



**FOM Hochschule für Ökonomie & Management**

Studienzentrum Frankfurt am Main

**Zusätzliche Beteiligung**

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 (Kai Herbst)	1
2	Versuchsumgebung (Kai Herbst)	2
3	Basismodell (Kai Herbst)	4
4	Pruning (Kai Herbst)	5
5	Performance-Evaluierung (Kai Herbst)	9
6	Tuning (Kai Herbst)	13
7	Evaluierung Rechenanforderung (Kai Herbst)	14
8	Fazit (Sebastian Viet)	18

# 1 Einleitung (Kai Herbst)

Im Rahmen des Moduls *Big-Data-Analyseprojekt* wird das *Pruning* von *Large Language Models* (LLM) theoretisch untersucht und anschließend in einem praktisch Versuch mit Hilfe des Python-Packages *LLM-Pruner* selbst am Llama-Modell *TinyLlama/TinyLlama-1.1B-Chat-v1.0* durchgeführt. Die aus dem Pruning resultierenden Modelle werden anschließend den Ergebnissen verschiedenen Tests evaluiert und mit dem originalen Modell verglichen.

Die dafür verwendeten Aufrufe der Skripte des *LLM-Pruners* und dazugehörige Konsolenausgabe sind in dieser Arbeit dokumentiert.

## 2 Versuchsumgebung (Kai Herbst)

Das Pruning wurde aufgrund der erforderlichen Rechenleistung auf einer EC2-Instanz in AWS (Amazon Web Services) durchgeführt. Dabei verwendet wurde eine EC2-Instanz vom Instanz-Typ *gd4n.xlarge* mit 16GiB Hauptspeicher und vier vCPU.

Zusätzlich wurde das von AWS bereitgestellte *Pytorch Deep Learning* AMI verwendet. Darauf vorinstalliert ist bereits eine Conda-Umgebung mit PyTorch und alle benötigten CUDA-Treiber, um PyTorch die Möglichkeit zu bieten die GPU zu verwenden.

In Tabelle 1 sind die genauen Spezifikationen der verwendeten Instanz aufgelistet.

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

Über einen in der AWS-Konsole für die Instanz erstellten Schlüssel kann sich anschließend über ssh mit der Instanz verbunden werden.

Anschließend wird die PyTorch-Umgebung aufgerufen und das GitHub-Repository des *LLM-Pruners* geklont.

---

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

---

Die in der *requirements.txt* spezifizierten Packages können sowohl mit Conda als

auch mit Pip installiert werden.

### 3 Basismodell (Kai Herbst)

Für das Pruning wurde das LLM *TinyLlama/TinyLlama-1.1B-Chat-v1.0* verwendet.

Mit folgendem Skript wurden genauere Informationen über die Architektur des Modells generiert:

---

```
1 from transformers import AutoModel
2
3 model = AutoModel.from_pretrained("TinyLlama/TinyLlama-1.1B-Chat-v1.0")
4 print(model)
```

---

Das Skript erzeugt folgende Ausgabe:

---

```
1 $ python printModel.py
2
3 LlamaModel(
4   (embed_tokens): Embedding(32000, 2048)
5   (layers): ModuleList(
6     (0-21): 22 x LlamaDecoderLayer(
7       (self_attn): LlamaSdpaAttention(
8         (q_proj): Linear(in_features=2048, out_features=2048, bias=False)
9         (k_proj): Linear(in_features=2048, out_features=256, bias=False)
10        (v_proj): Linear(in_features=2048, out_features=256, bias=False)
11        (o_proj): Linear(in_features=2048, out_features=2048, bias=False)
12        (rotary_emb): LlamaRotaryEmbedding()
13      )
14      (mlp): LlamaMLP(
15        (gate_proj): Linear(in_features=2048, out_features=5632, bias=False)
16        (up_proj): Linear(in_features=2048, out_features=5632, bias=False)
17        (down_proj): Linear(in_features=5632, out_features=2048, bias=False)
18        (act_fn): SiLU()
19      )
20      (input_layernorm): LlamaRMSNorm((2048,), eps=1e-05)
21      (post_attention_layernorm): LlamaRMSNorm((2048,), eps=1e-05)
22    )
23  )
24  (norm): LlamaRMSNorm((2048,), eps=1e-05)
25  (rotary_emb): LlamaRotaryEmbedding()
26 )
```

