

# Projektbericht: IMARA

---

## Domain-specific GraphRAG pipeline with model fine-tuning

**Modul:** Abschlussarbeit CAS Machine Learning for Software Engineers (ML4SE)

**Datum:** [Aktuelles Datum]

**Autoren:** Marco Allenspach, Lukas Koller, Emanuel Sovrano

### Abstract

Die Einführung von Retrieval-Augmented Generation (RAG) markierte einen bedeutenden Meilenstein in der Anwendung grosser Sprachmodelle (LLM), indem generative Fähigkeiten auf faktischen, externen Daten basierten, um Fehlinterpretationen zu vermeiden und die Relevanz zu erhöhen. Um die Schwächen von RAG der ersten Generation durch die Einführung strukturierter, relationaler Kontexte zu beheben, hat sich jedoch mit AI-Native GraphRAG ein weiterentwickeltes Paradigma etabliert.

Der rasante branchenweite Wandel hin zu graphenbasierten Architekturen ist eine notwendige Weiterentwicklung, die auf der Erkenntnis beruht, dass eine KI für effektives Denken ein Modell des Anwendungsbereichs benötigt, nicht nur eine Sammlung von Fakten. Der Fortschritt von unreflektierten LLMs zu grundlegenden RAGs löste das Problem der faktischen Fundierung, doch das Versagen rein vektorbasierter RAGs bei komplexen Anfragen zeigte, dass die Struktur des Wissens ebenso wichtig ist wie sein Inhalt. Ein Wissensgraph liefert diese Struktur und transformiert eine passive Dokumentensammlung in ein aktives, abfragefähiges Modell der Welt.

---

### 1. Management Summary

(Ca. 0.5 - 1 Seite) Zusammenfassung des gesamten Projekts: Problemstellung (Extraktion aus komplexen PDFs), gewählter Lösungsansatz (GraphRAG & Fine-tuning) und die wichtigsten Ergebnisse des Benchmarkings.

Traditionelle neuronale Netze eignen sich gut zur Kodierung linearer Beziehungen, doch Daten aus der realen Welt sind in der Regel komplex und multidimensional. Graphen sind besser geeignet, höherdimensionale Verbindungen darzustellen, in denen jeder Knoten mit jedem anderen Knoten in Beziehung steht. Dadurch eignen sich Graphen besser zur Speicherung komplexer Beziehungen aus der realen Welt.

### 2. Einleitung und Zielsetzung

Das IMARA-Projekt hat zum Ziel, aufzuzeigen wie die Genauigkeit der Abfrage eines graph-basierten RAG-Systems sich verbessert.

Um eine Grundlage für die Messbarkeit zu haben, wurde OpenRAGBench als Referenzdatensatz ausgewählt.

#### Defining the "AI-Native" GraphRAG Paradigm

AI-Native GraphRAG represents a specific and powerful subset of graph-based RAG systems. Solutions must automate the entire workflow from unstructured data to a natural language answer, abstracting complexities

of graph theory and database management.

## naives RAG

### Die inhärenten Einschränkungen von vektorbasierter RAG

Konventionelle RAG-Architekturen verlassen sich auf Vektorsimilaritätssuche über ein Korpus von geteiltem Text. Dieser Ansatz behandelt Wissen als eine Sammlung von unzusammenhängenden Fakten und hat Schwierigkeiten mit Fragen, die erfordern:

- Synthese von Informationen aus mehreren Quellen
- Verständnis nuancierter Beziehungen zwischen Entitäten
- Durchführung von Multi-Hop-Reasoning Der Kontext, der dem LLM bereitgestellt wird, ist oft eine Liste von Textausschnitten, die keine explizite Darstellung ihrer Verbindungen enthalten.

### Kontextuelle Fragmentierung und Blindheit

Das Chunking bricht den natürlichen Informationsfluss willkürlich. Relevanter Kontext kann über verschiedene Chunks, Dokumente oder Abschnitte verstreut sein. Die Vektorschre, die die Anfrage mit jedem Chunk einzeln vergleicht, versagt oft dabei, diesen vollständigen, verteilten Kontext abzurufen, was zu unvollständigen oder oberflächlichen Antworten führt. Sie versteht semantische Ähnlichkeit, ist jedoch blind für explizite Beziehungen wie Kausalität, Abhängigkeit oder Hierarchie.

