



FOM Hochschule für Ökonomie & Management

Studienzentrum Frankfurt am Main

Projektarbeit

zum Thema

Analyse der Auswirkungen von Pruning an *Large Language
Models*

von

Kai Daniel Herbst und Sebastian Viet

28.02.2025

Gutachter:	Prof. Dr. Rüdiger Buchkremer
Zeitraum:	01.09.2024 - 28.02.2025
Semester:	3
Modul:	Big-Data-Analyseprojekt

Inhaltsverzeichnis

1	Einleitung	1
1.1	Forschungsfragen	1
2	Theorie	2
3	Methodik	3
3.1	Literaturrecherche	3
3.2	Arbeitsorganisation	3
3.3	Entwicklungsumgebung	3
3.4	Verwendetes Large Language Model	4
3.5	Pruning	5
3.5.1	Verwendetes Framework	5
3.5.2	Durchführung des Prunings	6
3.6	Fine-Tuning	7
3.7	Evaluierung der Modelle	8
3.7.1	Modell-Benchmarks	8
3.7.2	Modellperformance	9
4	Ergebnisse	10
4.1	Evaluierung Basismodell	10
4.2	Tatsächliche Anzahl geprunter Parameter	11
4.3	Evaluierung Pruning zu 30%	12
4.4	Evaluierung Pruning zu 40%	13
4.5	Evaluierung Pruning zu 70%	14
4.6	Speicherverbrauch	15
5	Diskussion	16
6	Fazit	17

1 Einleitung

1.1 Forschungsfragen

Die vorliegende Projektarbeit setzt sich mit insgesamt vier Forschungsfragen auseinander, die im Verlauf der Arbeit untersucht und beantwortet werden sollen. Die erste Forschungsfrage beschäftigt sich mit dem aktuellen Stand der Forschung im Bereich des Prunings von Large Language Models (LLMs). Dabei wird analysiert, welche Methoden derzeit für das Pruning angewendet werden, welche Unterschiede zwischen diesen bestehen und welche Vor- und Nachteile die jeweiligen Ansätze mit sich bringen. Das Ziel dieser Untersuchung ist es, einen Überblick über den aktuellen Stand der Forschung in diesem Bereich zu geben und die relevantesten Methoden zu beschreiben.

Ein weiterer Aspekt dieser Arbeit ist die Analyse der verschiedenen Frameworks, die Pruning-Methoden für LLMs unterstützen. Hierbei soll untersucht werden, welche dieser Frameworks welche Methoden implementieren und wie sie in der Praxis angewendet werden können. Basierend auf dieser Analyse wird eine geeignete Auswahl getroffen und eines dieser Frameworks für die praktische Umsetzung festgelegt.

Neben der theoretischen Auseinandersetzung mit dem Thema wird zudem ein eigener praktischer Versuch durchgeführt. Hierzu wird ein geeignetes Basismodell ausgewählt, das anschließend in verschiedenen Stufen geprunt wird. Die Auswirkungen dieses Prunings werden untersucht und dokumentiert. Im Rahmen dieser Analyse werden standardisierte Benchmarks erstellt, um die geprunten Modelle zu bewerten. Dabei wird nicht nur die Performance hinsichtlich der Ergebnisse in den Tests betrachtet, sondern auch der Einfluss des Prunings auf den Speicherverbrauch und die benötigte Rechenzeit. Ziel ist es, Erkenntnisse darüber zu gewinnen, wie sich das Pruning auf verschiedene Aspekte der Modelle auswirkt.

Zusammenfassend ergeben sich daraus die folgenden Forschungsfragen:

1. Wie ist der aktuelle Stand der Forschung im Bereich des Prunings von Large Language Models, und welche Methoden werden derzeit eingesetzt?
2. Welche Frameworks bieten Unterstützung für das Pruning von LLMs, und wie lassen sich diese in der Praxis anwenden?
3. Welche Auswirkungen hat das Pruning eines Modells auf dessen Leistungsfähigkeit in Bezug auf standardisierte Benchmarks und Aufgaben?
4. Inwiefern beeinflusst das Pruning eines Modells dessen Speicherverbrauch sowie die benötigte Rechenzeit?

2 Theorie

Hier kommt alles zur Theorie hin!

