

Small & Specific Language Modeling: Before Large Language Models to Model Distillation

Tony DiRubbo
School of Information Studies, Syracuse University
IST 664: Natural Language Processing
Benjamin Nichols
May 2nd 2025

Abstract

Language modeling uses text representation to allow machine learning models to predict words given a context. Many different machine learning model types have been used to predict language patterns with the most prominent models being deep neural networks. In 2017 Google released the transformer model. From the transformer model came the BERT architecture which became a strong backbone for small and specific research models. The rise of large language models led to techniques such as model/knowledge distillation allowing for smaller student models to be trained from these teacher large language models. Open source architecture and platforms has made the training of these small and specific language models more accessible; however, a requirement of commercial computing power is still needed. Attached is an example of how to develop a small purpose language model with the purpose of identifying symptoms of Tarlov Cysts.

1. Small and Specific Language Modeling Before Large Language Models

Basics of Language Modeling - Definition and Development of the Transformer.

At the simplest level, language modeling refers to the ability of a machine learning model to predict missing or upcoming words in a sequence of text. While this concept has existed for multiple decades, the emergence and rise in popularity of large language models (LLMs) in recent years has significantly broadened its applications. Previously, language modeling was mainly seen in niche areas like customer service bots, educational tools for language learners, contract and policy simplification tools, and even healthcare-related applications such as clinical research summarization and symptom checking. However, the massive leap in performance of LLMs has pushed these technologies into more prominent roles across various industries. Large language models have pushed into the mainstream with the rise of ChatGPT and other platforms for these LLMs; however, multiple technologies have allowed for smaller language models for a specific purpose to exist for quite some time.

A major concept pertaining to language modeling is text representation, which refers to the process of converting natural human language into a numerical form that a machine can process. Koroteev (2021) explains that this representation can be viewed as a set of rules that translate text into machine-readable data. In earlier stages, this was achieved through categorical or one-hot encoding - where each word was represented by a sparse matrix filled mostly with zeros, except for one position indicating its index in a dictionary. While simple and computationally light, this method lacked the ability to reflect the semantic meaning of words, resulting in a representation that was both inefficient and limited in understanding contextual relationships between terms.

Recognizing these limitations, researchers from institutions such as Harvard, Yale, and Stanford developed more sophisticated techniques for converting words into numerical vectors. These methods aim to encode not just the identity of a word, but also its semantic and emotional connotation, allowing language models to better grasp the meaning and tone of the text. Computers can't understand the definition of "pizza", but the word vectorization allows a computer to quantify and understand the definition.

```

TargetWord = "Pizza"
doc = nlp(TargetWord)
Vector = doc.vector #converts the target word to a vector
print(Vector) #prints the vector

```

```

[ 0.0068727 -0.21634  0.27831 -0.26192  0.22884  0.89332
 0.4131  0.27377  0.22652  1.5041 -0.58059  0.56083
-0.18432  0.27738 -0.10709 -0.13519  0.023817  1.1765
-0.12659  0.043173  0.23242 -0.63213  0.40228 -0.20605
 0.46381 -0.12991 -0.68031 -0.010371  0.50033 -0.32266
 0.24053  0.40178  0.12051 -0.13791  0.40821  0.54735
-0.25946  0.020254  0.21249  0.91965 -0.21202  0.66568
 0.25879 -0.36124 -0.10977  0.87492 -0.089425  0.39184
-0.32589 -0.22331 -0.17504  0.074762  0.45271  0.085476
-0.079526 -0.23986 -0.010322  0.089974  0.29794  0.26672
-0.044288 -0.082716  0.20801  0.38404  0.15281 -1.1292
-0.094527  0.16901 -0.018155  0.31023 -0.095716  0.32587
-0.2225 -0.040376 -0.52201 -0.040547 -0.2473  0.059596
 0.31592  0.48751  0.14681 -0.29337  0.61309 -0.7844
-0.16297  0.042847  0.90914  0.70536 -0.44725 -0.3035
-0.26998 -0.32488  0.10539 -0.24494 -0.023413  0.51872
-0.0060798 -0.039611  0.28618  0.17071 -0.661 -0.1303
 0.59381  0.33792 -0.016678 -0.31536  0.92849 -0.19661
 0.14412  0.141  0.095604 -0.65534 -0.66278 -0.068806
 0.57471 -0.34244  0.30524  0.070219 -0.31053  0.25418
 0.16362  0.48417  0.15889  0.20571 -0.24816 -0.52146
 0.85366  0.029624 -0.20695  0.68848 -0.19801 -0.55261
 0.25334  0.23374 -0.47797 -0.58102 -0.30506 -0.24182
-0.089947  0.020074 -2.442 -0.38229  0.38005  0.45891
-0.10147 -0.46439 -0.47909  0.48057  0.66937  0.16773
 0.23094 -0.0037971  0.11692  0.19027  0.22866 -0.10451
 0.16913  0.1929  0.21792  0.3183 -0.69639 -0.039663
-0.034875 -0.34398 -0.033303 -0.44731  0.39323  0.28786
 0.41256 -0.042063  0.053043  0.032974  0.13665 -0.47123
-0.59784 -0.15469 -0.043701 -0.28768 -0.5192 -0.81342
 0.28083 -0.20795 -0.0063995 -0.2165 -0.18462 -0.12112
 0.16446  0.1074 -0.21256 -0.079728 -0.26936  0.1213
-0.41473 -0.3929  0.11391  0.017356  0.6225  0.39374
 0.18043  0.06208 -0.048457 -0.13303 -0.28215  0.23984
 0.19951  0.079811 -0.24321  0.52115  0.37684 -0.16641
 0.18813  0.38081 -0.068178  0.17925  0.26455 -0.19848
 0.23375 -0.075531 -0.26779  0.20576 -0.14164  0.18145
-0.27216 -0.27441  0.07846 -0.15033 -1.008  0.010385
-0.42617 -0.2697 -0.1037  0.23114  0.31361 -0.64973
 0.081318 -0.20336 -0.21733 -0.47194 -0.22844  0.3207
 0.47024  0.19194 -0.38945 -0.53102  0.30868  0.4516
 0.35045 -0.41879 -0.6142 -0.25174  0.1445 -0.37202
-0.29062 -0.24497  0.11433  0.51326 -0.10476  0.070229
-0.33251 -0.26181 -0.22115 -0.098211 -0.8274 -0.73995
-0.2146  0.61523 -0.36608 -0.043749 -1.0097  0.19351
-0.54799  0.28998 -0.07631  0.075785 -0.63011  0.084629
-0.021395 -0.060536  0.48363  0.16488  0.42662 -0.1786
 0.14382  0.3153  0.017293  0.6706  0.49765  0.34787
-0.62071  0.20983  0.71446  0.071448 -0.14597  0.50002
-0.42449 -0.3517  0.10134 -0.5875  0.11043  0.47559 ]

```

Figure 1.1: An example of vector text representation with the word "pizza".

Today's language models are built using deep learning neural networks, which represent a major shift in how natural language processing (NLP) tasks are approached. Unlike traditional machine learning approaches that relied on manual feature engineering - such as Support Vector Machines or logistic regression - deep learning models automatically learn latent features from data. According to Min (2023), these latent features form internal representations that are shared across different NLP tasks. This allows the model to generalize more effectively, especially when training data is limited. Before deep learning became dominant, simpler statistical approaches like Naive Bayes were more common. The development of template-based learning methods allowed researchers to align NLP tasks more closely with the structure of language modeling. As Min (2023) outlines, this approach leverages prompt templates - structured phrases with blank spaces - to convert tasks like classification or summarization into formats more compatible with language models. The model then fills in the blanks, effectively performing the task without requiring extensive fine-tuning or large datasets.

The design of neural networks made strong advancements after 2017 with the release of Google's groundbreaking paper, *Attention is All You Need*. This research introduced the transformer model, which replaced traditional sequence-processing methods like recurrent neural networks (RNNs) with a new mechanism based on self-attention. At the heart of this architecture lies the attention function, which, as Vaswani et al. (2017) describe, takes a query and a set of key-value pairs and produces an output by calculating weighted averages. The weights are determined by how well each key matches the query, allowing the model to dynamically focus on relevant words in the context. Methods such as Naive Bayes struggled to interpret such context, thus why the method was labeled as *Naive*. These transformer and neural networks developments allow for context which is essential for language models of today.

Unlike previous models that processed text sequentially, transformers can examine all parts of a sentence simultaneously. This parallel processing enables them to capture nuanced relationships between words, regardless of their position in the text. Vaswani et al. (2017) also introduced the concept of multi-head attention, which performs attention operations multiple times in parallel using different linear projections. This lets the model train and analyze different parts of the input text simultaneously and from various perspectives. The outputs from these multiple heads are then combined to form a richer and more comprehensive understanding of the input. Parallel computing has become a major component of the LLM arms race during this decade. Another innovation of the transformer is the encoder-decoder structure, particularly useful for tasks like translation or summarization. As described by Vaswani et al. (2017), encoder-decoder attention layers allow the decoder, the structure of the transformer which generates output, to access all encoded information from the input. This provides flexibility and accuracy in producing responses that are closely tied to the input sequence.

Finally, once attention scores are calculated, they are passed through a SoftMax function. This step converts the raw scores into a probability distribution, emphasizing the most relevant information while minimizing less important signals. As Sanderson (2024) notes, this normalization ensures that output values fall within a valid range - typically with the highest probabilities nearing 1 and the lowest approaching 0 - making them interpretable and actionable by the model.