---

## 4 Pruning (Kai Herbst)

Im Nachfolgenden sind die Befehle und Teile der Ausgabe aufgelistet, die für das Pruning des Basismodells verwendet wurde. Zu beachten ist, dass beim Pruning selbst in der Konsole bereits Angaben zu der verbleibenden Anzahl an Parametern generiert werden. Wie in der Hauptarbeit jedoch beschrieben, wurde die Anzahl an Parametern nach dem Pruning mit dem Evaluierungsskript im *LLM-Pruner* gemessen, das selbst die Anzahl an Parametern ausliest. Für die weitere Evaluierung wurden die Anzahl der Parameter verwendet, die durch das Evaluierungsskript generiert wurden, da dieses Skript anhand dieser Angabe die benötigte Rechenleistung berechnet.

Aufgelistet sind hier die verwendeten Befehle zum Prunen des Modells mit der Angabe das Modell um 30% mit den vier möglichen Methoden reduzieren. Für die 40% und 70% sind die Befehle identisch, außer, dass die *-pruning\_ratio* entsprechend auf 0.3 bzw. 0.7 gesetzt und der Speicherpfad angepasst wurde.

Pruning des Modells mit der Taylor-Methode:

---

```
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_0418_prune_log
12     --pruner_type taylor
13     --taylor param_first
14     --save_model
15     --max_seq_len 2048
16     --test_after_train
```

---

```
1  2025-01-04 15:57:27 - INFO :      We will use 90% of the memory on device 0
    for storing the model, and 10% for the buffer to avoid OOM. You can set
    'max_memory' in to a higher value to use more memory (at your own risk).
```

---

```

2 2025-01-04 15:57:29 - INFO :      Use taylor pruner...
3 2025-01-04 15:57:29 - INFO :      Pruning Attention Layer = [4, 5, 6, 7, 8, 9,
    10, 11, 12, 13, 14, 15, 16, 17]
4 2025-01-04 15:57:29 - INFO :      Pruning MLP Layer = [4, 5, 6, 7, 8, 9, 10,
    11, 12, 13, 14, 15, 16, 17]
5 2025-01-04 15:57:29 - INFO :      Start Pruning
6 2025-01-04 15:57:32 - INFO :      Start Backwarding in iterative steps = 0...
7 2025-01-04 15:57:32 - INFO :      Loss = 3.7723865509033203
8 2025-01-04 15:57:33 - INFO :      After Iter 1/1, #parameters: 921651200
9 2025-01-04 15:57:33 - INFO :      #Param before: 1100048384, #Param after:
    921651200, Ratio = 83.7828%
10 ...
11 2025-01-04 15:59:37 - INFO :      PPL after pruning: {'wikitext2':
    15.190999662372043, 'ptb': 43.19069160819034}
12 2025-01-04 15:59:37 - INFO :      Memory Requirement: 4130.3818359375 MiB

```

---

## Pruning des Modells mit der L1-Methode:

```

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_0418_l1_prune_log
12    --pruner_type l1
13    --save_model
14    --max_seq_len 2048
15    --test_after_train

```

---

```

1 2025-01-04 14:31:44 - INFO :      We will use 90% of the memory on device 0
    for storing the model, and 10% for the buffer to avoid OOM. You can set
    'max_memory' in to a higher value to use more memory (at your own risk).
2 2025-01-04 14:31:46 - INFO :      Use l1 pruner...
3 2025-01-04 14:31:46 - INFO :      Pruning Attention Layer = [4, 5, 6, 7, 8, 9,
    10, 11, 12, 13, 14, 15, 16, 17]
4 2025-01-04 14:31:46 - INFO :      Pruning MLP Layer = [4, 5, 6, 7, 8, 9, 10,
    11, 12, 13, 14, 15, 16, 17]
5 2025-01-04 14:31:46 - INFO :      Start Pruning
6 2025-01-04 14:31:47 - INFO :      After Iter 1/1, #parameters: 921651200
7 2025-01-04 14:31:47 - INFO :      #Param before: 1100048384, #Param after:
    921651200, Ratio = 83.7828%