### Empfindlichkeit gegenüber der Chunking-Strategie

Die Leistung ist hochgradig empfindlich gegenüber der Chunking-Strategie (z.B. Chunk-Grösse, Überlappung). Suboptimale Strategien können übermässiges Rauschen einführen (Chunks zu gross) oder kritischen Kontext verlieren (Chunks zu klein), was umfangreiche und brüchige Anpassungen erfordert.

### Unfähigkeit, Multi-Hop-Reasoning durchzuführen

Es gibt Schwierigkeiten, komplexe Fragen zu beantworten, die "Multi-Hop"-Reasoning erfordern. Zum Beispiel: "Welche Marketingkampagnen wurden von der in dem Q3-Bericht erwähnten Lieferkettenstörung betroffen?" erfordert die Verknüpfung von Störung → betroffene Produkte → Marketingkampagnen. Eine einfache Vektorschre ist unwahrscheinlich, diese Informationssprünge zu überbrücken.

**Analogie:** Vektorbasierte RAG bietet einem Forscher einen Stapel isolierter Karteikarten, während GraphRAG darauf abzielt, eine umfassende Mindmap zu erstellen und bereitzustellen, die entscheidende Verbindungen aufdeckt.

## 2.1 Projekttitel: IMARA

## 2.2 Problemstellung

Die Extraktion und Verarbeitung von Informationen aus unstrukturierten PDF-Dokumenten stellt eine Herausforderung für herkömmliche RAG-Systeme dar.

## 2.3 Projektziele

- Die Implementation von graphbasierten System und der Vergleich zu klassischen RAG-Systemen
- Der Vergleich zwischen verschiedenen graphbasierten RAG-Systemen

**Graph-basiertes RAG:** Aufbau einer Pipeline zur Erstellung dichter Wissensgraphen. In diesem Projekt werden drei ausprägungen deines graphbasierten RAG-Systems umgesetzt.

1. **LeanRAG** - basiert auf einem Knowledge Graph mit einer hierarchischen Aggregierung.
  2. **linearRAG** - relying on lightweight entity recognition and semantic linking
  3. **GraphMERT** - kompaktes, rein grafisches Encoder-Modell, das hochwertige KGs aus unstrukturierten Textkorpora und seinen eigenen internen Repräsentationen generiert
- 

**Model Fine-tuning:** Optimierung eines LLMs (z.B. Qwen) basierend auf dem Graph.

- 

**Automation:** End-to-End Automatisierung der Pipeline. Eine flexible Pipeline bauen, die bei der Evaluation der verschiedenen RAG-Systeme unterstützt.