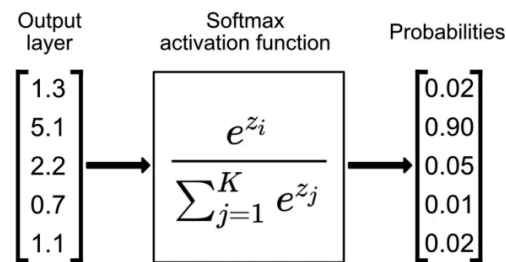


Figure 1.2: The SoftMax function takes in an output layer derived from the word vector and creates a distribution of probabilities for the next word.

From Google’s Transformer Models Comes Google’s BERT

In 2018, Google introduced BERT - Bidirectional Encoder Representations from Transformers. This would mark a significant advancement in the evolution of NLP, derived from Google’s own transformer model, and a base structure for many small language models. Developed by researchers at Google AI Language under the leadership of Jacob Devlin, BERT was designed to have detailed pre-train bidirectional representations of text, enabling more nuanced understanding for downstream machine learning tasks. One of its strengths lies in its simplicity of implementation - it requires only a minimal addition to existing neural architectures, yet it significantly outperforms previous models across a range of NLP tasks (Koroteev, 2021). Since its release, BERT has become a foundational model for developing more specialized NLP tools. As Gu (2021) points out, BERT is now “a standard building block” for constructing task-specific models across many applications. Its versatility stems from the fact that it can be adapted for various functions with minimal architecture changes.

BERT was trained on a massive corpus that includes Wikipedia and Google Books, allowing it to learn from a diverse and extensive range of text. This is a wide base plate for researchers wanting to add their own data to their own BERT models. The model operates on

both the sentence level and the token level, which enables it to capture broader contextual meanings and fine-grained word-level distinctions. BERT's training was based on two key techniques: Masked Language Modeling (MLM) and Next Sentence Prediction (NSP). MLM works by randomly masking certain words in a sentence and training the model to predict them using the surrounding context. NSP, on the other hand, involves predicting whether one sentence logically follows another, helping the model understand relationships between sentence pairs (Google Cloud Tech, 2023). Today's popular LLMs are known by the mass media for their prowess in next sentence prediction.

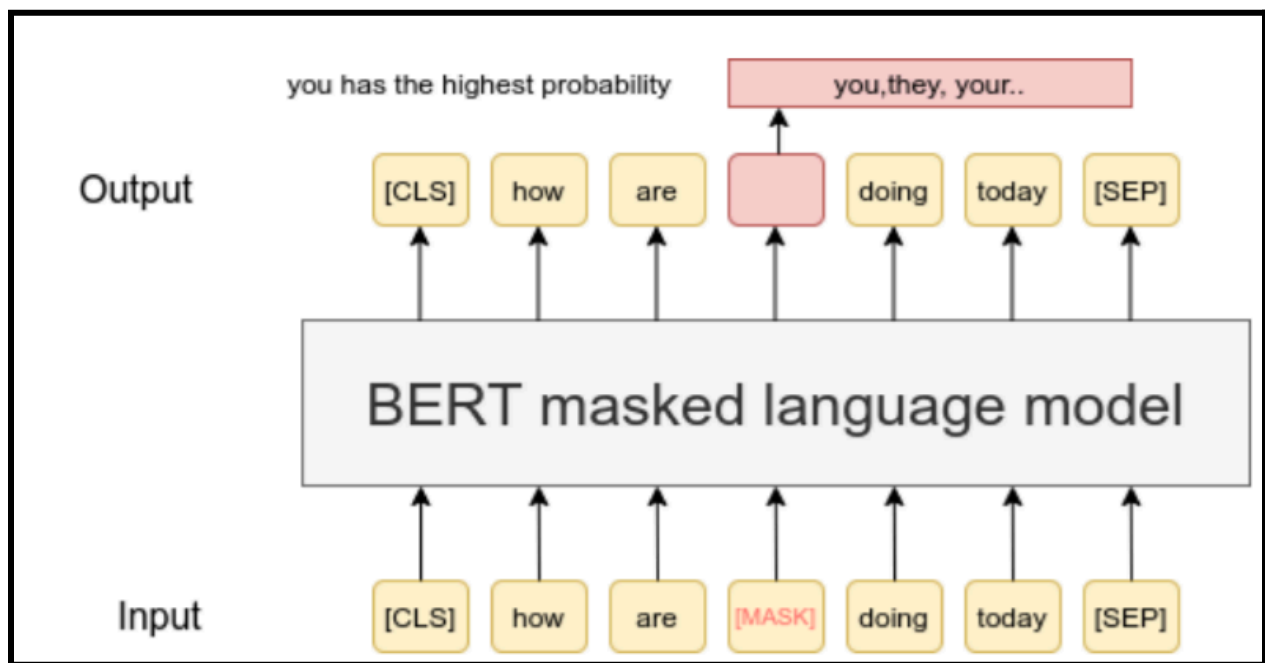


Figure 1.3: From Google's BERT resources how masked language modeling (MLM) works with the BERT ecosystem.

Structurally, BERT is built using the encoder component of the Transformer architecture, which allows it to learn complex contextual relationships through stacked encoder layers. Unlike autoregressive models that process text in one direction, BERT leverages masked self-attention mechanisms, enabling it to attend to both preceding and succeeding tokens simultaneously when

forming a representation for any given word (Min, 2023). This approach gives BERT a powerful contextual understanding that enhances performance in tasks like question answering, classification, and named entity recognition. Like traditional machine learning models, BERT training follows a two-phase process. First, the model undergoes pre-training on vast amounts of unlabeled data, learning generalized language patterns. This is followed by fine-tuning on smaller, labeled datasets tailored to specific tasks or domains. While the architecture remains consistent, fine-tuning allows for customization to suit applications, thereby increasing relevance and accuracy (Koroteev, 2021).

The prominent strength of BERT is its adaptability for domain-specific applications. Organizations and individuals can fine-tune BERT on their own data to develop tailored models for specialized tasks. In biomedical research, for example, models like BioBERT and tools like BLURB (Biomedical Language Understanding and Reasoning Benchmark) have emerged to evaluate and benchmark language models in the healthcare domain. Gu (2021) explains that BLURB was created to assess the performance of biomedical LMs, especially those pre-trained on PubMed and similar datasets. The benchmark includes a comprehensive suite of NLP tasks, such as named entity recognition, relation extraction, document classification, sentence similarity, evidence-based information extraction, and question answering. These tasks help researchers systematically evaluate how well models like BERT perform in medical contexts and guide improvements in model design for domain-specific needs.

Despite its many strengths, BERT is not without limitations. One commonly cited issue is its constrained focus in comparison to similar model archetypes such as GPT - Generative Pretrained Transformers. The GPT model is developed by OpenAI and utilized and commonly associated with the ChatGPT LLM ecosystem. While GPT models scan text in a left-to-right

manner, which restricts the contextual window to prior tokens, BERT incorporates context from both directions. As Koroteev (2021) notes, this bidirectional approach allows BERT to predict masked words by analyzing both the words that come before and after, creating a more holistic understanding of text. In contrast, GPT models are autoregressive, meaning they generate text based on the tokens that precede them, which limits their ability to fully capture bidirectional context. BERT's bidirectional nature enables it to overcome semantic limitations associated with unidirectional models, especially in cases where a word's meaning depends heavily on its surrounding context.

2. Building Small, Focused Models with Help from LLMs

Knowledge distillation is the idea of using a larger, more prominent large language model in order to train a smaller language model. This is done by leveraging the outputs or internal representations of larger, pre-trained models - often referred to as "teacher models." The goal is to transfer the teacher's knowledge to a more compact "student" model without significantly sacrificing performance. According to West (2021), this is achieved by encouraging the student model to match the predictions of the teacher model across a predefined output space. This is especially effective in classification tasks where both models generate probability distributions over a limited set of labels.

Multiple strategies have emerged for extracting knowledge from large models. Xu (2024) outlines several key techniques for knowledge elicitation. Labeling involves the teacher generating labeled outputs from inputs, while Expansion focuses on generating similar examples through in-context learning. Data curation refers to the teacher synthesizing training data based on metadata such as topics or entities. Feature extraction captures internal representations, like logits, to understand the model's inner workings. Feedback is a process where the teacher

evaluates or corrects the student’s outputs, ensuring accuracy. Finally, Self-knowledge encourages the student to evaluate its own responses for quality before refining them. These techniques offer diverse approaches for enhancing model understanding and performance.