3 Methodik

3.1 Literaturrecherche

3.2 Arbeitsorganisation

Für die Zusammenarbeit und Organisation der anstehenden Aufgaben und zu verfassenden Kapitel wurde das Notizenprogramm *Notion* verwendet. Notion erlaubt das Erstellen kollaborativer Notizbücher, die parallel von mehreren Personen bearbeitet werden können. Über die Datenbankfunktion lassen sich sowohl Meilensteinpläne als auch Task-Boards erstellen, die zur Planung der Arbeit verwendet wurden.

3.3 Entwicklungsumgebung

Da sowohl das Pruning von Large Language Models als auch die anschließende Evaluierung der geprunten Modelle erhebliche Rechenressourcen erfordern, wurde entschieden, diese Prozesse auf einem leistungsstarken Server in der Cloud durchzuführen. Insbesondere das Testen der resultierenden Modelle stellt eine hohe Belastung für die verfügbaren Ressourcen dar und ist auf herkömmlicher Hardware nicht effizient durchzuführen. Zu diesen Zweck wurden EC2-Instanzen des Cloud-Providers AWS (Amazon Web Services) aus der *g*-Instanzfamilie verwendet. Diese Instanzfamilie ist speziell auf rechenintensive Anwendungen ausgelegt und bietet GPUs (Graphics Processing Units), die für parallele Berechnungen optimiert sind und den Einsatz von CUDA, dem von *NVIDIA* entwickelten Toolkit zur Verwendung von GPUs in allgemeinen Rechenaufgaben, ermöglichen.

Im Detail wurden Instanzen des Typs *gd4n.xlarge* verwendet, die mit 16 GiB Hauptspeicher, einem leistungsstarken Prozessor der *Intel Xeon Family* sowie einer *NVIDIA T4 Tensor Core* Grafikeinheit ausgestattet sind. Für den Start der Instanzen wurde das speziell auf Deep-Learning-Anwendungen zugeschnittene *Deep Learning OSS Nvidia Driver AMI GPU PyTorch 2.5.1 (Ubuntu 22.04) 20241208* Ubuntu-Image verwendet. In diesem Image ist das benötigte Python-Package *PyTorch* bereits vorinstalliert und ist zusätzlich mit *CUDA* ausgestattet.

Verwendet wurden im Detail Instanzen des Types *gd4n.xlarge*, die mit 16GiB Hauptspeicher, einem physischen Prozessor der *Intel Xeon Family* und einer *NVIDIA T4 Tensor Core* Grafikeinheit kommen. Gestartet wurden die Instanzen mit dem *Deep Learning OSS Nvidia Driver AMI GPU PyTorch 2.5.1 (Ubuntu 22.04) 20241208* Ubuntu Image, das bereits das benötigte Python-Package *PyTorch* vorinstalliert hat. Diese, auf Deep Learning ausgerichteten, Images haben neben PyTorch bereits *CUDA* vorinstalliert und ermöglichen so einen einfachen Start in das Arbeiten mit Programmen, die die GPU benötigen.

Nachdem über SSH (Secure Shell) mit der jeweiligen Instanz verbundene wurde, wurde zunächst die vorinstallierte PyTorch-Umgebung gestartet und anschließenden das Reposi-

Verwendete EC2-Instanz	
Instanz-Typ	gd4n.xlarge
Hauptspeicher	16GiB
vCPU	4
Clock Speed	2.5GHz
GPU	nvidia t4 tensor core
CPU	Intel Xeon Family
AMI	Deep Learning OSS Nvidia Driver AMI GPU PyTorch 2.5.1
AMI-ID	ami-0fa5e5fd27b3e163a

Tabelle 1: Attribute der verwendeten Hardware

tory des verwendeten Frameworks geklont.

```
1 $ source activate pytorch
2 $ git clone https://github.com/horseee/LLM-Pruner.git
```

Da die auf dem Image vorinstallierte PyTorch-Umgebung mit dem Conda-Package-Manager bereitgestellt wird, wurden die im Repository in der *requirements.txt*-Datei angegebenen Package über Conda installiert.

```
1 $ conda install transformers sentencepiece datasets wandb ...
```

Damit ist die Versuchsumgebung identisch zu der in dieser Arbeit verwendeten Umgebung.