### 3. Datenbasis und Vorverarbeitung

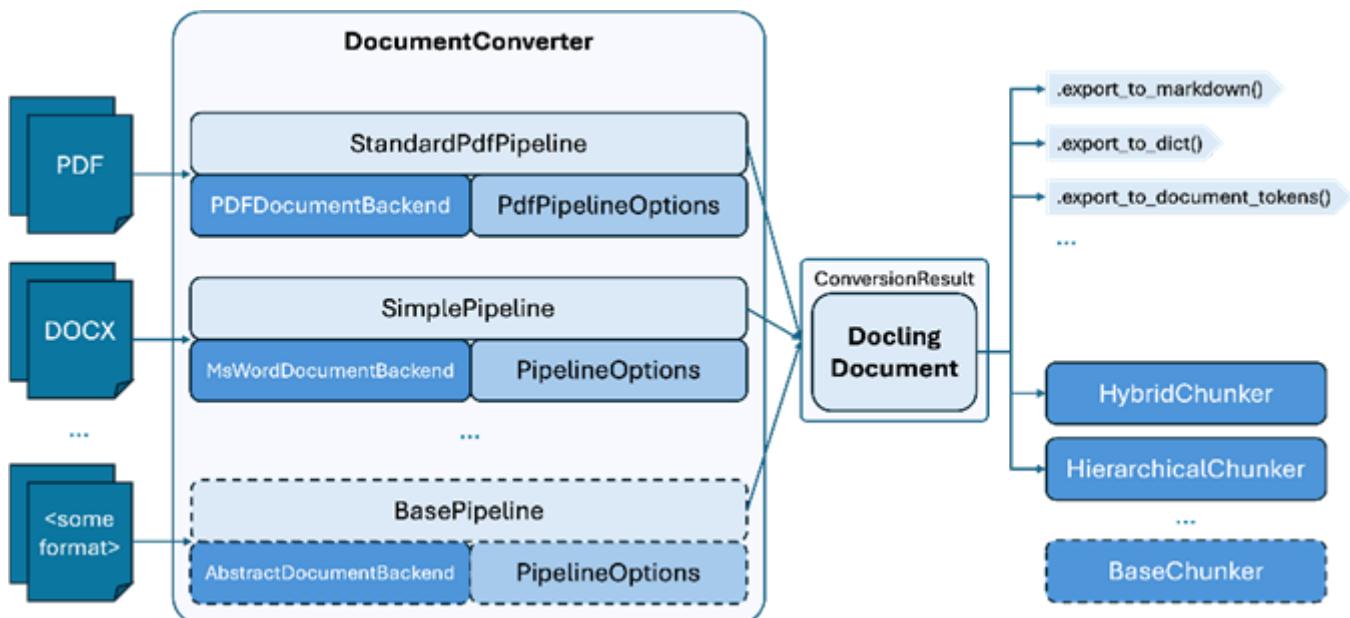
#### 3.1 Datenquellen

Beschreibung der verwendeten Datensätze, wie z.B. der **Open RAG Bench Dataset** (Arxiv-Kategorien) oder **PubMedQA**.

### 4. Methodik und Architektur

#### 4.1 PDF-Extraktion mit Docling

Einsatz des **Docling Toolkits** zur effizienten Konvertierung von Dokumenten in maschinenlesbare Formate (Markdown/JSON).



#### 4.2 Graph-Konstruktion

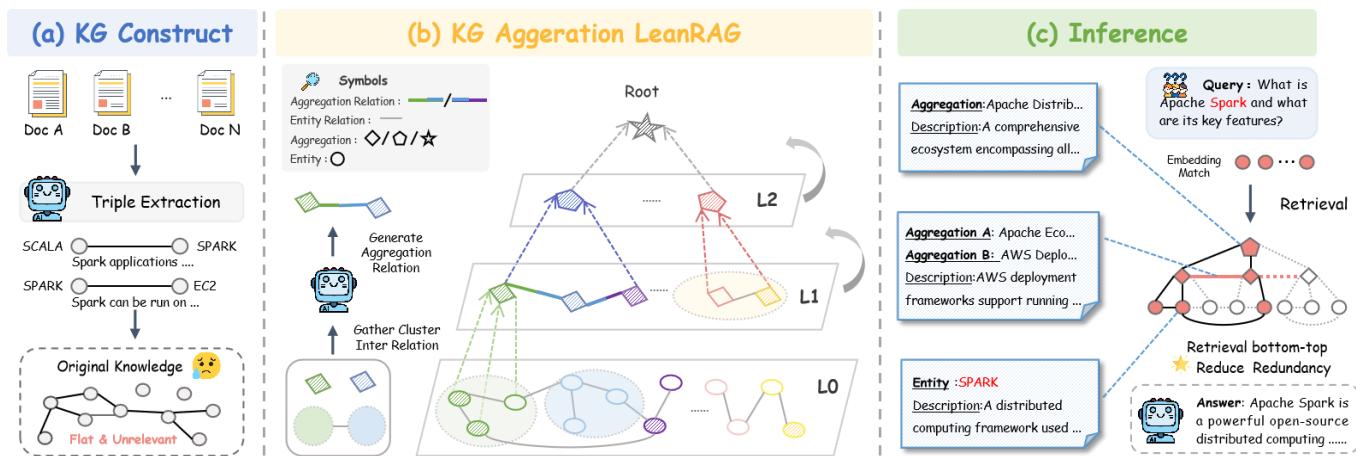
##### 4.1.1 LeanRAG Ansatz

Detaillierung der Triple-Extraktion und der hierarchischen Retrieval-Struktur.