```



```
8 ...
9 2025-01-04 14:35:14 - INFO :      PPL after pruning: {'wikitext2':
      281.63868765896893, 'ptb': 4992.165187767088}
10 2025-01-04 14:35:14 - INFO :      Memory Requirement: 1783.6884765625 MiB
```

---

## Pruning des Modells mit der L2-Methode:

---

```
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_0418_12_prune_log
12    --pruner_type l2
13    --save_model
14    --max_seq_len 2048
15    --test_after_train
```

---

```
1 2025-01-04 14:38:13 - INFO :      We will use 90% of the memory on device 0
      for storing the model, and 10% for the buffer to avoid OOM. You can set
      'max_memory' in to a higher value to use more memory (at your own risk).
2 2025-01-04 14:38:15 - INFO :      Use l2 pruner...
3 2025-01-04 14:38:15 - INFO :      Pruning Attention Layer = [4, 5, 6, 7, 8, 9,
      10, 11, 12, 13, 14, 15, 16, 17]
4 2025-01-04 14:38:15 - INFO :      Pruning MLP Layer = [4, 5, 6, 7, 8, 9, 10,
      11, 12, 13, 14, 15, 16, 17]
5 2025-01-04 14:38:16 - INFO :      Start Pruning
6 2025-01-04 14:38:16 - INFO :      After Iter 1/1, #parameters: 921651200
7 2025-01-04 14:38:16 - INFO :      #Param before: 1100048384, #Param after:
      921651200, Ratio = 83.7828%
8 ...
9 2025-01-04 14:41:14 - INFO :      PPL after pruning: {'wikitext2':
      42.190180314697045, 'ptb': 167.51849280886478}
10 2025-01-04 14:41:14 - INFO :      Memory Requirement: 1783.6904296875 MiB
```

---

## Pruning des Modells mit der Random-Methode:

---

```
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_0418_random_prune_log
12         --pruner_type random
13         --save_model
14         --max_seq_len 2048
15         --test_after_train

```

---

```

1 2025-01-04 15:50:55 - INFO :      We will use 90% of the memory on device 0
    for storing the model, and 10% for the buffer to avoid OOM. You can set
    'max_memory' in to a higher value to use more memory (at your own risk).
2 2025-01-04 15:50:57 - INFO :      Use random pruner...
3 2025-01-04 15:50:57 - INFO :      Pruning Attention Layer = [4, 5, 6, 7, 8, 9,
    10, 11, 12, 13, 14, 15, 16, 17]
4 2025-01-04 15:50:57 - INFO :      Pruning MLP Layer = [4, 5, 6, 7, 8, 9, 10,
    11, 12, 13, 14, 15, 16, 17]
5 2025-01-04 15:50:57 - INFO :      Start Pruning
6 2025-01-04 15:50:57 - INFO :      After Iter 1/1, #parameters: 921651200
7 2025-01-04 15:50:57 - INFO :      #Param before: 1100048384, #Param after:
    921651200, Ratio = 83.7828%
8 ...
9 2025-01-04 15:53:04 - INFO :      PPL after pruning: {'wikitext2':
    40.892124929554, 'ptb': 166.21485355527523}
10 2025-01-04 15:53:04 - INFO :      Memory Requirement: 1793.31396484375 MiB

```

---

## 5 Performance-Evaluierung (Kai Herbst)

Auch hier wurde für die entsprechenden Evaluierungen die nahezu identischen Aufrufe verwendet. Im folgenden aufgelistet sind daher der Befehl zur Evaluierung des Basismodells und der Befehl zur Evaluierung eines geprüften Modells.