3.4 Verwendetes Large Language Model

Aufgrund von finanziellen und zeitlichen Einschränkungen wurde für diese Analyse das TinyLlama-Modell (*TinyLlama/TinyLlama-1.1B-Chat-v1.0*) als Basismodell gewählt. Mit 1,1 Milliarden Parametern ist es im Vergleich zu modernen, größeren Modellen wie *Llama 3.1* – das mit 405 Milliarden Parametern deutlich umfangreicher ist – relativ klein. Das verwendete TinyLlama wurde auf einem Datensatz von drei Milliarden Token vortrainiert und basiert auf der gleichen Architektur wie die *Llama 2*-Modelle. Diese Wahl ermöglichte es, innerhalb der verfügbaren Ressourcen eine Analyse durchzuführen, während gleichzeitig die Komplexität des Modells berücksichtigt wurde.

Der nachfolgende Auszug zeigt die detaillierte Architektur des Modells:

```
1 LlamaModel(
2   (embed_tokens): Embedding(32000, 2048)
3   (layers): ModuleList(
4     (0-21): 22 x LlamaDecoderLayer(
5       (self_attn): LlamaSdpaAttention(
6         (q_proj): Linear(in_features=2048, out_features=2048, bias=False)
```

```

7         (k_proj): Linear(in_features=2048, out_features=256, bias=False)
8         (v_proj): Linear(in_features=2048, out_features=256, bias=False)
9         (o_proj): Linear(in_features=2048, out_features=2048, bias=False)
10        (rotary_emb): LlamaRotaryEmbedding()
11    )
12    (mlp): LlamaMLP(
13        (gate_proj): Linear(in_features=2048, out_features=5632, bias=False)
14        (up_proj): Linear(in_features=2048, out_features=5632, bias=False)
15        (down_proj): Linear(in_features=5632, out_features=2048, bias=False)
16        (act_fn): SiLU()
17    )
18    (input_layernorm): LlamaRMSNorm((2048,), eps=1e-05)
19    (post_attention_layernorm): LlamaRMSNorm((2048,), eps=1e-05)
20 )
21 )
22 (norm): LlamaRMSNorm((2048,), eps=1e-05)
23 (rotary_emb): LlamaRotaryEmbedding()
24 )

```

Die Architektur besteht aus einer Embedding-Schicht, über die die Eingabesequenzen in Vektoren der Dimension 2048 umwandelt werden. Die Token-Embeddings basieren auf einem Vokabular von 32.000 Wörtern. Darauf folgt eine Modul-Liste mit 22 LlamaDecoderLayer, die jeweils aus mehreren Submodulen bestehen. Die Decoder-Schicht enthält einen Self-Attention-Mechanismus (LlamaSdpaAttention).

Um die grundlegenden Funktionen der ersten und letzten Schichten nicht zu beeinflussen, wurden sowohl im Attention-Abschnitt als auch im MLP-Abschnitt ausschließlich die Layer 4 bis 18 für das Pruning verwendet. Die Schichten der Architektur außerhalb des Decoder-Layers wie bspw. die Embedding-Schichten werden grundsätzlich nicht berührt, da sie für die Umwandlung der Texteingaben in die korrekte Repräsentation nötig sind.

3.5 Pruning

Das nachfolgende Kapitel beinhaltet die Vorgehensweise und verwendeten Technologien, die für das Pruning des TinyLlama-Modells verwendet wurde.

3.5.1 Verwendetes Framework

Für die Durchführung des Prunings standen zwei Frameworks zur Auswahl: *LLM-Pruner* und *Wanda* (Pruning by Weights and Activations). Während der *LLM-Pruner* ausschließlich das im vorherigen Kapitel beschriebene strukturierte Pruning unterstützt, bietet *Wanda* zusätzlich die Möglichkeit, unstrukturiertes Pruning durchzuführen. Trotz dieser zusätzlichen Funktionalität wurde für die Umsetzung dieser Untersuchung aus verschiedenen Gründen ausschließlich der *LLM-Pruner* verwendet.

Ein wesentlicher Faktor für diese Entscheidung war die Aktualität der Projekte. Die letzten Updates im *Wanda*-Projekt wurden Ende 2023 vorgenommen. Dies lag zum Zeitpunkt des Verfassens dieser Arbeit bereits mehr als ein Jahr zurück. Im Gegensatz dazu wird der *LLM-*

Pruner weiterhin aktiv weiterentwickelt, was eine aktuellere und besser unterstützte Basis bietet. Ein weiterer wichtiger Aspekt ist die Unterstützung spezifischer Modelle. Während beim *LLM-Pruner* explizit die Kompatibilität mit dem *TinyLlama*-Modell hervorgehoben wird, fehlt eine solche Aussage im Fall von *Wanda*. Es ist zwar anzunehmen, dass die Methoden von *Wanda* aufgrund der Unterstützung von *Llama-2* – dessen Architektur dem *TinyLlama* ähnlich ist – ebenfalls auf das *TinyLlama*-Modell anwendbar sein könnten. Allerdings bleibt dies ungewiss, da eine direkte Unterstützung nicht garantiert wird.