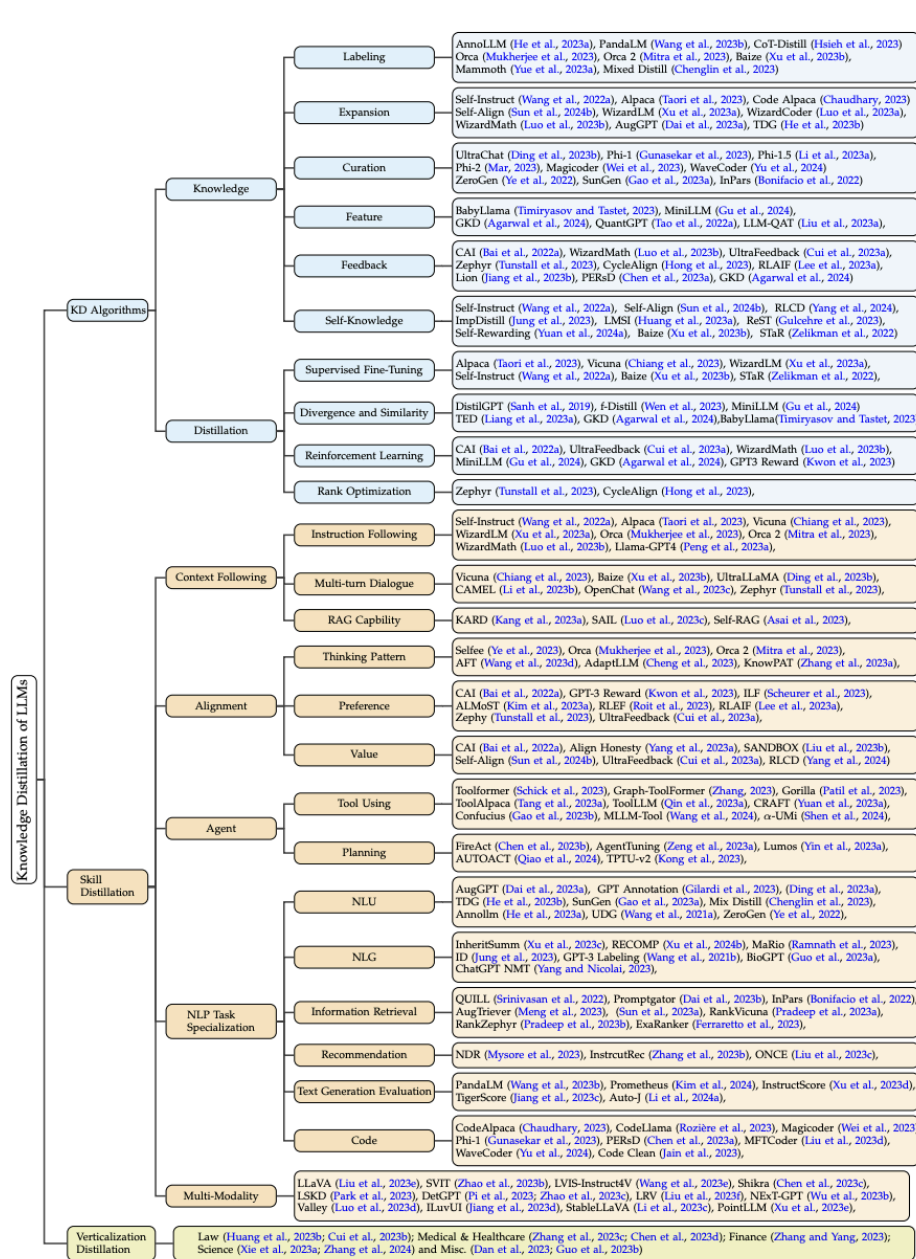


Fig. 3: Taxonomy of Knowledge Distillation of Large Language Models. The detailed taxonomy of Verticalization Distillation is shown in Figure 7.

Figure 2.1: Xu’s (2024) roadmap of possible knowledge distillation approaches

This process has become increasingly popular in the open-source language model community, where developers have access to the weights and training methods of models like Meta's Llama, DeepSeek, and Mistral. These models are often distilled into smaller, more task-specific systems, which reduces computational requirements and energy consumption. Xu (2024) notes that this "enhances the environmental sustainability of AI operations" while also making advanced AI tools more accessible to smaller organizations and independent researchers. As a result, knowledge distillation plays a key role in democratizing access to AI, leveling the playing field across industries and academia. This democratization is not just a desirable outcome but increasingly a practical necessity. Open-source LLMs allow developers to tailor models to their own needs without relying on the APIs of commercial providers, which can be expensive, restricted, or subject to terms that limit innovation. Xu (2024) also emphasizes the growing trend of using multiple teacher models to guide the training of a single student model. This multi-source distillation can improve performance by incorporating diverse patterns of reasoning and knowledge from a variety of large models, further increasing the robustness of the student model.

The controversy surrounding DeepSeek-R1, released in early 2025, highlighted some of the ethical and legal gray areas in knowledge distillation. Although the model received praise from the AI community for its performance, OpenAI accused its developers of improperly distilling knowledge from proprietary models like ChatGPT. According to Porter (2024), OpenAI alleged that large quantities of data had been exfiltrated through developer accounts tied to DeepSeek, raising concerns about potential intellectual property violations. Microsoft's security team reportedly identified suspicious activity linked to this alleged data scraping. Some prompts even led DeepSeek to refer to itself as ChatGPT, but this could be explained by broader

training on internet data, where "ChatGPT" has become synonymous with AI assistants, much like "Kleenex" is often used generically for tissues. While this naming confusion might reflect natural data bias, OpenAI maintains that its terms of service prohibit using its API outputs to train rival commercial models (Porter, 2024). OpenAI permits developers to integrate its models into applications but explicitly forbids using the outputs to train new competing systems.

As the ecosystem matures, symbolic knowledge distillation has emerged as a new paradigm, especially in the generation of knowledge graphs for commonsense reasoning. West (2021) introduces this framework as a method for automatically generating high-quality symbolic knowledge by combining outputs from various models with novel algorithms. While promising, this method - and knowledge distillation in general - can be expensive and resource-intensive. West also notes that the machine-to-corpus-to-machine pipeline, where a model generates new training data that is then re-ingested, may offer better long-term efficiency over purely human-authored datasets.

A prime example of how knowledge distillation addresses scalability is with BERT. As a large pre-trained model with millions of parameters - 110 million for BERT-Base and 330 million for BERT-Large - it has proven effective across numerous NLP benchmarks, including GLUE and SQuAD (Sun, 2019). However, BERT's size can hinder its real-world deployment, especially in environments with limited computing power. Distillation allows BERT to be compressed into smaller, shallower models that maintain high performance while consuming fewer resources. This technique has since been applied to other models like XLNet (Yang et al., 2019) and RoBERTa (Liu et al., 2019), enabling more efficient alternatives for a variety of applications. Model distillation has become a new pathway for small purpose and specific

language models to be developed. Developers can now take the relevant points of an LLM for the required task, and fine tune it for the business objective using modern computing.

3. Example of Knowledge Distillation in Practice

A Jupyter Notebook is provided at the end of this paper demonstrating an example of how to perform knowledge distillation. Many try and except clauses are utilized throughout the notebook to allow errors to be more easily read and recorded. This provides an added layer of transparency when training and troubleshooting the model, as it gives users feedback on what parts of the pipeline may be breaking and why. This type of error handling is especially useful when working with custom-built training loops or experimental model configurations such as what is utilized in this notebook, as debugging can otherwise be time-consuming and vague. Print() statements are also utilized to show completion of steps. The goal of the notebook is that another individual can bring their own data into the notebook and then train their own model.

The example utilizes Llama 3.2 1-billion parameter and BERT-base-uncased models. The idea is for the model to leverage the NSP capability of the Llama architecture and the MLM functionality from BERT. As mentioned in the first section, NSP has been a huge focus of LLMs and MLM was a major driver behind the development of BERT. These models were loaded directly from HuggingFace - a platform for sharing open-source language models and associated resources. The structure of the platform is similar to GitHub, however it is more based on sharing code and resources devoted to large language models and artificial intelligence tools. Llama originally uploaded its materials to GitHub before migrating over to HuggingFace. HuggingFace also provides its own API, titled transformers, which contains high-level functions for importing models, tokenizers, and datasets into a Python environment with minimal boilerplate. Note that

this transformers API is not the same entity as BERT's transformer structure previously mentioned.

PyTorch is utilized over TensorFlow and Keras due to its native compatibility with Llama models. Since Llama was developed using PyTorch, integrating and modifying its components becomes significantly easier. While TensorFlow was partially developed by Google, similar to the transformer architecture and BERT, it is easier to utilize BERT with PyTorch than vice versa.

BERT's corresponding tokenizer is used for the student model. The tokenizer determines how input text is split into tokens, effectively shaping the vocabulary that the student model learns from. Because the student model uses the BERT tokenizer, its vocabulary size and token embeddings are aligned with BERT's, rather than Llama's. This mismatch between tokenizer and model architecture is addressed in a later section through a custom training class that aligns token outputs from both teacher models with the student's output space.

The `student_config_obj` is where the architecture of the student neural network is defined. This object configures important model attributes such as tokenizer type, sequence length, hidden size, number of layers, and overall layer dimensions. For this particular model, the student architecture includes four layers with a hidden size of 256 and a maximum sequence length of 512. These values were chosen to balance computational feasibility and model capacity, given the constraints of the hardware used, this will be addressed when discussing how the model performed on two inference examples.

For loading the data which will be inputted into the student model. Functions were created which could load links as well as strings of text into the training corpus. This small model developed for the example will be trained with knowledge of Tarlov Cysts. From the American Association of Neurological Surgeons (AANS) - "Tarlov cysts are fluid-filled nerve

root cysts found most commonly at the sacral level of the spine – the vertebrae at the base of the spine. These cysts typically occur along the posterior nerve roots. Cysts can be valved or non-valved. The main feature that distinguishes Tarlov cysts from other spinal lesions is the presence of spinal nerve root fibers within the cyst wall or in the cyst cavity itself.” This is a rare medical condition and a case example of what a small language model might be trained upon. The student model would be trained upon data, symptoms, potentially images of Tarlov Cysts, and then the model could be utilized to identify Tarlov Cysts based on symptoms, images or other medical criteria. Sources for this specific example model are included in the notebook with APA formatting.