---

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

---

Dazugehöriger relevanter Ausschnitt der Konsolenausgabe:

---

```
1 {
2   "results": {
3     "arc_challenge": {
4       "acc": 0.30119453924914674,
5       "acc_stderr": 0.01340674176784762,
6       "acc_norm": 0.3250853242320819,
7       "acc_norm_stderr": 0.013688147309729119
8     },
9     "boolq": {
10      "acc": 0.6131498470948012,
11      "acc_stderr": 0.00851818834084475
12    },
13    "openbookqa": {
14      "acc": 0.242,
15      "acc_stderr": 0.019173085678337164,
16      "acc_norm": 0.376,
17      "acc_norm_stderr": 0.021683827539286125
18    },
19    "winogrande": {
20      "acc": 0.6029992107340174,
21      "acc_stderr": 0.013751092519806702
22    },
23    "hellaswag": {
24      "acc": 0.4614618601872137,
25      "acc_stderr": 0.004974937803907464,
26      "acc_norm": 0.5827524397530373,
27      "acc_norm_stderr": 0.004920967192255295
```

```

28     }
29 },
30 "versions": {
31     "arc_challenge": 0,
32     "boolq": 1,
33     "openbookqa": 0,
34     "winogrande": 0,
35     "hellaswag": 0
36 },
37 "config": {
38     "model": "hf-causal-experimental",
39     "model_args": "pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0",
40     "num_fewshot": 0,
41     "batch_size": null,
42     "device": "cuda:0",
43     "no_cache": true,
44     "limit": null,
45     "bootstrap_iters": 100000,
46     "description_dict": {}
47 }
48 }
49 hf-causal-experimental (pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0), limit:
      None, provide_description: False, num_fewshot: 0, batch_size: None
50 |   Task      | Version | Metric | Value | | Stderr |
51 |-----|-----|-----|-----|---|-----|
52 | arc_challenge |      0 | acc     | 0.3012 | +- | 0.0134 |
53 |               |        | acc_norm | 0.3251 | +- | 0.0137 |
54 | boolq         |      1 | acc     | 0.6131 | +- | 0.0085 |
55 | openbookqa    |      0 | acc     | 0.2420 | +- | 0.0192 |
56 |               |        | acc_norm | 0.3760 | +- | 0.0217 |
57 | winogrande    |      0 | acc     | 0.6030 | +- | 0.0138 |
58 | hellaswag     |      0 | acc     | 0.4615 | +- | 0.0050 |
59 |               |        | acc_norm | 0.5828 | +- | 0.0049 |

```

---

Aufruf zur Evaluierung eines geprunten Modells:

```

1  $ export PYTHONPATH='.'
2  $ python lm-evaluation-harness/main.py
3      --model hf-causal-experimental
4      --model_args
5          checkpoint=prune_log/tinyllama_30_0418_prune_log/pytorch_model.bin,
6          config_pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0
7      --tasks
9          openbookqa,arc_easy,winogrande,hellaswag,arc_challenge,piqa,boolq
8      --device cuda:0
9      --no_cache

```

10           --output\_path results/tinyllama\_30\_0418

---

```
1  {
2    "results": {
3      "hellaswag": {
4        "acc": 0.379008165704043,
5        "acc_stderr": 0.00484148671685576,
6        "acc_norm": 0.48376817367058356,
7        "acc_norm_stderr": 0.004987151381091178
8      },
9      "arc_challenge": {
10       "acc": 0.24658703071672355,
11       "acc_stderr": 0.012595726268790125,
12       "acc_norm": 0.29180887372013653,
13       "acc_norm_stderr": 0.013284525292403511
14     },
15     "winogrande": {
16       "acc": 0.5422257300710339,
17       "acc_stderr": 0.014002284504422438
18     },
19     "openbookqa": {
20       "acc": 0.232,
21       "acc_stderr": 0.01889619359195206,
22       "acc_norm": 0.362,
23       "acc_norm_stderr": 0.021513662527582404
24     },
25     "boolq": {
26       "acc": 0.555045871559633,
27       "acc_stderr": 0.008691897543539214
28     }
29   },
30   "versions": {
31     "hellaswag": 0,
32     "arc_challenge": 0,
33     "winogrande": 0,
34     "openbookqa": 0,
35     "boolq": 1
36   },
37   "config": {
38     "model": "hf-causal-experimental",
39     "model_args":
40       "checkpoint=prune_log/tinyllama_10_0616_prune_log/pytorch_model.bin,
41       config_pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0",
42     "num_fewshot": 0,
43     "batch_size": null,
44     "device": "cuda:0",
```

```

44     "no_cache": true,
45     "limit": null,
46     "bootstrap_iters": 100000,
47     "description_dict": {}
48 }
49 }
50 hf-causal-experimental
    (checkpoint=prune_log/tinyllama_30_0418_prune_log/pytorch_model.bin,
51 config_pretrained=TinyLlama/TinyLlama-1.1B-Chat-v1.0), limit: None,
    provide_description: False, num_fewshot: 0, batch_size: None
52 |   Task      | Version | Metric | Value | | Stderr |
53 |-----|-----:|-----|-----:|---|-----:|
54 | hellaswag   |        | 0|acc   | 0.3790|+- | 0.0048|
55 |             |        | |acc_norm| 0.4838|+- | 0.0050|
56 | arc_challenge |       | 0|acc   | 0.2466|+- | 0.0126|
57 |             |        | |acc_norm| 0.2918|+- | 0.0133|
58 | winogrande  |        | 0|acc   | 0.5422|+- | 0.0140|
59 | openbookqa  |        | 0|acc   | 0.2320|+- | 0.0189|
60 |             |        | |acc_norm| 0.3620|+- | 0.0215|
61 | boolq       |        | 1|acc   | 0.5550|+- | 0.0087|