## ◆ Features

- **Semantic Aggregation:** Clusters entities into semantically coherent summaries and constructs explicit relations to form a navigable aggregation-level knowledge network.
- **Hierarchical, Structure-Guided Retrieval:** Initiates retrieval from fine-grained entities and traverses up the knowledge graph to gather rich, highly relevant evidence efficiently.
- **Reduced Redundancy:** Optimizes retrieval paths to significantly reduce redundant information—LeanRAG achieves ~46% lower retrieval redundancy compared to flat retrieval baselines (based on benchmark evaluations).
- **Benchmark Performance:** Demonstrates superior performance across multiple QA benchmarks with improved response quality and retrieval efficiency.

## 🏛️ Architecture Overview



LeanRAG's processing pipeline follows these core stages:

### 1. Semantic Aggregation

- Group low-level entities into clusters; generate summary nodes and build adjacency relations among them for efficient navigation.

### 2. Knowledge Graph Construction

- Construct a multi-layer graph where nodes represent entities and aggregated summaries, with explicit inter-node relations for graph-based traversal.

### 3. Query Processing & Hierarchical Retrieval

- Anchor queries at the most relevant detailed entities ("bottom-up"), then traverse upward through the semantic aggregation graph to collect evidence spans.

### 4. Redundancy-Aware Synthesis

- Streamline retrieval paths and avoid overlapping content, ensuring concise evidence aggregation before generating responses.

### 5. Generation

- Use retrieved, well-structured evidence as input to an LLM to produce coherent, accurate, and contextually grounded answers.

- **Extraktion:** Umwandlung von Text in Entitäten und Relationen.

## leanRAG Workflow

```
file_chunk.py
```

1. chunk raw input token-based with 512 Tokens and 64 Tokens overlap

### Method 1: CommonKG

```
CommonKG/create_kg.py
```

2. create a list of match words (entities) for each chunk
3. create a list of "all entities" based on the match words without duplicates
4. "new triples" have "subject, predicate, object" triples init with corresponding reference to the chunk of origin
5. "next layer entities"
6. "new triples descriptions"

```
CommonKG/deal_triple.py
```

7. summarize descriptions => relation.jsonl

### Method 2: GraphRAG

```
GraphExtraction/chunk.py
```

2. loads the chunks
3. performs a "triple extraction" => entity.jsonl, relation.jsonl

```
GraphExtraction/deal_triple.py
```

4. deal with duplicates of entries and relations

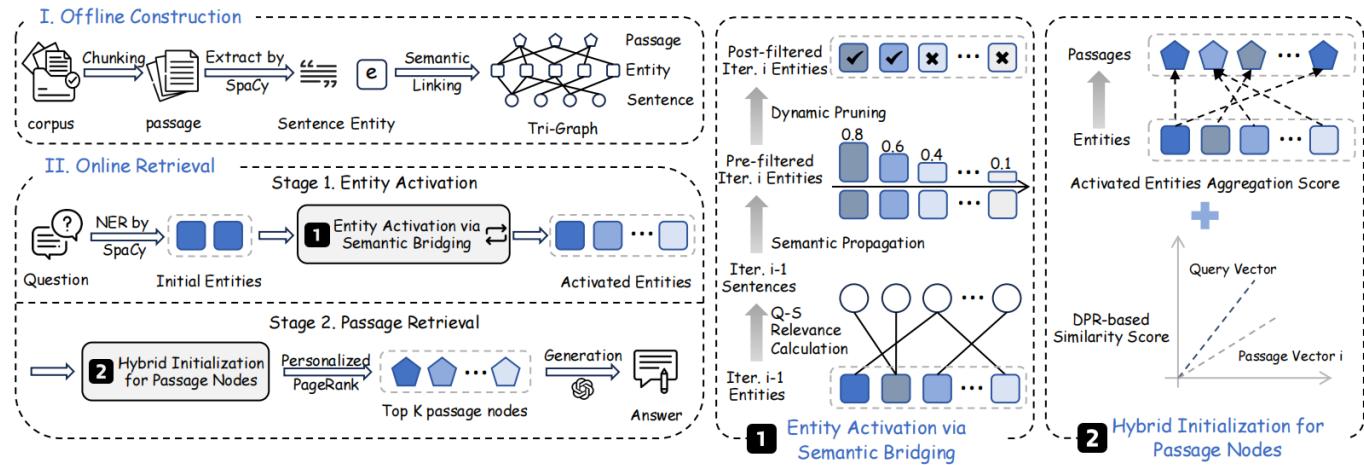
```
build_graph.py
```

5. generating embeddings
6. clustering labels (based on the embeddings)
7. layer 1 clustering
8. layer 2 clustering

## 9. building vector DB

### 4.1.2 LinearRAG

LinearRAG: Linear Graph Retrieval-Augmented Generation on Large-scale Corpora - A relation-free graph construction method for efficient GraphRAG.