Ein weiterer entscheidender Punkt war, dass trotz mehrerer Versuche mit *Wanda* kein erfolgreiches Pruning durchgeführt werden konnte. Angesichts dieser Schwierigkeiten und unter Berücksichtigung der bereits genannten Argumente fiel die Entscheidung, sich ausschließlich auf den *LLM-Pruner* zu fokussieren.

3.5.2 Durchführung des Prunings

Das Pruning wurde, wie bereits erläutert, anhand des *TinyLlama*-Modells getestet und untersucht. Dabei wurden drei unterschiedliche Pruning-Ratios gewählt, um die Auswirkungen verschiedener Verkleinerungsgrade des Modells zu analysieren. Da das *LLM-Pruner*-Projekt bereits eigene Ergebnisse für das *TinyLlama*-Modell mit einer Pruning-Ratio von 20% veröffentlicht hatte, wurde dieses Verhältnis in der vorliegenden Untersuchung nicht erneut evaluiert. Stattdessen wurden die bereits existierenden Ergebnisse mit den Verhältnissen von 30%, 40% und schließlich 70% verglichen. Obwohl eine Verkleinerung um 70% bei einem ohnehin bereits kompakten Modell wie *TinyLlama* als unrealistisch erscheint, wurde dennoch im Rahmen der Untersuchung überprüft, welche Ergebnisse das Modell bei einer solch extremen Reduktion in den Tests liefert.

Das Pruning erfolgt stets über den folgenden Befehl, der ausgeführt werden kann, sobald sich im *LLM-Pruner*-Repository auf der höchsten Ebene befindet:

```
1 $ python llama3.py
2     --base_model TinyLlama/TinyLlama-1.1B-Chat-v1.0
3     --pruning_ratio [PRUNING_RATIO]
4     --device cuda
5     --eval_device cuda
6     --block_wise
7     --block_mlp_layer_start [START_MLP_LAYER]
8     --block_mlp_layer_end [END_MLP_LAYER]
9     --block_attention_layer_start [START_ATTENTION_LAYER]
10    --block_attention_layer_end [END_ATTENTION_LAYER]
11    --save_ckpt_log_name [SAVE_PATH]
12    --pruner_type [PRUNER_TYPE]
13    [--taylor param_first]
14    --save_model
15    --max_seq_len 2048
16    --test_after_train
```

Für das Pruning des *TinyLlama*-Modells wurde das Skript *llama3.py* verwendet, das speziell für das Pruning von *Llama3*-Modellen entwickelt wurde. In dieser Untersuchung

diente stets das zuvor beschriebene Modell *TinyLlama/TinyLlama-1.1B-Chat-v1.0* als *base_model*. Wie bereits erwähnt, wurden für die Pruning-Ratio die Verhältnisse 0.3, 0.5 und 0.7 ausgewählt, um unterschiedliche Stufen der Modellreduktion zu analysieren.

Alle Prozesse wurden in einer CUDA-fähigen Umgebung ausgeführt, weshalb die Anweisung, die GPU für die Berechnungen zu verwenden, stets mitgegeben wurde. Von den insgesamt 22 verfügbaren MLP- und Attention-Layern des Modells wurden für das Pruning jeweils die Layer vier bis 18 berücksichtigt.

Bezüglich der verfügbaren Pruner-Typen bot der *LLM-Pruner* vier verschiedene Optionen an: *Taylor*, *L1*, *L2* und *random*. Jede dieser vier Methoden wurde für jede der drei gewählten Pruning-Ratios getestet und evaluiert, um ihre jeweiligen Auswirkungen auf die Modellleistung zu untersuchen.

Zusätzlich wurde bei allen Experimenten das Argument *test_after_train* verwendet. Dadurch wurde nach jedem Pruning automatisch die Perplexity des Modells ermittelt, basierend auf den beiden Testdatensätzen *wikidata2* und *ptb* (*Penn Treebank*).