A custom class is defined to initiate the distilled student model, which is then connected to the architecture defined in `student_config_obj`. The model includes an embedding layer, several linear layers with ReLU activations aiming to imitate transformer encoder layers in a simplified form, and a final decoder layer that maps the internal hidden representations back to vocabulary-size logits for prediction tasks.

Distillation loss is calculated in a slightly non-standard way due to the presence of two teacher models. A function is created that takes in the output logits from both teacher models and the student model. KL (Kullback-Leibler) divergence loss is then computed to align the student's output distribution with those of the teachers. A temperature parameter is included to control the sharpness of the distributions during this loss computation. KL divergence is calculated separately between the student and each teacher model, and the final loss is the average of the two. This dual-distillation approach encourages the student model to learn shared knowledge from both teachers. Standard training parameters such as `batch_size`, `epochs`, and `learning_rate` are still used for tuning performance.

$$\begin{aligned}
\log \mathbf{p}_s &= \log \left(\text{softmax} \left(\frac{\mathbf{s}}{T} \right) \right) \\
\mathbf{p}_{\text{llama}} &= \text{softmax} \left(\frac{\mathbf{t}_{\text{llama}}}{T} \right) \\
\mathbf{p}_{\text{bert}} &= \text{softmax} \left(\frac{\mathbf{t}_{\text{bert}}}{T} \right) \\
\text{KL}_{\text{llama}} &= \sum_i \mathbf{p}_{\text{llama},i} (\log \mathbf{p}_{\text{llama},i} - \log \mathbf{p}_{s,i}) \\
\text{KL}_{\text{bert}} &= \sum_i \mathbf{p}_{\text{bert},i} (\log \mathbf{p}_{\text{bert},i} - \log \mathbf{p}_{s,i}) \\
\mathcal{L}_{\text{distill}} &= \frac{1}{2} (T^2 \cdot \text{KL}_{\text{llama}} + T^2 \cdot \text{KL}_{\text{bert}})
\end{aligned}$$

Figure 3.1: A mathematical representation of how the double loss function occurs.

Two inference tasks were used to test the model post-training: text generation and masked language modeling. While the student model was able to perform both tasks, the results were largely incoherent and disappointing. This is likely due to limitations in available processing power. Despite these models being open-source, effective use often still requires commercial-grade hardware. A Tesla T4 GPU was used for this notebook, running with CUDA optimization on Google's Colab servers. To make the Llama 3.2 1B model run on this environment, quantization techniques and the BitsAndBytes API had to be used. Llama 3.2 1B is itself a significantly compressed version of the full model family, and memory constraints meant that parameters such as sequence length, batch size, and epoch count had to be tightly controlled to avoid memory leaks or runtime crashes.

The main takeaway from this notebook is that model distillation - even with the use of two teacher models - is accessible and viable for creating small, purpose-built, language models. This can be an excellent option for researchers or developers who want a targeted model rather than relying on general-purpose large language models. With more powerful hardware and a

higher-parameter version of Llama, the inference results from this student model would most likely be significantly improved. This is the direction that small specific language models are heading towards. Instead of developing their own independent architectures like before utilizing neural networks and frameworks such as BERT, it is now better performance wise to train a student model based on a prominent large language model to use as a teaching base, for the pursuit of developing a smaller language model suited for a specific task.

Model Distillation: Training a Smaller Student Model from LLaMA and BERT

This accompanying notebook demonstrates how to train a smaller "student" language model using model distillation. I use two larger "teacher" models (LLaMA 3.2 1B and BERT-base-uncased) to guide the training of a simpler student model.

Symbolic knowledge distillation - a new conceptual framework towards high-quality automatic knowledge graphs for commonsense, leveraging state-of-the-art models and novel methodology.

Key Steps:

1. **Setup:** Install libraries, log in to Hugging Face, import necessary modules.
2. **Configuration:** Set up quantization, load teacher models, initialize tokenizer, define student model configuration.
3. **Dataset:** Define and load the training/evaluation dataset (using a small sample for demonstration).
4. **Student Model:** Define the architecture of the student model.
5. **Distillation Loss:** Define the loss function (KL divergence between student and teacher outputs).
6. **Training:** Set up `TrainingArguments` and a custom `DistillerTrainer` to handle the distillation process.
7. **Execute Training:** Run the training loop.
8. **Evaluation:** Evaluate the trained student model.
9. **Inference Examples:** Show how to use the trained student model for generation and MLM-like tasks.
10. **Cleanup:** Release resources.

Note: Running LLaMA models, even quantized, requires significant GPU memory. This notebook was generated in Google Colab utilizing GPU resources over the standard CPU. You will also need a Hugging Face token with access granted to the Llama 3.2 model.

```
In [1]: # Install Required Libraries
!pip install torch transformers datasets evaluate bitsandbytes accelerate sc
!pip install -U bitsandbytes
```

Requirement already satisfied: bitsandbytes in /usr/local/lib/python3.11/dist-packages (0.45.5)

Requirement already satisfied: torch<3,>=2.0 in /usr/local/lib/python3.11/dist-packages (from bitsandbytes) (2.6.0+cu124)

Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from bitsandbytes) (2.0.2)

Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.18.0)

Requirement already satisfied: typing-extensions>=4.10.0 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (4.13.1)

Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.4.2)

Requirement already satisfied: jinja2 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.1.6)

Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (2024.12.0)

Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)

Requirement already satisfied: nvidia-cuda-runtime-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)

Requirement already satisfied: nvidia-cuda-cupti-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)

Requirement already satisfied: nvidia-cudnn-cu12==9.1.0.70 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (9.1.0.70)

Requirement already satisfied: nvidia-cublas-cu12==12.4.5.8 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.5.8)

Requirement already satisfied: nvidia-cufft-cu12==11.2.1.3 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (11.2.1.3)

Requirement already satisfied: nvidia-curand-cu12==10.3.5.147 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (10.3.5.147)

Requirement already satisfied: nvidia-cusolver-cu12==11.6.1.9 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (11.6.1.9)

Requirement already satisfied: nvidia-cusparselt-cu12==12.3.1.170 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.3.1.170)

Requirement already satisfied: nvidia-cusparse-cu12==12.3.1.170 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.3.1.170)

Requirement already satisfied: nvidia-cusparselt-cu12==0.6.2 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (0.6.2)

Requirement already satisfied: nvidia-nccl-cu12==2.21.5 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (2.21.5)

Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)

Requirement already satisfied: nvidia-nvjitlink-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)

Requirement already satisfied: triton==3.2.0 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.2.0)

Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (1.13.1)

Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch<3,>=2.0->bitsandbytes) (1.3.0)

Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2->torch<3,>=2.0->bitsandbytes) (3.0.2)

In [2]: *# Log in to Hugging Face (You need to request Llama Access for each individual user)*
!huggingface-cli login

```

_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |
_ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ | _ |

```

A token is already saved on your machine. Run `huggingface-cli whoami` to get more information or `huggingface-cli logout` if you want to log out.

Setting a new token will erase the existing one.

To log in, `huggingface_hub` requires a token generated from <https://huggingface.co/settings/tokens>.

Enter your token (input will not be visible):

Add token as git credential? (Y/n) n

Token is valid (permission: fineGrained).

The token `test` has been saved to /root/.cache/huggingface/stored_tokens

Your token has been saved to /root/.cache/huggingface/token

Login successful.

The current active token is: `test`

```

In [3]: # importing necessary libraries
import os
import numpy as np
import copy
from sklearn.metrics import accuracy_score, precision_recall_fscore_support
from torch.utils.data import Dataset, DataLoader
# PyTorch is being utilized over tensorflow due to its development in Llama
import torch
import torch.nn as nn
import torch.nn.functional as F
import gc # Garbage collector for memory management

from transformers import (
    AutoTokenizer, AutoModelForMaskedLM, AutoModelForCausalLM, AutoConfig,
    TrainingArguments, Trainer, PretrainedConfig, BitsAndBytesConfig
)
from torch import bfloat16 # Potentially useful for mixed precision, though

print("Imports successful.")
# Check for GPU availability
if torch.cuda.is_available():
    print(f"GPU detected: {torch.cuda.get_device_name(0)}")
    device = torch.device("cuda")
else:
    print("No GPU detected, using CPU. Warning: Training will be very slow.")
    device = torch.device("cpu")

```

Imports successful.