```

---

## 6 Tuning (Kai Herbst)

Wie in der Hausarbeit beschrieben, wurde eines der Modelle nach dem Pruning nachtrainiert. Dazu wurde folgender Befehl verwendet. Damit die Verbindung mit Wandb funktioniert, muss sich entsprechend vorher ein Account erstellt werden und sich mit diesem über das Terminal angemeldet werden.

---

```
1  python post_training.py --prune_model
    prune_log/tinyllama_30_0418_prune_log/pytorch_model.bin
2  --data_path yahma/alpaca-cleaned
3  --lora_r 8
4  --num_epochs 2
5  --learning_rate 1e-4
6  --batch_size 64
7  --output_dir tune_log/tl30
8  --wandb_project llama_tune
```

---

## 7 Evaluierung Rechenanforderung (Kai Herbst)

Im folgenden Text sind erneut die Befehle für die Evaluierung des Basismodells aufgelistet und für eines der 30% Modelle sowie den dazugehörigen Output im Terminal.

Evaulierung des Basismodells:

---

```
1 $ python test_speedup.py
2   --model_type pretrain
3   --base_model TinyLlama/TinyLlama-1.1B-Chat-v1.0
```

---

Dazugehörige Ausgabe:

---

```
1 LlamaForCausalLM(
2   1195.99 M, 94.805% Params, 76.94 GMac, 99.509% MACs,
3   (model): LlamaModel(
4     1130.46 M, 89.610% Params, 72.75 GMac, 94.084% MACs,
5     (embed_tokens): Embedding(0, 0.000% Params, 0.0 Mac, 0.000% MACs, 32000,
6       2048)
7     (layers): ModuleList(
8       (0-21): 22 x LlamaDecoderLayer(
9         (self_attn): LlamaAttention(
10          16.78 M, 1.330% Params, 1.09 GMac, 1.411% MACs,
11          (q_proj): Linear(4.19 M, 0.332% Params, 268.44 MMac, 0.347% MACs,
12            in_features=2048, out_features=2048, bias=False)
13          (k_proj): Linear(4.19 M, 0.332% Params, 268.44 MMac, 0.347% MACs,
14            in_features=2048, out_features=2048, bias=False)
15          (v_proj): Linear(4.19 M, 0.332% Params, 268.44 MMac, 0.347% MACs,
16            in_features=2048, out_features=2048, bias=False)
17          (o_proj): Linear(4.19 M, 0.332% Params, 268.44 MMac, 0.347% MACs,
18            in_features=2048, out_features=2048, bias=False)
19          (rotary_emb): LlamaRotaryEmbedding(0, 0.000% Params, 0.0 Mac, 0.000%
20            MACs, )
21        )
22      (mlp): LlamaMLP(
23        34.6 M, 2.743% Params, 2.21 GMac, 2.865% MACs,
24        (gate_proj): Linear(11.53 M, 0.914% Params, 738.2 MMac, 0.955% MACs,
25          in_features=2048, out_features=5632, bias=False)
26        (down_proj): Linear(11.53 M, 0.914% Params, 738.2 MMac, 0.955% MACs,
27          in_features=5632, out_features=2048, bias=False)
28        (up_proj): Linear(11.53 M, 0.914% Params, 738.2 MMac, 0.955% MACs,
29          in_features=2048, out_features=5632, bias=False)
```



```

22         (act_fn): SiLU(0, 0.000% Params, 360.45 KMac, 0.000% MACs, )
23     )
24     (input_layernorm): LlamaRMSNorm(2.05 k, 0.000% Params, 262.14 KMac,
        0.000% MACs, )
25     (post_attention_layernorm): LlamaRMSNorm(2.05 k, 0.000% Params, 262.14
        KMac, 0.000% MACs, )
26 )
27 )
28     (norm): LlamaRMSNorm(2.05 k, 0.000% Params, 262.14 KMac, 0.000% MACs, )
29 )
30     (lm_head): Linear(65.54 M, 5.195% Params, 4.19 GMac, 5.425% MACs,
        in_features=2048, out_features=32000, bias=False)
31 )
32 Computational complexity:          77.32 GMac
33 Number of parameters:              1261.53 M
34 GPU Memory Requirement: 2427.3115234375 MiB