Um das TinyLlama-Modell zu 30% mit der *Taylor*-Methode zu prunen sieht der Befehl demnach beispielhaft wie folgt aus:

```
1 $ python llama3.py
2     --base_model TinyLlama/TinyLlama-1.1B-Chat-v1.0
3     --pruning_ratio 0.3
4     --device cuda
5     --eval_device cuda
6     --block_wise
7     --block_mlp_layer_start 4
8     --block_mlp_layer_end 18
9     --block_attention_layer_start 4
10    --block_attention_layer_end 18
11    --save_ckpt_log_name tinyllama_30_0616_prune_log
12    --pruner_type taylor
13    --taylor param_first
14    --save_model
15    --max_seq_len 2048
16    --test_after_train
```

3.6 Fine-Tuning

Für das nach dem Pruning stattfindende Fine-Tuning wird vom LLM-Pruner *PEFT* (Parameter-Efficient Fine-Tuning) verwendet. PEFT stellt eine Methode dar, um große vortrainierte Modelle an spezifische Aufgaben anzupassen, ohne den gesamten Parameter-raum des Modells zu optimieren. Stattdessen wird nur ein kleinerer Teil der Parameter während des Trainings modifiziert.

In den Evaluierungen der Modelle, die direkt im Anschluss an das Pruning durchgeführt wurden, hat sich bereits gezeigt, dass über die Pruning-Methode *Taylor* die vielversprechendsten Ergebnisse erzielt werden konnten. Diese geprunten Modelle konnten die höchsten

Werte in den Evaluierungen erreichen. Das rechen- und kostenintensive Fine-Tuning wurde daher nur testweise für das Modell durchgeführt, das zu 30% mit der *Taylor*-Methode geprunt wurden.

Verwendet wurde dafür das im *LLM-Pruner* vorhandene Skript *post_training.py* über den folgenden Befehl:

```
1 $ python post_training.py
2     --prune_model prune_log/tinyllama_30_0418_prune_log/pytorch_model.bin
3     --data_path yahma/alpaca-cleaned \
4     --lora_r 8 \
5     --num_epochs 2 \
6     --learning_rate 1e-4 \
7     --batch_size 64 \
8     --output_dir tune_log/tinyllama_30_tuned \
9     --wandb_project tinyllama_30_tune
```

3.7 Evaluierung der Modelle

Als Basis jeglicher durchgeführter Tests dienten jeweils die Ergebnisse des selbst durchgeführten Benchmarks des TinyLlama-Basismodells. Alle erhobenen Werte und Ergebnisse werden in Relation dazu bewertet.

3.7.1 Modell-Benchmarks

Die Modelle, die durch das Pruning und das anschließende Fine-Tuning entstanden sind, wurden anschließend auf ihre verbliebenen Fähigkeiten hin untersucht. Zu diesem Zweck wurden sie anhand verschiedener Datensätze evaluiert, die jeweils unterschiedliche Aspekte der Leistungsfähigkeit von LLMs testen. Welche spezifischen Aspekte dabei geprüft wurden, wurde bereits in den vorherigen Abschnitten detailliert beschrieben.

- *openbookqa*
- *winogrande*
- *hellaswag*
- *arc_challenge*
- *boolq*

Die Evaluierung anhand dieser Datensätze wurde, wie beim Pruning, mit dem immer gleichen Befehl - angepasst an das jeweilige Modell - durchgeführt.

```
1 $ export PYTHONPATH='.'
2 $ python lm-evaluation-harness/main.py
3     --model hf-causal-experimental
4     --model_args checkpoint=[PATH_TO_MODEL]/pytorch_model.bin,
5                     config_pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0
6     --tasks openbookqa,winogrande,hellaswag,arc_challenge,boolq
```

```
7 --device cuda:0
8 --no_cache
9 --output_path [PATH_TO_RESULTS]
```

Wie im Befehl ersichtlich, wurde das *lm-evaluation-harness*-Framework von OpenAI verwendet, das bereits im Repository enthalten ist. Dieses Framework dient der Evaluierung von Sprachmodellen anhand verschiedener Benchmarks bzw. *Tasks*. Durch die Nutzung des Frameworks in Kombination mit den definierten *Tasks* wird eine standardisierte Bewertung der Modelle ermöglicht, was wiederum den Vergleich unterschiedlicher Modelle erleichtert.