GPU detected: Tesla T4

Here is the configuration for model loading:

- **Quantization:** 8-bit quantization (`BitsAndBytesConfig`) is used for the large LLaMA teacher model to reduce memory footprint.
- **Teacher Models:** Loading the pre-trained teacher models:
 - `meta-llama/Llama-3.2-1B` : A large causal language model (decoder-only). Loaded with quantization.
 - `bert-base-uncased` : A masked language model (encoder-only).
- Models are moved to the GPU if available (`.to(device)` or `device_map="auto"`).

```
In [4]: quantization_config = BitsAndBytesConfig(load_in_8bit=True) # Using 8-bit qu
print("Quantization config created (8-bit for LLaMA).")

try:
    print("Loading teacher LLaMA model with quantization...")
    # device_map="auto" helps distribute the model across available GPUs/CPU
    teacher_llama = AutoModelForCausalLM.from_pretrained(
        'meta-llama/Llama-3.2-1B',
        quantization_config=quantization_config,
        device_map="auto" # Recommended for large models
    )
    teacher_llama.eval() # Set to evaluation mode
    print(f"LLaMA model loaded with quantization. Device map: {teacher_llama
except Exception as e:
    print(f"Error loading quantized LLaMA model: {e}")
    print("Ensure you have accepted the license terms for Llama-3.2 on Huggi
    raise

try:
    print("Loading teacher BERT model...")
    teacher_bert = AutoModelForMaskedLM.from_pretrained('bert-base-uncased')
    # Move BERT to the primary device if not using device_map
    if torch.cuda.is_available() and not hasattr(teacher_bert, 'hf_device_ma
        teacher_bert.to(device) # Use the device defined earlier
    teacher_bert.eval() # Set to evaluation mode
    print(f"BERT model loaded. On device: {next(teacher_bert.parameters()).c
except Exception as e:
    print(f"Error loading BERT model: {e}")
    raise
```

Quantization config created (8-bit for LLaMA).

Loading teacher LLaMA model with quantization...

/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:

The secret `HF_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and restart your session.

You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

warnings.warn(

LLaMA model loaded with quantization. Device map: {'': 0}

Loading teacher BERT model...

Some weights of the model checkpoint at bert-base-uncased were not used when initializing BertForMaskedLM: ['bert.pooler.dense.bias', 'bert.pooler.dense.weight', 'cls.seq_relationship.bias', 'cls.seq_relationship.weight']

- This IS expected if you are initializing BertForMaskedLM from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining model).
- This IS NOT expected if you are initializing BertForMaskedLM from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

BERT model loaded. On device: cuda:0

- **Tokenizer:** BERT's corresponding Tokenizer will be utilized. This determines the vocabulary the student model will use. **Important:** Because the student uses the BERT tokenizer, its vocabulary size will match BERT's, but *not* LLaMA's. The custom trainer will need to handle this mismatch.
- **Student Config:** The architecture of the student model is defined using `PretrainedConfig`. This student is much smaller than the teachers (fewer layers, smaller hidden size, fewer attention heads).
- The `MAX_SEQ_LENGTH` is altered to limit memory usage.

```
In [5]: # Using BERT's tokenizer for the student model
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
print(f"Tokenizer loaded: {tokenizer.name_or_path}, Vocab size: {tokenizer.v

# Configuring the Student Model
MAX_SEQ_LENGTH = 512 # Reduced sequence length for efficiency
student_config_obj = PretrainedConfig(
    model_type="distilled_student", # Custom model type name
    vocab_size=tokenizer.vocab_size, # Use BERT's vocab size
    n_layer=4, # Number of layers (BERT-base has 12)
    hidden_size=256, # Hidden dimension size (BERT-base has 768)
    intermediate_size=512, # Feed-forward layer size (BERT-base has 3072)
    num_attention_heads=4, # Number of attention heads (BERT-base has 12)
    max_position_embeddings=MAX_SEQ_LENGTH, # Max sequence length
    pad_token_id=tokenizer.pad_token_id,
    # Add any other relevant config parameters if needed
)
print(f"Student model config created with {student_config_obj.n_layer} layer
```

Tokenizer loaded: bert-base-uncased, Vocab size: 30522

Student model config created with 4 layers, hidden size 256, and max sequence length 512.

Several functions are used to load and process the training data.

- `load_medical_dataset`: A simple function to load text data (here, from a list).
- `MedicalDataset`: A PyTorch `Dataset` class that tokenizes the text using the specified tokenizer and max length.

- `create_dataset` : A helper function to instantiate the dataset.

For this example, the data is a very small, hardcoded dataset, about Tarlov cysts. For replication on another scenario, you would replace the `txt_list` with loading from a file or a Hugging Face dataset.

The sources from the data are below:

American Association of Neurological Surgeons. (n.d.). Tarlov cyst. AANS. Retrieved April 11, 2025, from <https://www.aans.org/patients/conditions-treatments/tarlov-cyst/>

Cleveland Clinic. (n.d.). Tarlov cyst. Cleveland Clinic. Retrieved April 11, 2025, from <https://my.clevelandclinic.org/health/diseases/tarlov-cyst>

National Institute of Neurological Disorders and Stroke. (n.d.). Tarlov cysts. National Institutes of Health. Retrieved April 11, 2025, from <https://www.ninds.nih.gov/health-information/disorders/tarlov-cysts>

National Organization for Rare Disorders. (n.d.). Tarlov cysts. Rare Diseases. Retrieved April 11, 2025, from <https://rarediseases.org/rare-diseases/tarlov-cysts/>

The Tarlov Cyst Foundation. (n.d.). Homepage. The Tarlov Cyst Foundation. Retrieved April 11, 2025, from <https://www.tarlovcystfoundation.org/>

Wikipedia. (n.d.). Tarlov cyst. Wikipedia, the Free Encyclopedia. Retrieved April 11, 2025, from https://en.wikipedia.org/wiki/Tarlov_cyst

WebMD Editorial Contributors. (n.d.). What to know about cysts at the base of the spine. WebMD. Retrieved April 11, 2025, from <https://www.webmd.com/skin-problems-and-treatments/what-to-know-cysts-base-spine>

```
In [6]: def load_medical_dataset(file_path=None, text_list=None):
        if text_list:
            # Ensure texts are non-empty strings
            return [text.strip() for text in text_list if text and isinstance(text, str)]
        else:
            return [] # Return empty list if no text_list provided

        class MedicalDataset(Dataset):
            def __init__(self, txt_list, tokenizer, max_length=512):
                self.tokenizer=tokenizer; self.max_length=max_length
                # Filter out potential None or non-string entries just in case
                self.txt_list=[text for text in txt_list if text and isinstance(text, str)]
                self.encodings=self.tokenizer(
                    self.txt_list,
                    max_length=self.max_length,
                    padding='max_length', # Pad to max_length
                    truncation=True,      # Truncate longer sequences
                    return_tensors='pt'   # Return PyTorch tensors
```

```

    )
def __len__(self):
    return len(self.txt_list)
def __getitem__(self, idx):
    # Return a dictionary slice for the given index
    return {key: val[idx].clone().detach() for key, val in self.encoding.items()}

def create_dataset(tokenizer, max_length):
    # Using a small sample dataset for demonstration
    txt_list = [
        # Definitions & Basic Characteristics
        "Tarlov cysts are fluid-filled sacs.",
        "These cysts are found on nerve roots.",
        "They are nerve root cysts.",
        "Tarlov cysts are most commonly found in the sacral spine region.",
        "The sacral level of the spine is the most common location.",
        "They form on the nerve root sheath.",
        "Often located in the lower spine.",
        "Also known as meningeal cysts.",
        "Also known as perineural cysts.",
        "Also known as sacral nerve root cysts.",
        "They are cerebrospinal fluid-filled sacs.",
        "CSF (cerebrospinal fluid) is the typical fluid within the cysts.",
        "Located in the spinal canal.",
        "A key feature is the presence of spinal nerve root fibers.",
        "Nerve fibers can be within the cyst wall.",
        "Nerve fibers can be found in the cyst cavity itself.",
        "This presence of nerve fibers distinguishes Tarlov cysts from other types.",
        "Tarlov cysts are classified as type II innervated meningeal cysts.",
        "They typically appear at the bottom of the spine (the sacrum).",
        "Cysts form in the roots of nerves exiting the spinal cord.",
        "They affect the nerve roots of the spine.",
        "The base of the spine (sacral region) is a common site.",
        "A person may have multiple Tarlov cysts.",
        "These multiple cysts can vary in size.",
        "Most frequently located in the sacral region (S1-S5) spinal canal.",
        "Less commonly found in the cervical spine.",
        "Less commonly found in the thoracic spine.",
        "Less commonly found in the lumbar spine.",
        "Cysts are formed within the nerve-root sheath.",
        "Formation occurs specifically at the dorsal root ganglion.",
        "Tarlov cysts are pockets of fluid.",
        "They form around the nerves that make up your spinal cord.",
        "The sacrum or lower back area is the most frequent location.",
        "Tarlov cysts can be valved or nonvalved.", # From AANS link text provided
        "They sit on the nerve roots extending out from the spinal cord.",

        # Asymptomatic Nature
        "Many cases of Tarlov cysts are asymptomatic.",
        "Most individuals with Tarlov cysts do not experience symptoms.",
        "The majority of Tarlov cysts do not cause symptoms.",
        "You may not know you have one unless found incidentally.",
        "Asymptomatic cysts usually do not require treatment.",

        # Symptoms - General
        "Symptoms vary greatly among patients.",

```


"Symptom presentation depends on cyst size.",
"Symptom presentation depends on cyst location.",
"Symptoms may flare up and then subside.",
"Larger cysts have a higher likelihood of causing symptoms.",
"Symptomatic cysts can cause pain.",
"Symptoms are based on the locations of the cysts along the spine.",
"Symptoms generally follow the pathology of spinal injury.",