```

---

```

1     $ python test_speedup.py
2         --model_type pruneLLM
3         --ckpt prune_log/tinyllama_30_0418_prune_log/pytorch_model.bin

```

---

## Dazugehörige Ausgabe:

```

1 LlamaForCausalLM(
2     856.02 M, 92.879% Params, 54.79 GMac, 99.988% MACs,
3     (model): LlamaModel(
4         790.49 M, 85.769% Params, 50.6 GMac, 92.334% MACs,
5         (embed_tokens): Embedding(0, 0.000% Params, 0.0 Mac, 0.000% MACs, 32000,
            2048)
6         (layers): ModuleList(
7             (0-3): 4 x LlamaDecoderLayer(
8                 44.04 M, 4.778% Params, 2.82 GMac, 5.144% MACs,
9                 (self_attn): LlamaAttention(
10                    9.44 M, 1.024% Params, 603.98 MMac, 1.102% MACs,
11                    (q_proj): Linear(4.19 M, 0.455% Params, 268.44 MMac, 0.490% MACs,
                        in_features=2048, out_features=2048, bias=False)
12                    (k_proj): Linear(524.29 k, 0.057% Params, 33.55 MMac, 0.061% MACs,
                        in_features=2048, out_features=256, bias=False)
13                    (v_proj): Linear(524.29 k, 0.057% Params, 33.55 MMac, 0.061% MACs,
                        in_features=2048, out_features=256, bias=False)
14                    (o_proj): Linear(4.19 M, 0.455% Params, 268.44 MMac, 0.490% MACs,
                        in_features=2048, out_features=2048, bias=False)
15                )
16            (mlp): LlamaMLP(
17                34.6 M, 3.754% Params, 2.21 GMac, 4.042% MACs,