- Context-Preserving: Relation-free graph construction, relying on lightweight entity recognition and semantic linking to achieve comprehensive contextual comprehension.
- Complex Reasoning: Enables deep retrieval via semantic bridging, achieving multi-hop reasoning in a single retrieval pass without requiring explicit relational graphs.
- High Scalability: Zero LLM token consumption, faster processing speed, and linear time/space complexity.

### Graphbuilding:

1. => load data
2. chunking data
3. get named entities - SpacyNER (Named Entity Recognition)
4. sentence splitting
5. get passages
6. get embeddings(sentences, entities, passages)
7. build graph => LinearRAG.graphml => ner\_results.json => passage\_embedding.parquet => dentence\_embedding.parquet => entity\_embedding.parquet

### Retreival:

1. retrieval\_results = qa(question)

### 4.1.3 GraphMERT

GraphMERT: Effiziente und skalierbare Gewinnung zuverlässiger Wissensgraphen aus unstrukturierten Daten

Ein einfaches Beispiel für eine Testimplementierung des Princeton GraphMERT-Papers.

<https://arxiv.org/abs/2510.09580>

Seit fast drei Jahrzehnten erforschen Wissenschaftler Anwendungen neurosymbolischer künstlicher Intelligenz (KI), da symbolische Komponenten Abstraktion und neuronale Komponenten Generalisierung ermöglichen. Die Kombination beider Komponenten verspricht rasante Fortschritte in der KI. Dieses Potenzial konnte das