Das Argument *--model hf-causal-experimental* wird übergeben, um die Nutzung der GPU während der Tests zu gewährleisten. Zusätzlich ist die Angabe des Pfads zum geprüften bzw. nachtrainierten Modell erforderlich, ebenso wie die Grundarchitektur des Basismodells. Die im Befehl spezifizierten *Tasks* entsprechen dabei den zuvor beschriebenen.

Ein beispielhafter Aufruf zur Evaluierung des Modells, das zu 30% geprünt wurde, sieht wie folgt aus:

```
1 $ export PYTHONPATH='.'
2 $ python lm-evaluation-harness/main.py
3   --model hf-causal-experimental
4   --model_args checkpoint=prune_log/tinyllama_30_0418_11_prune_log/pytorch_model.bin,
5     config_pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0
6   --tasks openbookqa,winogrande,hellaswag,arc_challenge,boolq
7   --device cuda:0
8   --no_cache
9   --output_path results/tinyllama_30_0418_11
```

3.7.2 Modellperformance

Neben den verbliebenen Fähigkeiten des geprünten LLMs soll zusätzlich dessen Performance getestet werden. Die Messung der Performance erfolgt durch die Durchführung der zuvor beschriebenen Benchmarks. Während dieser Tests werden sowohl die benötigte Zeit als auch der in Anspruch genommene Speicher berücksichtigt, um eine Bewertung der Effizienz des Modells zu ermöglichen.

Für diese Messungen wird der in der *zsh*-Shell integrierte *time*-Befehl verwendet, der sowohl Angaben zur verbrauchten Zeit als auch zum genutzten Speicher liefert. Um diese Daten zu erfassen, wird der *time*-Befehl einfach den zuvor beschriebenen Befehlen zur Durchführung der Benchmarks vorangestellt. Ein entsprechender Aufruf zur Messung der Performance sieht wie folgt aus:

```
1 $ export PYTHONPATH='.'
2 $ time python lm-evaluation-harness/main.py
3   --model hf-causal-experimental
4   ...
```

4 Ergebnisse

Im folgenden Kapitel werden die Ergebnisse der Evaluierungen und Tests der geprunten sowie teilweise anschließend trainierten Modelle vorgestellt. Dabei liegt der Fokus insbesondere auf den ermittelten Perplexity-Werten sowie den Ergebnissen der Tests mit dem *lm-evaluation-harness*-Framework. Die verschiedenen Pruning-Stufen und die dabei verwendeten Methoden werden einzeln analysiert, bevor abschließend ein umfassender Gesamtüberblick über die Ergebnisse gegeben wird.

4.1 Evaluierung Basismodell

Um eine Vergleichsbasis zu schaffen, muss zunächst das Basismodell in allen Tests evaluiert werden. Als Basismodell dient hierbei das *TinyLlama*-Modell (*TinyLlama/TinyLlama-1.1B-Chat-v1.0*). Die Darstellung der Ergebnisse orientiert sich an den vom *LLM-Pruner* für andere Modelle bereitgestellten Resultaten, um eine einheitliche Vergleichbarkeit zu ermöglichen.

Die Spalte *Average* gibt den Durchschnitt der getesteten *Tasks* wieder. Die Ergebnisse für *WikiText2* und *PTB* werden dabei nicht in diese Berechnung einbezogen, da bei diesen Benchmarks ein niedrigerer Score eine bessere Leistung des Modells widerspiegelt.

Die zusammengefassten Ergebnisse sind in der folgenden Tabelle dargestellt:

Pruning Ratio	Method	WikiText2	PTB	BoolQ	HellaSwag	WinoGrande	ARC-c	OBQA	Average
Pruned 0%	–	7.97	20.76	61.31	46.15	60.30	30.12	24.20	39,58

Tabelle 2: Evaluierung des Basismodells

Wie hier zu erkennen ist, weist das Basismodell bereits einen vergleichsweise niedrigen Score von *39,58* auf. Zum Vergleich: In den Ergebnissen des *LLM-Pruners* für das Modell *Llama7B* ergibt sich ein Durchschnittswert von *68,59*. Allerdings wurden in diesen Tests noch weitere *Tasks* berücksichtigt, die in der vorliegenden Analyse nicht enthalten sind.

Da sich diese Untersuchung jedoch auf die relative Verschlechterung im Vergleich zum Basismodell konzentriert, stellt der niedrigere Ausgangswert hier kein Problem dar.