```

```

18         (gate_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
19             in_features=2048, out_features=5632, bias=False)
20         (up_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
21             in_features=2048, out_features=5632, bias=False)
22         (down_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
23             in_features=5632, out_features=2048, bias=False)
24         (act_fn): SiLU(0, 0.000% Params, 360.45 KMac, 0.001% MACs, )
25     )
26     (input_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac, 0.000% MACs,
27         (2048,), eps=1e-05)
28     (post_attention_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac,
29         0.000% MACs, (2048,), eps=1e-05)
30 )
31 (4-17): 14 x LlamaDecoderLayer(
32     31.3 M, 3.396% Params, 2.0 GMac, 3.656% MACs,
33     (self_attn): LlamaAttention(
34         7.08 M, 0.768% Params, 452.98 MMac, 0.827% MACs,
35         (q_proj): Linear(3.15 M, 0.341% Params, 201.33 MMac, 0.367% MACs,
36             in_features=2048, out_features=1536, bias=False)
37         (k_proj): Linear(393.22 k, 0.043% Params, 25.17 MMac, 0.046% MACs,
38             in_features=2048, out_features=192, bias=False)
39         (v_proj): Linear(393.22 k, 0.043% Params, 25.17 MMac, 0.046% MACs,
40             in_features=2048, out_features=192, bias=False)
41         (o_proj): Linear(3.15 M, 0.341% Params, 201.33 MMac, 0.367% MACs,
42             in_features=1536, out_features=2048, bias=False)
43     )
44     (mlp): LlamaMLP(
45         24.22 M, 2.628% Params, 1.55 GMac, 2.829% MACs,
46         (gate_proj): Linear(8.07 M, 0.876% Params, 516.69 MMac, 0.943% MACs,
47             in_features=2048, out_features=3942, bias=False)
48         (up_proj): Linear(8.07 M, 0.876% Params, 516.69 MMac, 0.943% MACs,
49             in_features=2048, out_features=3942, bias=False)
50         (down_proj): Linear(8.07 M, 0.876% Params, 516.69 MMac, 0.943% MACs,
51             in_features=3942, out_features=2048, bias=False)
52         (act_fn): SiLU(0, 0.000% Params, 252.29 KMac, 0.000% MACs, )
53     )
54     (input_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac, 0.000% MACs,
55         (2048,), eps=1e-05)
56     (post_attention_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac,
57         0.000% MACs, (2048,), eps=1e-05)
58 )
59 (18-21): 4 x LlamaDecoderLayer(
60     44.04 M, 4.778% Params, 2.82 GMac, 5.144% MACs,
61     (self_attn): LlamaAttention(
62         9.44 M, 1.024% Params, 603.98 MMac, 1.102% MACs,
63         (q_proj): Linear(4.19 M, 0.455% Params, 268.44 MMac, 0.490% MACs,
64             in_features=2048, out_features=2048, bias=False)

```

```

50         (k_proj): Linear(524.29 k, 0.057% Params, 33.55 MMac, 0.061% MACs,
           in_features=2048, out_features=256, bias=False)
51         (v_proj): Linear(524.29 k, 0.057% Params, 33.55 MMac, 0.061% MACs,
           in_features=2048, out_features=256, bias=False)
52         (o_proj): Linear(4.19 M, 0.455% Params, 268.44 MMac, 0.490% MACs,
           in_features=2048, out_features=2048, bias=False)
53     )
54     (mlp): LlamaMLP(
55         34.6 M, 3.754% Params, 2.21 GMac, 4.042% MACs,
56         (gate_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
           in_features=2048, out_features=5632, bias=False)
57         (up_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
           in_features=2048, out_features=5632, bias=False)
58         (down_proj): Linear(11.53 M, 1.251% Params, 738.2 MMac, 1.347% MACs,
           in_features=5632, out_features=2048, bias=False)
59         (act_fn): SiLU(0, 0.000% Params, 360.45 KMac, 0.001% MACs, )
60     )
61     (input_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac, 0.000% MACs,
           (2048,)), eps=1e-05)
62     (post_attention_layernorm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac,
           0.000% MACs, (2048,)), eps=1e-05)
63 )
64 )
65 (norm): LlamaRMSNorm(0, 0.000% Params, 0.0 Mac, 0.000% MACs, (2048,)),
           eps=1e-05)
66 (rotary_emb): LlamaRotaryEmbedding(0, 0.000% Params, 0.0 Mac, 0.000% MACs, )
67 )
68 (lm_head): Linear(65.54 M, 7.111% Params, 4.19 GMac, 7.654% MACs,
           in_features=2048, out_features=32000, bias=False)
69 )
70 Computational complexity:      54.8 GMac
71 Number of parameters:         921.65 M
72 GPU Memory Requirement: 1793.30126953125 MiB

```

---

## 8 Fazit (Sebastian Viet)

Mittlerweile gibt es verschiedene Packages um Large Language Models zu prunen. Leider funktionieren viele dieser Packages nicht einwandfrei und sind teilweise auf bestimmte Modelle zugeschnitten/eingeschränkt. Mit Hilfe des *LLM-Pruners* lassen sich zumindest die Llama-Modelle vergleichsweise einfach prunen, wenn auch hier zahlreiche Komplikationen in Bezug auf Package-Versionen entstanden sind, die zu verschiedenen Fehlermeldungen geführt haben. Dennoch lässt sich, unter der Verwendung der richtigen Package-versionen, das TinyLlama-Modell unkompliziert prunen und anschließend auf verschiedene Aspekte hin evaluieren.