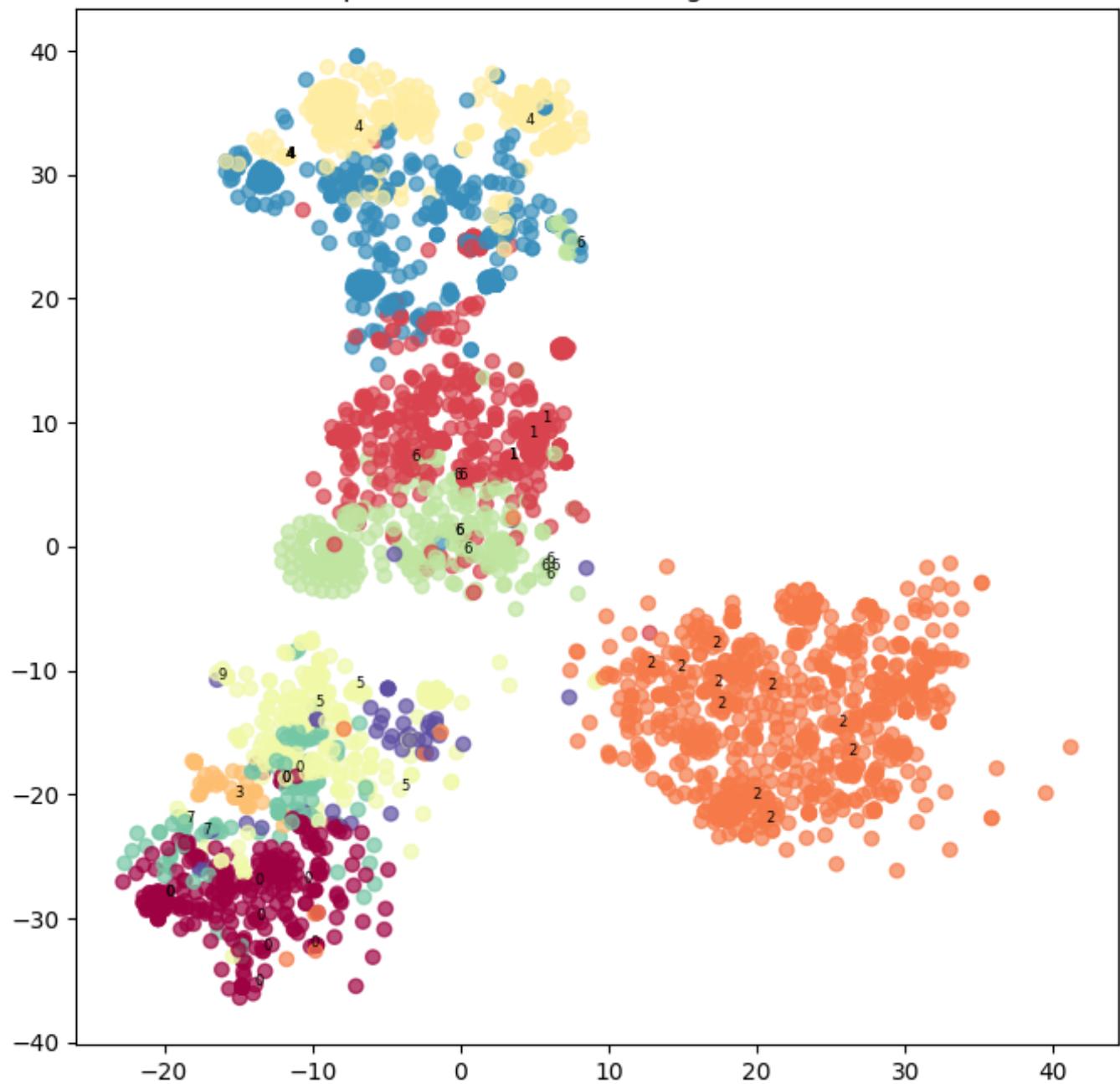
Feld jedoch bisher nicht ausschöpfen, da die meisten neurosymbolischen KI-Frameworks nicht skalierbar sind. Zudem schränken die impliziten Repräsentationen und das approximative Schliessen neuronaler Ansätze Interpretierbarkeit und Vertrauen ein. Wissensgraphen (KGs), die als Goldstandard für die Repräsentation expliziten semantischen Wissens gelten, können die symbolische Seite abdecken. Die automatische Ableitung zuverlässiger KGs aus Textkorpora stellt jedoch weiterhin eine Herausforderung dar. Wir begegnen diesen Herausforderungen mit GraphMERT, einem kompakten, rein grafischen Encoder-Modell, das hochwertige KGs aus unstrukturierten Textkorpora und seinen eigenen internen Repräsentationen generiert.

GraphMERT und sein äquivalenter Wissensgraph bilden einen modularen neurosymbolischen Stack: neuronales Lernen von Abstraktionen; symbolische Wissensgraphen für verifizierbares Schliessen. GraphMERT + Wissensgraph ist das erste effiziente und skalierbare neurosymbolische Modell, das höchste Benchmark-Genauigkeit und überlegene symbolische Repräsentationen im Vergleich zu Basismodellen erzielt.

Konkret streben wir zuverlässige domänenspezifische Wissensgraphen (KGs) an, die sowohl (1) faktisch korrekt (mit Herkunftsachweis) als auch (2) valide (ontologiekonsistente Relationen mit domänenspezifischer Semantik) sind. Wenn ein grosses Sprachmodell (LLM), z. B. Qwen3-32B, domänenspezifische KGs generiert, weist es aufgrund seiner hohen Sensitivität, seiner geringen Domänenexpertise und fehlerhafter Relationen Defizite in der Zuverlässigkeit auf. Anhand von Texten aus PubMed-Artikeln zum Thema Diabetes erzielt unser GraphMERT-Modell mit 80 Millionen Parametern einen KG mit einem FActScore von 69,8 %; ein LLM-Basismodell mit 32 Milliarden Parametern erreicht hingegen nur einen FActScore von 40,2 %. Der GraphMERT-KG erzielt zudem einen höheren ValidityScore von 68,8 % gegenüber 43,0 % beim LLM-Basismodell.