4.2 Tatsächliche Anzahl geprunter Parameter

Anhand der Modellevaluierungen, die nach dem dem Pruning über das Evaluierungsskript erstellt wurden, lässt sich schnell erkennen, dass die im Befehl angegebene Menge an zu prunenden Parametern bzw. das definierte Verhältnis nicht exakt so vom *LLM-Pruner* umgesetzt wird.

Am Beispiel des 30%-Prunings sieht dies wie folgt aus: Bei einer angegebenen Pruning-Ratio von 30% wäre zu erwarten, dass nach dem Pruning noch 70% der ursprünglichen Parameter des Basismodells erhalten bleiben. Tatsächlich sind es in diesem Fall jedoch 73,06%, also 3,06% mehr als ursprünglich vorgesehen.

Noch deutlicher wird diese Abweichung bei einer Pruning-Ratio von 70%. Hier sollten theoretisch nur noch 30% der Parameter im Modell verbleiben, tatsächlich sind es jedoch noch 55,08%, was einer Differenz von 25,08% entspricht.

Diese Diskrepanz muss bei der Interpretation der nachfolgenden Ergebnisse berücksichtigt werden, da die angegebenen Verhältnisse nicht mit den tatsächlichen übereinstimmen.

Specified ratio	#Parameters before	#Parameters after	Actual ratio
30%	1261.53 Mio	921.65 Mio	73.06%
40%	1261.53 Mio	873.22 Mio	69.21%
70%	1261.53 Mio	694.83 Mio	55,08%

Tabelle 3: Anzahl der vorhandenen Parameter nach dem Pruning

4.3 Evaluierung Pruning zu 30%

In Tabelle 4 sind die Ergebnisse der geprunten Modelle sowie die des Basismodells dargestellt. Getestet wurden – wie in Kapitel 3 beschrieben – die *Tasks* BoolQ, HellaSwag, WinoGrande, Arc-Challenge und OpenBookQA. Zusätzlich wurde die *Perplexity* anhand der Datensätze WikiText2 und PTB ermittelt.

Pruning Ratio	Method	WikiText2	PTB	BoolQ	HellaSwag	WinoGrande	ARC-c	OBQA	Average
Pruned 0%	–	7.97	20.76	61.31	46.15	60.30	30.12	24.20	39.58
Pruned 30%	Taylor	15.19	43.19	55.50	37.90	54.22	24.66	23.20	39.10
	L1	281.63	4992.16	46.54	28.53	48.46	21.33	14.00	31.77
	L2	42.19	167.51	60.89	35.97	54.30	22.95	18.40	38.50
	Random	40.89	166.21	59.63	33.18	53.99	20.99	18.00	37.17
Pruned 30% (tuned)	Taylor	/	/	45.75	39.29	56.67	25.26	22.00	37.79

Tabelle 4: Evaluierungen bis 30% Pruning

Es ist wichtig zu beachten, dass die tatsächliche Anzahl an Parametern im Vergleich zum Basismodell nur um *16,22%* reduziert wurde – und nicht, wie ursprünglich erwartet, um volle *30%*. Die Tests wurden unmittelbar nach dem Pruning durchgeführt, ohne dass ein erneutes Fine-Tuning erfolgte.

Betrachtet man ausschließlich die Ergebnisse der *Tasks* und deren Durchschnittswerte, so schneiden die *Taylor*- und *L2*-Methoden am besten ab, da ihre Werte am nächsten an denen des Basismodells liegen. Allerdings zeigt sich bei der *L2*-Methode eine deutlich höhere und damit schlechtere *Perplexity* in beiden Datensätzen im Vergleich zur *Taylor*-Methode. Hervorzuheben ist dennoch, dass sie in der *BoolQ*-Task um *5,39* Prozentpunkte besser abgeschnitten hat als die *Taylor*-Methode.

Die *Random*-Methode weist mit dessen Resultaten Ähnlichkeiten zu der *L2*-Methode auf und schneidet somit ebenfalls schlechter als *Taylor*-Methode ab. Hier sind somit keine weiteren Auffälligkeiten Hervorzuheben.

Am schlechtesten hat die *L1*-Methode abgeschnitten: Das daraus resultierende Modell weist extrem schlechte *Perplexity*-Werte auf, und auch der Durchschnitt der *Tasks* liegt im Vergleich zu den anderen Methoden deutlich niedriger.