Symptoms - Pain

"Pain can occur in the area of the nerves affected by the cysts.",
"Pain is often felt in the buttocks.",
"Lower back pain is a common symptom.",
"Pain can spread from the lower back to the buttocks and legs.",
"Sciatica (nerve pain radiating down the leg) can be caused by these",
"Secondary Sciatica is a documented symptom.",
"Tertiary sciatica can result from secondary piriformis muscle dysfunction",
"Perineal pain (pain in the area between the genitals and anus) can",
"Coccydynia (tailbone pain) is a possible symptom.",
"Sacral radiculopathy (nerve root pain originating from the sacrum)",
"Radicular pain (nerve root irritation) is a symptom.",
"Headaches can be associated with Tarlov cysts.",
"Chronic headaches have been reported in affected individuals.",
"Pressure behind the eyes can occur.",
"Vulvar pain has been reported.",
"Testicular pain has been reported.",
"Rectal pain has been reported.",
"Pelvic pain has been reported.",
"Abdominal pain has been reported.",
"Severe lower abdominal pain is a noted symptom.",
"Pain in the genitals can occur.",
"Leg cramps are a possible symptom.",
"Discomfort or pain when sitting may occur.",
"Discomfort or pain when standing may occur.",
"Difficulty sitting for prolonged periods is common.",
"Difficulty standing for extended periods can occur.",
"Neurogenic claudication (pain caused by walking limited distance) is",
"Pain can occur in the neck if cysts are in the cervical spine.",
"Pain can occur in the shoulders if cysts are higher up.",
"Pain can occur in the chest if cysts are in the thoracic spine.",
"Pain can occur in the arms if cysts are higher up.",
"Pain can occur in the legs from cysts higher up the spine.",

Symptoms - Sensory & Neurological

"Numbness is a potential symptom.",
"Loss of sensation on the skin can occur.",
"Altered sensation (paresthesia) may be experienced.",
"Paresthesias (strange sensations like tingling) in legs and feet can",
"Hypesthesia (decreased sensation) is possible.",
"A shocking sensation down the legs may occur.",
"A burning sensation down the legs may occur.",
"Loss of reflexes can be a sign.",
"Diminished reflexes may be observed.",
"Changes in sensation over the buttocks may occur.",
"Changes in sensation over the perineal area may occur.",
"Changes in sensation over the lower extremity may occur.",
"Vaginal paresthesia can occur.",

"Penile paresthesia can occur.",
"Persistent genital arousal disorder (PGAD) has been linked.",
"PGAD involves unwanted, unrelenting genital sensory awareness, itch",
"PGAD symptoms can persist for days, months, or even years.",
"Blurred vision has been reported.",
"Dizziness is a potential symptom.",
"Tingling or prickling sensation on the skin is reported by some.",
"Paresthesias can occur in the neck with cervical cysts.",
"Paresthesias can occur in the shoulders with higher cysts.",
"Paresthesias can occur in the chest with thoracic cysts.",
"Paresthesias can occur in the arms with higher cysts.",

Symptoms - Motor Function

"Weakness of muscles can occur.",
"Leg weakness is reported, though sometimes cited as rare.",
"Motor disorders in the lower limbs can occur.",
"Motor disorders in the genital area are possible.",
"Motor disorders in the perineal area are possible.",
"Motor disorders in the lumbosacral areas are possible.",
"Dragging of the foot when walking (foot drop) can occur.",
"Foot drop can be due to weakness of muscles in the ankles and feet.",
"Spasticity (muscle stiffness) may occur.",
"Hypertonia (increased muscle tone) may occur.",
"Muscular dysfunction or weakness can be present.",
"Secondary piriformis muscle dysfunction can occur.",

Symptoms - Bowel & Bladder

"Bowel or bladder dysfunction is a significant potential symptom.",
"Changes in bowel function can occur.",
"Constipation is a possible bowel symptom.",
"Bowel incontinence may occur.",
"Intestinal motility disorders, like constipation, are noted.",
"Changes in bladder function may occur.",
"Increased urinary frequency is possible.",
"Urinary incontinence is possible.",
"Loss of bladder control can happen.",
"Loss of bowel regulation can happen.",
"Neurogenic bladder (malfunctioning bladder due to nerve issues) can",
"Dysuria (painful urination) has been reported.",

Symptoms - Sexual Function

"Changes in sexual function can be a symptom.",
"Sexual dysfunction is noted in multiple sources.",
"Impotence can occur.",
"Retrograde ejaculation has been associated with Tarlov cysts.",
"Vaginismus (involuntary vaginal muscle spasms) can occur.",

Symptoms - Other Syndromes

"Cauda equina syndrome (a serious condition requiring urgent attention)",

Causes & Formation Theories

"The exact cause remains unknown.",
"Experts are not precisely sure what causes them.",
"Several hypotheses exist regarding their formation.",
"One theory involves trauma.",
"Trauma may cause cerebrospinal fluid (CSF) leakage into the area.",

"A traumatic spinal injury could lead to CSF leak.",
"Inflammation within the nerve root sheath is a proposed mechanism."
"Inflammation of the protective cover (sheath) around a nerve root n
"Congenital factors might play a role.",
"Some suggest a developmental origin.",
"An abnormal congenital connection between the subarachnoid space ar
"Hemorrhagic infiltration (bleeding into) of spinal tissue is propos
"Inflammation within nerve root cysts followed by fluid inoculation
"Arachnoidal proliferation along/around the exiting sacral nerve roc
"Breakage of venous drainage in perineuria/epineurium is suggested."
"This venous drainage issue might be secondary to hemosiderin deposi
"Accidents or falls involving the tailbone area might trigger sympto
"Previously undiagnosed Tarlov cysts can flare up after trauma.",
"Shock or trauma to the spine can cause CSF in cysts to build up.",
"Exertion might cause CSF build-up in cysts.",
"Physical strain, like heavy lifting, is suggested as a possible cor
"Automobile accidents are mentioned as potential triggers.",
"Childbirth is suggested as a potential contributing factor.",
"An increase in pressure in or on the cysts can worsen symptoms.",
"Increased pressure may cause nerve damage.",
"Normal fluctuations in CSF pressure might enlarge cysts.",
"Enlargement due to CSF pressure changes can lead to symptoms.",
"Some experts believe certain individuals are born with a higher ris

Diagnosis

"Diagnosis can be difficult.",
"Difficulty arises from limited knowledge about the condition.",
"Symptoms often mimic other disorders, complicating diagnosis.",
"Diagnosis may be suspected based on a thorough clinical evaluation.
"A neurological evaluation is typically performed.",
"A physical exam is part of the diagnostic process.",
"Confirmation is usually achieved through imaging.",
"MRI (Magnetic Resonance Imaging) scans are commonly used.",
"MRI helps visualize the cysts.",
"MRI shows the relationship of cysts to nerve roots.",
"MRI is considered the imaging study of choice for identification.",
"CT (Computed Tomography) scans can also detect Tarlov cysts.",
"CT scans are another examination method used.",
"Both CT and MRI are good for detecting extradural spinal masses lik
"Follow-up radiological studies may be recommended.",
"CT myelography is sometimes used for follow-up.",
"A myelogram uses contrast dye and X-rays.",
"Myelography helps visualize the subarachnoid space and CSF flow.",
"If bladder problems exist, standard urological tests are used.",
"Urological tests help determine if a neurogenic bladder is present.
"A Tarlov cyst might be discovered incidentally during an MRI for ot

Treatment - General & Conservative

"Treatment depends on whether the cyst causes symptoms.",
"Asymptomatic Tarlov cysts typically require no treatment.",
"If cysts don't cause symptoms, treatment is usually unnecessary.",
"Observation may be appropriate for asymptomatic or mildly symptomat
"Treatment for symptomatic cysts ranges from conservative to surgica
"Treatment options vary for symptomatic cysts.",
"Conservative management is often tried first.",
"Medications may be used for pain management.",

"Pain relief medication is a common treatment aspect.",
"Steroid injections might be administered.",
"Injections aim to manage pain or inflammation.",
"Physical therapy may be recommended.",
"Physical therapy aims to reduce symptoms and improve function.",
"Transcutaneous electrical nerve stimulation (TENS) is a non-invasive.",
"TENS uses electrical currents to relieve pain.",
"Lifestyle modifications might be suggested (e.g., avoiding certain

Treatment - Interventional (Non-Surgical)

"Interventional treatment is considered the only way to potentially",
"This is because cysts often refill after simple drainage.",
"Lumbar drainage of cerebrospinal fluid (CSF) is a non-surgical therapy.",
"CT scanning-guided cyst aspiration is performed.",
"Aspiration involves draining the fluid from the cyst.",
"Draining the fluid aims to relieve pressure.",
"A newer technique involves CSF removal followed by filling the space.",
"Fibrin glue injection is used in this newer technique.",
"Fibrin glue therapy aims to prevent the cyst from refilling.",
"A combination of draining and filling with fibrin glue is used.",

Treatment - Surgical

"Surgical intervention may be required for some cysts.",
"Surgery is considered for very large cysts.",
"Surgery may be needed if conservative/interventional treatments fail.",
"Surgery aims to decompress affected nerve roots.",
"Several surgical approaches exist.",
"Tarlov cyst surgery typically involves exposing the spine region.",
"Laminectomy (removal of part of a vertebra) may be part of the approach.",
"Laminectomy with wrapping of the cyst is one surgical technique.",
"During surgery, the cyst is often opened.",
"The fluid inside the cyst is drained.",
"Steps are taken to prevent the fluid from returning.",
"The cyst may be occluded (blocked off) with fibrin glue.",
"Other matter might be used to occlude the cyst.",
"One technique involves collapsing the cyst wall.",
"The collapsed cyst wall may then be sutured.",
"The cyst cavity might be packed with another substance after draining.",
"Another surgical procedure uses a muscle flap.",
"The muscle flap is used to fill the drained cyst space.",
"Using a muscle flap can help prevent recurrence.",
"Microfenestration (creating small openings) is a technique.",
"Surgical sleeving of the cysts is another technique.",
"Microfenestration and sleeving aim to reduce CSF accumulation.",
"These techniques also aim to decrease compression of the spine and",
"Surgical removal of the cyst is an option.", *# Note: Often complex*
"Surgical outcomes can vary.",

Prognosis & Other Info

"Pressure on nerves from the cysts can cause pain.",
"Pressure on nerves can also lead to bone deterioration.",
"The Tarlov Cyst Disease Foundation exists to help patients.",
"The foundation promotes research and education.",
"The foundation advocates for patients with symptomatic Tarlov cysts.",
"Symptoms can be life-altering for some individuals."

```

print(f"Loading dataset with {len(txt_list)} samples.")
loaded_texts = load_medical_dataset(text_list=txt_list)
if loaded_texts:
    return MedicalDataset(loaded_texts, tokenizer, max_length=max_length)
else:
    print("Warning: No texts were loaded.")
    return None

print("Dataset functions defined.")

