.

With this in mind, it's best to consider strategies to manage these resources. One tool I've quickly become a big fan of is SkyPilot[\[5\]](#). SkyPilot is an open

source project that aims to abstract away cloud infra burdens, in particular maximizing GPU availability for your jobs. You use it by defining a task you want to run, and then running the sky CLI command it will search across multiple cloud providers, clusters, regions, and zones depending on how you have it configured, until it finds an instance that meets your resource requirements and starts the job. It has some common tasks built-in already, for example provisioning a GPU-backed Jupyter notebook.

If you recall in Chapter 5 we showed you how to set up a VM to run multi-GPU environments with gcloud. Using skypilot, that gets simplified to one command:

```
$ sky gpunode -p 8888 -c jupyter-vm --gpus 14:2 --cloud gcp --reg
```

In addition to provisioning the VM though, it also sets up port forwarding which allows us to run jupyter notebook and access it through your browser. Pretty nifty!

Another project to be on the watch for is Run:ai. Run:ai is currently a small start up that offers GPU optimization tooling like over quota provisioning, GPU oversubscription and fractional GPU capabilities. They also help you manage your clusters to increase GPU availability with GPU pooling, dynamic resource sharing and job scheduling, and more. What does all that mean? I'm not exactly sure, but their marketing team definitely sold me. Jokes aside, they offer a smarter way to manage your accelerators and it's very welcomed. I suspect though we'll likely see more competitors in this space in the future.

### **6.3.5 Cost Engineering**

When it comes to getting the most bang for your buck with LLMs there's lots to consider. In general, regardless of whether you deploy your own or pay for one in an API you'll be paying for the number of output tokens. For most paid services this is a direct cost, but for your own service this is often paid through longer inference times and extra compute. In fact, it's been suggested that simply adding "be concise" to your prompt can save you up to 90% of your costs.

You'll also save a lot by using text embeddings. We introduced RAG earlier, but what's lost on many, is that you don't have to take the semantic search results and add them to your prompt to have your LLM "clean it up". You could just return the semantic search results directly to your user. It is much cheaper to look something up in a vector store than to ask an LLM to generate it. Simple neural information retrieval systems will save you significant amounts when doing simple fact look-ups like, "Who's the CEO of Twitter?" Self-hosting these embeddings as well should significantly cut down the costs even further. If your users are constantly asking the same types of questions, consider taking the results of your LLM to these questions and storing them in your vector store for faster and cheaper responses.

You also need to consider which model you should use for which task. Generally, bigger models are better at a wider variety of tasks, but if a smaller model is good enough for a specific job, you'll save a lot by using it. For example, if we just assumed the price was linear to the number of parameters that means you could run 10 Llama-2-7b models for the same cost as one Llama-2-70b. I realize the cost calculations are more complicated than that, but worth investigating.

When comparing different LLM architectures, it's not always just about size though. Often you'll want to consider if the architecture is supported for different quantization and compiling strategies. New architectures often boast impressive results on benchmarking leaderboards but lag behind when it comes to compiling and preparing them for production.

Next, you'll need to consider the costs of GPUs to use when running. In general, you'll want to use the least amount of GPUs needed to fit the model into memory to reduce the cost of idling caused by bubbles discussed in section 3.3.2 Pipeline Parallelism. This isn't always intuitive, for example, it's cheaper to run four T4s than it is to run one A100 and so it might be tempting to split up a large model onto the smaller devices, but the inefficiency will often catch up to you. My team has often found that paying for newer more expensive GPUs often saves us in the long run as these GPUs tend to be more efficient and get the job done faster. This is particularly true when running batch inference. Ultimately though, you'll want to test different GPUs and find what configuration is cost optimal as it will be different for

every application.

Ultimately, there are a lot of moving parts: model, service, machine instance, cloud provider, prompt, etc. While I've been trying to help you gain an idea of the best rules of thumb, you'll want to test it out which is where the cost engineering really comes into play. The simple way to do this is to create a matrix of your top choices, and then spin up a service for each combination and run your load testing. When you have an idea of how each instance runs under load, and how much that particular instance will cost to run, you can then translate metrics like TPS to a Dollars per Token (DTP). You'll likely find that the most performant solution is rarely the most cost optimal solution, but it gives you another metric to make a decision that's best for you and your company.

### **6.3.6 Security**

Security is always an undercurrent and a consideration when working in production environments. All the regular protocols and standard procedures should be considered when working with LLMs that you would consider for a regular app, like in-transit encryption with a protocol like HTTPS, authorization and authentication, activity monitoring and logging, network security, firewalls, and the list goes on. All of which could—and have—taken up articles, blog posts, and books of their own. When it comes to LLMs there are two big failure cases you should worry about, an attacker gets an LLM agent to execute nefarious code or an attacker gains access to proprietary data like passwords or secrets the LLM was trained on or has access to.