In der letzten Zeile der Tabelle sind zusätzlich die Ergebnisse des nachtrainierten Taylor-Modells zu sehen. Leider konnten hier nicht ähnliche Ergebnisse wie in der Dokumentation des *LLM-Pruners* erzielt werden. Nach einem dreistündigen Post-Training ist der Durchschnittswert von 39,85 auf 37,79 gesunken. Das Modell hat damit nach dem Post-Training schlechtere Ergebnisse geliefert als davor. Allein in den Evaluierungsaufgaben *HellaSwag* und *ARC-c* konnte eine leichte Steigerung erzielt werden.

4.4 Evaluierung Pruning zu 40%

In Tabelle 5 sind die Inhalte der Tabelle 4 ergänzt um die Ergebnisse der Tests zu den Modellen, die mit der Angabe 40% geprunt wurden zu sehen. Es wurden erneut jeweils die vier möglichen Methoden angewendet und getestet.

Pruning Ratio	Method	WikiText2	PTB	BoolQ	HellaSwag	WinoGrande	ARC-c	OBQA	Average
Pruned 0%	–	7.97	20.76	61.31	46.15	60.30	30.12	24.20	39.58
Pruned 30%	Taylor	15.19	43.19	55.50	37.90	54.22	24.66	23.20	39.10
	L1	281.63	4992.16	46.54	28.53	48.46	21.33	14.00	31.77
	L2	42.19	167.51	60.89	35.97	54.30	22.95	18.40	38.50
	Random	40.89	166.21	59.63	33.18	53.99	20.99	18.00	37.17
Pruned 30% (tuned)	Taylor	/	/	45.75	39.29	56.67	25.26	22.00	37.79
Pruned 40%	Taylor	18.43	53.33	59.39	34.84	52.88	23.55	18.80	37.89
	L1	441.35	14000.87	48.99	28.13	49.96	21.42	14.80	32.66
	L2	86.23	229.87	60.83	34.18	51.62	21.08	19.60	37.59
	Random	86.91	364.46	57.89	30.82	52.33	20.56	17.00	35.72

Tabelle 5: Evaluierungen bis 40% Pruning

Auch hier zeigt sich, dass die Taylor-Methode die vielversprechendsten Ergebnisse erzielt.

4.5 Evaluierung Pruning zu 70%

Pruning Ratio	Method	WikiText2	PTB	BoolQ	HellaSwag	WinoGrande	ARC-c	OBQA	Average
Pruned 0%	–	7.97	20.76	61.31	46.15	60.30	30.12	24.20	39.58
Pruned 30%	Taylor	15.19	43.19	55.50	37.90	54.22	24.66	23.20	39.10
	L1	281.63	4992.16	46.54	28.53	48.46	21.33	14.00	31.77
	L2	42.19	167.51	60.89	35.97	54.30	22.95	18.40	38.50
	Random	40.89	166.21	59.63	33.18	53.99	20.99	18.00	37.17
Pruned 30% (tuned)	Taylor	/	/	45.75	39.29	56.67	25.26	22.00	37.79
Pruned 40%	Taylor	18.43	53.33	59.39	34.84	52.88	23.55	18.80	37.89
	L1	441.35	14000.87	48.99	28.13	49.96	21.42	14.80	32.66
	L2	86.23	229.87	60.83	34.18	51.62	21.08	19.60	37.59
	Random	86.91	364.46	57.89	30.82	52.33	20.56	17.00	35.72
Pruned 70%	Taylor	83.25	274.04	52.39	28.43	48.70	19.11	17.00	33.12
	L1	37762.14	11607.13	46.91	25.98	50.43	20.39	16.40	32.02
	L2	394.08	783.72	60.49	26.84	51.62	20.31	13.20	34.49
	Random	4138.65	4074.48	54.74	26.64	51.38	20.99	14.20	33.59

Tabelle 6: Evaluierungen bis 70% Pruning

4.6 Speicherverbrauch

Pruning ratio	Computation Complexity	Ratio	GPU Memory Requirement	Ratio	Ratio Parameters
Pruned 0%	77.32 GMac	100%	2427.31 MiB	100%	100%
Pruned 30%	54.8 GMac	70.87%	1793.30 MiB	73.86%	73.06%
Pruned 40%	51.7 GMac	66.86%	1709.30 MiB	70.31%	69.21%
Pruned 70%	40.28 GMac	52.09%	1337.05 MiB	55.08%	55.08%

Tabelle 7: Speicheransprüche nach dem Pruning

5 Diskussion

Hier kommt alles zur Diskussion hin!

6 Fazit

Hier kommt das Fazit hin!