```

Dataset functions defined.

This next chunk defines the `DistilledModel` class. It's a simple sequential model:

- An embedding layer.
- A series of linear layers with ReLU activations (simulating transformer layers in a very basic way).
- A final linear layer (decoder) to map back to vocabulary size for predictions.

It also includes a basic `generate` method for text generation, similar to causal LMs, although its structure is not truly autoregressive like a GPT or LLaMA. It performs simple multinomial sampling. This is so there can be a CLM and MLM example later on.

```

In [7]: class DistilledModel(nn.Module):
    def __init__(self, config: PretrainedConfig):
        super().__init__()
        self.config = config
        self.embeddings = nn.Embedding(config.vocab_size, config.hidden_size)
        self.layers = nn.ModuleList([nn.Linear(config.hidden_size, config.hidden_size) for _ in range(config.num_hidden_layers)])
        self.activation = nn.ReLU()
        self.decoder = nn.Linear(config.hidden_size, config.vocab_size)

    def forward(self, input_ids, attention_mask=None, **kwargs):
        hidden_states = self.embeddings(input_ids)
        for layer in self.layers:
            hidden_states = self.activation(layer(hidden_states))
        return self.decoder(hidden_states)

    def generate(self, input_ids, max_length=50, temperature=1.0, no_repeat_device = next(self.parameters()).device # Get device from one of the parameters
        generated_ids = input_ids.to(device) # Ensure input is on the same device as the model

        for _ in range(max_length - len(input_ids[0])): # Generate for `max_length` tokens
            # Get the logits for the next token
            logits = self.forward(generated_ids) # Assumes forward returns logits to the last token
            logits = logits[:, -1, :] # Get logits for the last token

            # Apply temperature and softmax to get probabilities
            logits = logits / temperature
            probabilities = F.softmax(logits, dim=-1)

            # Sample or pick the top token (e.g., greedy or top-k sampling)

```

```

        next_token = torch.multinomial(probabilities, 1) # Using multinomial
        generated_ids = torch.cat([generated_ids, next_token], dim=1)

    return generated_ids

student_model = DistilledModel(student_config_obj)
if torch.cuda.is_available():
    student_model.to(torch.device("cuda")) # Manually move student model

print(f"Student model instantiated. On device: {next(student_model.parameters())}")

```

Student model instantiated. On device: cuda:0

- **TrainingArguments:** Configure training parameters like epochs, batch size, learning rate (defaults used here), logging, saving strategy, gradient accumulation, and mixed-precision training (fp16).
`gradient_accumulation_steps` helps simulate a larger batch size without using more memory per step.
- **DistillerTrainer:** Defines a custom Trainer class that overrides `compute_loss`. This is crucial for:
 - Getting logits from both teacher models within the training loop (using `torch.no_grad()` for teachers).
 - **Handling Vocabulary Mismatch:** Checking if the student's vocab size matches each teacher's vocab size before attempting to calculate the respective KL divergence loss. This prevents errors if, for instance, LLaMA's vocab is different from BERT's (which it is).
 - Calculating the combined distillation loss using KL divergence.
 - Includes extensive print statements and memory management (`gc.collect`, `torch.cuda.empty_cache`).

```

In [8]: # --- Training Arguments ---
training_args = TrainingArguments(
    output_dir='./results_distil_memopt',
    num_train_epochs=5,
    per_device_train_batch_size=1,
    gradient_accumulation_steps=8,
    warmup_steps=10,
    weight_decay=0.01,
    logging_dir='./logs_distil_memopt',
    logging_steps=5,
    save_strategy="epoch",
    save_total_limit=1,
    fp16=torch.cuda.is_available(),
    report_to="tensorboard",
    dataloader_num_workers=2,
)
print("TrainingArguments configured.")

```

```

class DistillerTrainer(Trainer):
    """Custom Trainer, skips loss for mismatched vocab sizes."""
    def __init__(self, *args, teacher_llama, teacher_bert, **kwargs):
        super().__init__(*args, **kwargs)
        self.teacher_llama = teacher_llama
        self.teacher_bert = teacher_bert
        # It's safer to check the config of the model passed to the trainer
        self.student_vocab_size = self.model.config.vocab_size # self.model
        self.llama_vocab_size = teacher_llama.config.vocab_size
        self.bert_vocab_size = teacher_bert.config.vocab_size
        print("DistillerTrainer initialized.")
        print(f"Vocab Sizes - Student: {self.student_vocab_size}, LLaMA: {se
        # Print the primary device the Trainer will use
        print(f"Trainer primary device: {self.args.device}")

    def compute_loss(self, model, inputs, num_items_in_batch=None, return_o
        """Computes distillation loss, skipping components with mismatched v
        inputs = {k: v.to(self.args.device) for k, v in inputs.items()}

        # Get student logits (model call should handle internal device place
        student_outputs = model(**inputs) # This calls DistilledModel.forward
        student_logits = student_outputs # if hasattr(student_outputs, 'logi

        loss_components = []
        llama_logits_for_loss = None # Initialization for clarity
        bert_logits_for_loss = None # Initialization for clarity

        #Get Teacher Logits & Calculate Loss (conditional)
        with torch.no_grad():
            # LLaMa
            if self.llama_vocab_size == self.student_vocab_size:
                try:
                    # LLaMA call
                    teacher_llama_outputs = self.teacher_llama(**inputs)
                    teacher_llama_logits = teacher_llama_outputs.logits
                    # Move back to primary device
                    teacher_llama_logits = teacher_llama_logits.to(self.args
                    print("Got LLaMA logits.")
                    llama_logits_for_loss = teacher_llama_logits # Keep refe
                    del teacher_llama_outputs # Free memory
                except Exception as e:
                    print(f"Error getting LLaMA logits: {e}")
                    llama_logits_for_loss = None # Indicate failure
            else:
                print(f"Skipping LLaMA: Vocab size mismatch ({self.llama_voc
                llama_logits_for_loss = None # Indicate failure

            # BERT
            if self.bert_vocab_size == self.student_vocab_size:
                try:
                    # BERT device check/move
                    if hasattr(self.teacher_bert, 'device') and self.teacher
                        self.teacher_bert.to(self.args.device) # Move BERT i
                    # BERT call
                    teacher_bert_outputs = self.teacher_bert(**inputs)

```



```

        teacher_bert_logits = teacher_bert_outputs.logits
        teacher_bert_logits = teacher_bert_logits.to(self.args.device)
        print("Got BERT logits.")
        bert_logits_for_loss = teacher_bert_logits # Keep reference
        del teacher_bert_outputs
    except Exception as e:
        print(f"Error getting BERT logits: {e}")
        bert_logits_for_loss = None
    else:
        print(f"Skipping BERT: Vocab size mismatch ({self.bert_vocab_size} vs {teacher_bert_vocab_size})")
        bert_logits_for_loss = None

# Calculate final loss (outside no_grad to ensure student grads)
temp = 2.0
soft_log_probs_student = F.log_softmax(student_logits / temp, dim=-1)
seq_len = student_logits.shape[1] # Use student's seq len

if llama_logits_for_loss is not None:
    # Ensure float and correct seq len
    t_llama = llama_logits_for_loss[:, :seq_len, :].float()
    soft_probs_llama = F.softmax(t_llama / temp, dim=-1)
    loss_llama = F.kl_div(soft_log_probs_student, soft_probs_llama, reduction='sum')
    loss_components.append(loss_llama)
    del llama_logits_for_loss, t_llama, soft_probs_llama

if bert_logits_for_loss is not None:
    # Ensure float and correct seq len
    t_bert = bert_logits_for_loss[:, :seq_len, :].float()
    soft_probs_bert = F.softmax(t_bert / temp, dim=-1)
    loss_bert = F.kl_div(soft_log_probs_student, soft_probs_bert, reduction='sum')
    loss_components.append(loss_bert)
    del bert_logits_for_loss, t_bert, soft_probs_bert

# Combine losses
if loss_components:
    final_loss = torch.mean(torch.stack(loss_components))
else:
    print("No compatible teachers found for loss calculation! Returning zero loss")
    final_loss = torch.tensor(0.0, device=self.args.device, requires_grad=False)

# CUDA Cleanup
del inputs, student_logits, soft_log_probs_student
gc.collect()
if torch.cuda.is_available():
    torch.cuda.empty_cache()

return (final_loss, student_outputs) if return_outputs else final_loss

print("DistillerTrainer class defined.")