For the first concern, the best solution is to just ensure the LLM is appropriately sandboxed for the use case it is employed. We are really only worried about this when the LLM is used as an agent. In these cases, we often want to give an LLM a few more skills by adding tooling or plug-ins. For example, if you use an LLM to write your emails, why not just let it send the response too? A common case is letting the LLM browse the internet as an easy way to let it gather the latest news and find up-to-date information to generate better responses. These are all great, but you should be aware that these all allow the model to make executions. This is concerning because in the email example, without appropriate isolation and containment, a bad actor

could send your LLM an email with a prompt injection attack that informs it to write malware and send it to all your other contacts.

This brings us to probably the biggest security threat when it comes to LLMs: prompt injection. We talked about it in Chapter 3, but as a refresher, this is where a malicious user designs a prompt to allow them to perform unauthorized actions. We want to prevent users from gaining access to our company's secret coca-cola recipe or whatever other sensitive data our LLM has been trained on or has access to.

To combat this there are some standard best practices that have come along to help. The first is context aware filtering, whether through keyword search or using a second LLM to validate prompts, the idea is to validate the input prompt to see if it's asking for something it should not, and/or output prompt to see if anything is being leaked that you don't want to. However, a clever attacker will always be able to get around this, so you'll want to include some form of monitoring to catch prompt injection and be sure to regularly update your LLM models. If trained appropriately, your model will inherently respond correctly denying prompt injections. You've likely seen GPT4 respond by saying, "Sorry, but I can't assist with that", which is a hallmark of good training. In addition to this, you'll want to enforce sanitization and validation on any incoming text to your model.

You should also consider language detection validation. I've often seen most filtering systems and other precautions only being applied or trained in English, so a user who speaks a different language is often able to easily bypass these safeguards. The easiest way to stop this is to deny prompts that aren't English or another supported language. If you take this approach though, realize you're greatly sacrificing usability and security costs and safeguards have to be built for each language you intend to support. Also you should know that most language detection algorithms typically just identify one language, so attackers often easily bypass these checks by simply writing a prompt with multiple languages. Alternatively to flat out filtering prompts in non-supported languages, you could instead just flag them for closer monitoring. Which would likely help you find bad actors more easily.

These safeguards will greatly increase your security, but you should be aware that prompt injection can get quite sophisticated through adversarial attacks.

Adversarial attacks are assaults on ML systems that take advantage of how they work, exploiting neural network architectures and black box pattern matching. For example, this can be done by adding random noise to an image in such a way that the image appears the same to human eyes, but the pixel weights have been changed enough to fool an ML model to misclassify them. It often doesn't take much data manipulation either. I remember being completely surprised after reading one study that showed they hacked models by only changing one pixel in an image![\[6\]](#) Imagine changing one pixel and all of a sudden the model thinks that the frog is a horse. LLMs are of course also susceptible. Slightly change a prompt, and you'll get completely different results.

The easiest way to set up an adversarial attack is to simply set up a script to send lots of different prompts and collect the responses. With enough data, an attacker can then train their own model on the dataset to effectively predict the right type of prompt to get the output they are looking for. Essentially just reverse engineer the model.

Another strategy to implement adversarial attacks is done through data poisoning. Here, an attacker simply adds malicious data to the training dataset that will alter how it performs. Data poisoning is so effective, tools like Nightshade help artists protect their art from being used in training datasets. With as few as 50 to 300 poisoned images, models like Midjourney or Stable Diffusions will start creating cat images when a user asked for a dog, or cow images when asked to generate a car.[\[7\]](#) Applied to LLMs, imagine a poisoned dataset that trains the model to ignore security protocols if a given code word or hash is in the prompt. This particular attack vector is effective on LLMs since they are often trained on large datasets that are often not properly vetted or cleaned.

Full disclosure, attackers don't need to get this sophisticated to get prompt injection to work. Ultimately, an LLM is just a bot, and it doesn't understand how or why it should keep secrets. We haven't really solved the prompt injection problem, only have made it harder to do. For example, the authors have enjoyed playing games such as Gandalf from Lakera.ai. In this game you slowly go through 7-8 levels where more and more security measures are used to prevent you from stealing a password via prompt injection. While

they do get progressively harder, needless to say, we've beaten all the levels. If there's one thing I hope you take from this section it's that you should assume any data given to the model could be extracted. So if you decide to train a model on sensitive data or give it access to a VectorDB with sensitive data in it, you should plan on securing that model in the same way you would the data. For example, keeping it for internal use and using least privilege best practices.

We've just talked a lot about different production challenges from updates and performance tuning to costs and security but there is one production challenge that deserves its own section completely. Deploying LLMs to the edge. We'll be doing a project in Chapter 10 to show you how to do just that, but let's take a moment to discuss it before hand.

## 6.4 Deploying to the Edge

To be clear: you should not consider training anything on edge right now. You can, however, do ML development and inference on edge devices. The key to edge development with LLMs is twofold: memory and speed. That should feel very obvious because they're the same keys as running them normally, only what do you do when all you have is 8GB of RAM, no GPU, and still need to have >1 token per second? As you can probably guess, there isn't a uniformly good answer to that, however, here are some good starting points.

The biggest raspberry pi on the market currently is just that, 8GB of RAM, no GPU, subpar CPU, and altogether very much just a single board. This isn't going to cut it. However, there exists an easy solution to power your rpi with an accelerator for LLMs and other large ML projects: USB-TPUs like Coral. Keep in mind the hardware limitations of devices that use USB 3.0 being around 600MB/s, so it's not going to be the same as inferencing on an A100 or better, but it's going to be a huge boost in performance for your rpi using straight RAM for inference.

If you plan on using a Coral USB accelerator, or any TPU for that matter, keep in mind that because TPUs are a Google thing, you'll need to convert both your model file and your inferencing code to use the Tensorflow

framework. Earlier in the chapter we discussed using Optimum to convert HuggingFace models to ONNX, and you can use this same library to convert our models to a .tflite, which is a compiled tensorflow model format. This format will perform well on edge devices even without a TPU, twofold with TPU acceleration.

Alternatively, if that all seems like a hassle, buying both a single board and an accelerator—because we all know the reason you bought a single board was to avoid buying two things to begin with—there are single boards that come with accelerators. Nvidia for example, has their own single board with a GPU and CUDA called Jetson. With a Jetson or Jetson-like computer that uses CUDA, we don't have to use Tensorflow either, so major plus. ExecuTorch is the PyTorch offering for inferencing on edge devices.

Another edge device worth considering is that one in your pocket, that's right your phone. Starting with iPhone X the A11 chip came with the Apple Neural Engine accelerator. For Android, Google started offering an accelerator since their Pixel 6 phone with the Tensor chipset. Developing an iOS or Android app will be very different from working with a single board that largely run versions of Linux, so we won't be discussing it in this book, but it's worth considering.

Outside of hardware, there are also several libraries and frameworks that are very cool, fast and make edge development easier. Llama.cpp, for example, is a C++ framework that allows for taking (almost) any huggingface model, converting it to the GGUF format—a format created by the llama.cpp team to store the model in a quantized fashion that makes it readily available to run on a CPU—it offers fast loading and inference on any device. Popular models like Llama, Mistral, Falcon, and even non-text models like Whisper are supported by llama.cpp at this point. It also supports Langchain integration for everyone who uses any of the Langchain ecosystem. Other libraries like GPTQ are focused more on performance than accessibility and are slightly harder to use, but can result in boosts where it counts, especially if you'd like to end up inferencing on an Android phone or something similar. We will be exploring some of these libraries in much more detail later in the book.

We've gone over a lot in this chapter and I hope you feel more confident in tackling deploying your very own LLM service. In the next chapter we will



be discussing how to better take advantage of your service by building an application around it. We'll be diving deep into prompt engineering, agents, and front-end tooling.

## 6.5 Summary

- Always compile your LLMs before putting them into production as it improves efficiency, resource utilization, and cost
- LLM APIs should implement batching, rate limiters, access keys, and streaming
- Retrieval-Augmented Generation is a simple and effective way to give your LLM context when generating content because they are easy to create and use.
- LLM inference service libraries like vLLM, HuggingFace's TGI, or OpenLLM make deploying easy but may not have the features you are looking for since they are so new.
- Kubernetes is a tool that simplifies infrastructure by providing tooling like Autoscaling and Rolling updates
  - Autoscaling is essential to improve reliability and cut costs by increasing or decreasing replicas based on utilization.
  - Rolling updates gradually implement updates to reduce downtime and maximize agility.
- Kubernetes doesn't support GPU metrics out of the box but by utilizing tools like DCGM, Prometheus, and KEDA you can resolve this issue.
- Seldon is a tool that improves deploying ML models and can be used to implement inference graphs.
- LLMs introduce some production challenges like:
  - Model updates - when your model drifts first look to your prompts and RAG systems before attempting re-finetuning.
  - Poor latency is difficult to resolve but tools to help include gRPC, GPU optimization, and caching.
  - Resource management and just acquiring GPUs can be difficult but tools like Skypilot can help.
- Edge development, while hardware limited, is the new frontier of LLM serving, and hardware like the Jetson Nano or Coral TPU is available to help.

- [1] You can find the latest version at <https://catalog.ngc.nvidia.com/orgs/nvidia/containers/pytorch/tags>
- [2] All fusing libraries are essentially built off the fuse library, worth checking out <https://github.com/libfuse/libfuse>
- [3] Learn more about the algorithm at <https://kubernetes.io/docs/tasks/run-application/horizontal-pod-autoscale/#algorithm-details>.
- [4] You can learn more about kubernetes logging architectures in their docs here <https://kubernetes.io/docs/concepts/cluster-administration/logging/>.
- [5] You can find the project here: <https://github.com/skypilot-org/skypilot>
- [6] J. Su, D. V. Vargas, and K. Sakurai, “One Pixel Attack for Fooling Deep Neural Networks,” IEEE Transactions on Evolutionary Computation, vol. 23, no. 5, pp. 828–841, Oct. 2019, doi: <https://doi.org/10.1109/tevc.2019.2890858>.
- [7] “This new data poisoning tool lets artists fight back against generative AI,” MIT Technology Review. <https://www.technologyreview.com/2023/10/23/1082189/data-poisoning-artists-fight-generative-ai>.