```

TrainingArguments configured.
DistillerTrainer class defined.

Now the actual dataset instances can be created and the `DistillerTrainer` can be initialized.

```
In [9]: # Prepare Dataset and Trainer
train_dataset = create_dataset(tokenizer=tokenizer, max_length=MAX_SEQ_LENGTH)
eval_dataset = create_dataset(tokenizer=tokenizer, max_length=MAX_SEQ_LENGTH)

if train_dataset is None:
    print("Failed to create training dataset. Exiting.")
    exit()

# Initialize Trainer
trainer = DistillerTrainer( # Initialize with the class
    model=student_model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=eval_dataset,
    tokenizer=tokenizer, # Often needed for saving etc.
    teacher_llama=teacher_llama,
    teacher_bert=teacher_bert
)
print("Trainer initialized.") # Saying the trainer initialized
```

Loading dataset with 240 samples.

Loading dataset with 240 samples.

<ipython-input-8-45f504bd2f5a>:23: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `DistillerTrainer.__init__`. Use `processing_class` instead.

```
    super().__init__(*args, **kwargs)
```

DistillerTrainer initialized.

Vocab Sizes - Student: 30522, LLaMA: 128256, BERT: 30522

Trainer primary device: cuda:0

Trainer initialized.

```
In [10]: # Training the Model
print("Starting Training with memory optimizations and device fix...")
training_successful = False
try:
    train_result = trainer.train()
    print("Training finished successfully!")
    trainer.save_model()
    trainer.log_metrics("train", train_result.metrics)
    trainer.save_metrics("train", train_result.metrics)
    trainer.save_state()
    print("Training finished!")
    training_successful = True # Set flag on success
except Exception as e:
    print(f"Error during training: {e}") # Print an error message
```

[illegible]

Step Training Loss

5	38516.768800
10	39823.162500
15	39647.196900
20	40038.768800
25	39046.768800
30	39704.306300
35	37804.365600
40	39919.209400
45	40035.390600
50	38673.571900
55	40849.181300
60	39027.440600
65	40046.753100
70	39373.206300
75	38740.093800
80	38602.487500
85	39814.934400
90	39192.712500
95	39598.646900
100	38776.853100
105	38756.465600
110	40304.296900
115	37881.950000
120	39863.028100
125	38891.143800
130	39385.859400
135	39205.559400
140	38071.550000
145	39767.203100
150	39385.209400

```

Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Skipping LLaMA: Vocab size mismatch (128256 vs 30522)
Got BERT logits.
Training finished successfully!
***** train metrics *****
epoch                =          5.0
total_flos            =          0GF
train_loss            = 39291.4695
train_runtime         = 0:07:25.58
train_samples_per_second =      2.693
train_steps_per_second  =      0.337
Training finished!

```

Inference Examples

These two examples show how the combination model can perform CLM and MLM tasks

Example 1: Causal Language Modeling (Text Generation)

```

In [17]: print("\n--- Running Text Generation Example ---")
if 'student_model' in locals() and 'tokenizer' in locals():
    prompt = "Tarlov cysts are fluid filled "
    # Device handling for inference input
    inputs = tokenizer(prompt, return_tensors='pt').to(device) # Ensure inp

```

```

# Generate predictions using the student model (Distilled model)
max_length = len(inputs['input_ids'][0]) + 50
try:
    generated_ids = student_model.generate(
        inputs['input_ids'],
        max_length=max_length,
        temperature=1.0,
        top_p=0.92,
        top_k=50
    )

    # Decode the generated ids back to text
    generated_text = tokenizer.decode(generated_ids[0], skip_special_tokens=True)

    print(f"Prompt: {prompt}")
    print(f"Generated text: {generated_text}")
except Exception as e:
    print(f"Error during generation: {e}")
else:
    print("Skipping generation example: student_model or tokenizer not found")

```

--- Running Text Generation Example ---

Prompt: Tarlov cysts are fluid filled

Generated text: tarlov cysts are fluid filled buttons greenfield current exhaled campaigning clint suzy antics seamen crawley 303 giggling emphasize walnut gray congress flames hire invasion organizers arbitrary beaux deportation etudes sr filed [unused187] alexandre popularlais judge gripped版 [unused893] ico male neglect magician expandedods filmmakers borrowed stranddict blessing informally constituted electron abbreviated patrolling

Example 2: Masked Language Modeling (MLM) Prediction

Simulating an MLM task with the student model

```

In [15]: print("\n--- Running MLM Prediction Example ---")
if 'student_model' in locals() and 'tokenizer' in locals():

    # Example prompt for MLM
    prompt = "Tarlov cysts are fluid-filled nerve root cysts that can cause"

    # Tokenize the prompt
    inputs = tokenizer(prompt, return_tensors='pt')

    # Mask tokens randomly (simulating an MLM task)
    input_ids = inputs['input_ids'].clone()
    labels = input_ids.clone() # For MLM, labels are the same as input_ids

    # Randomly select positions to mask in the input
    # Need to import random if not already done
    import random
    num_masked_tokens = int(input_ids.size(1) * 0.3)
    # Simple random masking
    mask_positions = random.sample(range(input_ids.size(1)), num_masked_tokens)
    masked_input_ids = input_ids.clone() # Clone again to be safe
    masked_input_ids[0, mask_positions] = tokenizer.mask_token_id # Mask the

```

```

model_inputs = {k: v.to(device) for k, v in inputs.items()} # Move input
model_inputs['input_ids'] = masked_input_ids.to(device) # Ensure masked

# Print the masked input
print(f"Original text: {prompt}")
masked_input_text = tokenizer.decode(masked_input_ids[0], skip_special_t
print(f"Masked input: {masked_input_text}")

# Forward pass through the model (MLM)
try:
    with torch.no_grad(): # Good practice for inference
        outputs = student_model(**model_inputs) # Pass prepared inputs
        logits = outputs

    # Get the predictions for the masked tokens
    predictions = torch.argmax(logits, dim=-1)

    # Decode the predicted tokens to text
    predicted_tokens_text = tokenizer.decode(predictions[0], skip_specia

    print(f"Predicted text (full sequence): {predicted_tokens_text}")

    # Optional: Show predictions specifically for masked positions
    print("\nPredicted tokens for [MASK] positions (approximate):")
    full_predicted_ids = masked_input_ids.clone()
    for i in mask_positions:
        pred_id = predictions[0, i].item()
        pred_token = tokenizer.decode([pred_id])
        orig_id = input_ids[0, i].item() # Original input_ids before ma
        orig_token = tokenizer.decode([orig_id])
        print(f"- Pos {i}: Predicted='{pred_token}' (Original='{orig_to
        full_predicted_ids[0, i] = pred_id
    filled_text = tokenizer.decode(full_predicted_ids[0], skip_special_t
    print(f"\nFull sentence with predictions filled in:\n{filled_text}")

except Exception as e:
    print(f"Error during MLM prediction: {e}")
else:
    print("Skipping MLM example: student_model or tokenizer not found.")

```

--- Running MLM Prediction Example ---

Original text: Tarlov cysts are fluid-filled nerve root cysts that can cause various symptoms.

Masked input: tarlov cysts are filled root cysts cause various.

Predicted text (full sequence):

Predicted tokens for [MASK] positions (approximate):

- Pos 6: Predicted='.' (Original='fluid')
- Pos 17: Predicted='.' (Original='symptoms')
- Pos 13: Predicted='.' (Original='that')
- Pos 7: Predicted='.' (Original='-')
- Pos 14: Predicted='.' (Original='can')
- Pos 9: Predicted='.' (Original='nerve')

Full sentence with predictions filled in:

tarlov cysts are.. filled. root cysts.. cause various..

This notebook was converted with convert.ploomber.io

References

- Google Cloud Tech. (2023, June 5). Transformer models and BERT model: Overview [Video]. YouTube. http://www.youtube.com/watch?v=t45S_MwAcOw
- Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., ... & Poon, H. (2021). Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1), 1-23.
- Koroteev, M. V. (2021). BERT: a review of applications in natural language processing and understanding. *arXiv preprint arXiv:2103.11943*.
- Min, B., Ross, H., Sulem, E., Veyseh, A. P. B., Nguyen, T. H., Sainz, O., ... & Roth, D. (2023). Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys*, 56(2), 1-40.
- Porter, J. (2024, May 28). OpenAI cites DeepSeek as evidence that AI distillation works. **The Verge**. <https://www.theverge.com/news/601195/openai-evidence-deepseek-distillation-ai-data>
- Sanderson, G. (2024, April 1). Transformers (how LLMs work) explained visually | DL5 [Video]. YouTube. <https://www.youtube.com/watch?v=wjZofJX0v4M>
- Sun, S., Cheng, Y., Gan, Z., & Liu, J. (2019). Patient knowledge distillation for bert model compression. *arXiv preprint arXiv:1908.09355*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, 30.
- West, P., Bhagavatula, C., Hessel, J., Hwang, J. D., Jiang, L., Bras, R. L., ... & Choi, Y. (2021). Symbolic knowledge distillation: from general language models to commonsense models. *arXiv preprint arXiv:2110.07178*.
- Xu, X., Li, M., Tao, C., Shen, T., Cheng, R., Li, J., ... & Zhou, T. (2024). A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*.