### **GraphMERT Node Embeddings (t-SNE View)**

### GraphMERT Node Embeddings (t-SNE View)



### GraphMERT Semantic Graph Visualization

## GraphMERT Semantic Graph Visualization



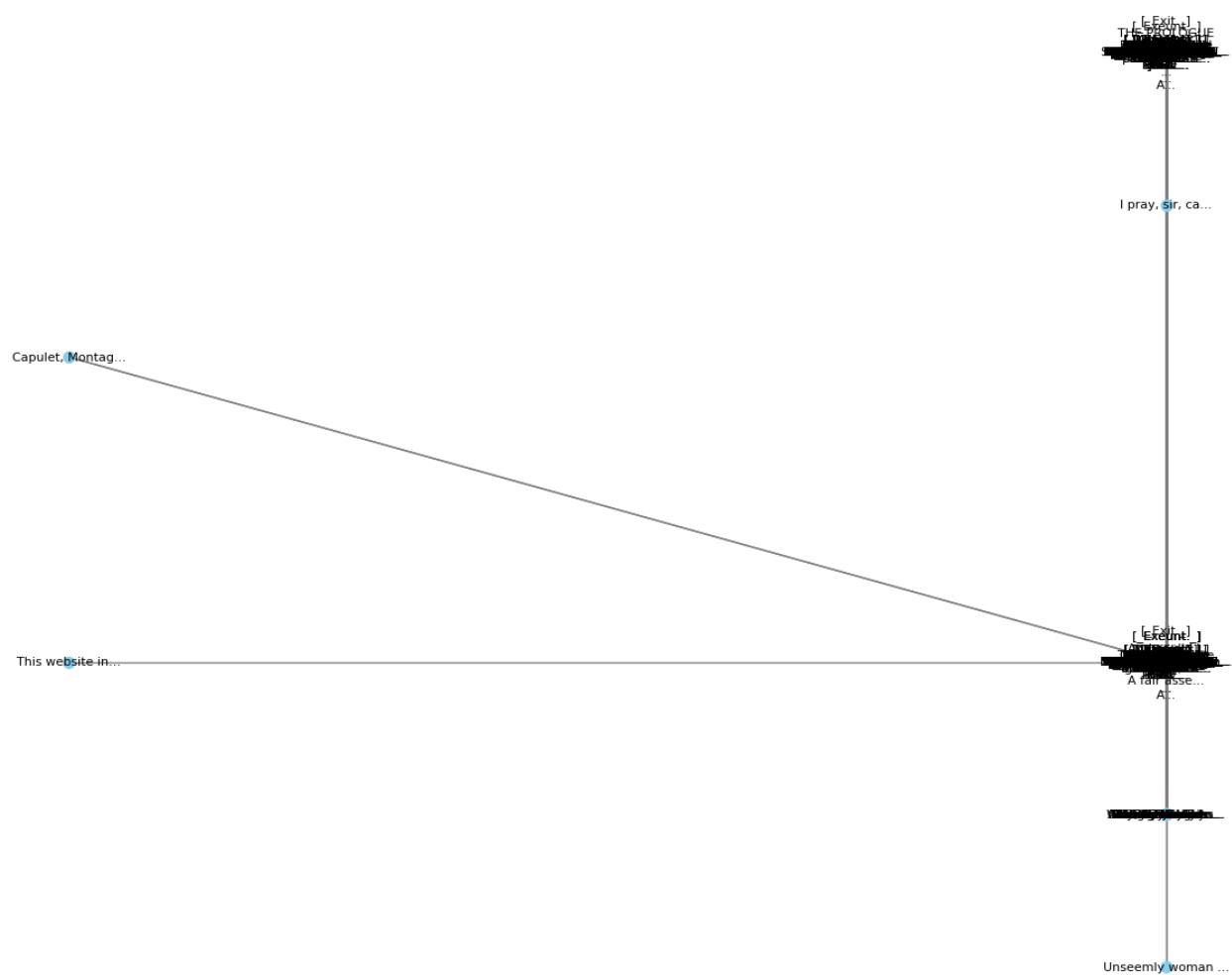
**Query search on the graphs results** Das ist es, was wir wollen, da die Suche im Graphen linear ist und auf verkettetem Wissen basiert, wobei die Knoten Daten über sich selbst enthalten.

***Ein perfektes Resultat***

## Graph Visualization from GraphMERT Model Output Embeddings



**Ein fast perfektes Resultat**



- **Extraktion:** Umwandlung von Text in Entitäten und Relationen.
- 

**Aggregation:** Semantische Aggregation zur Reduzierung von Redundanz.

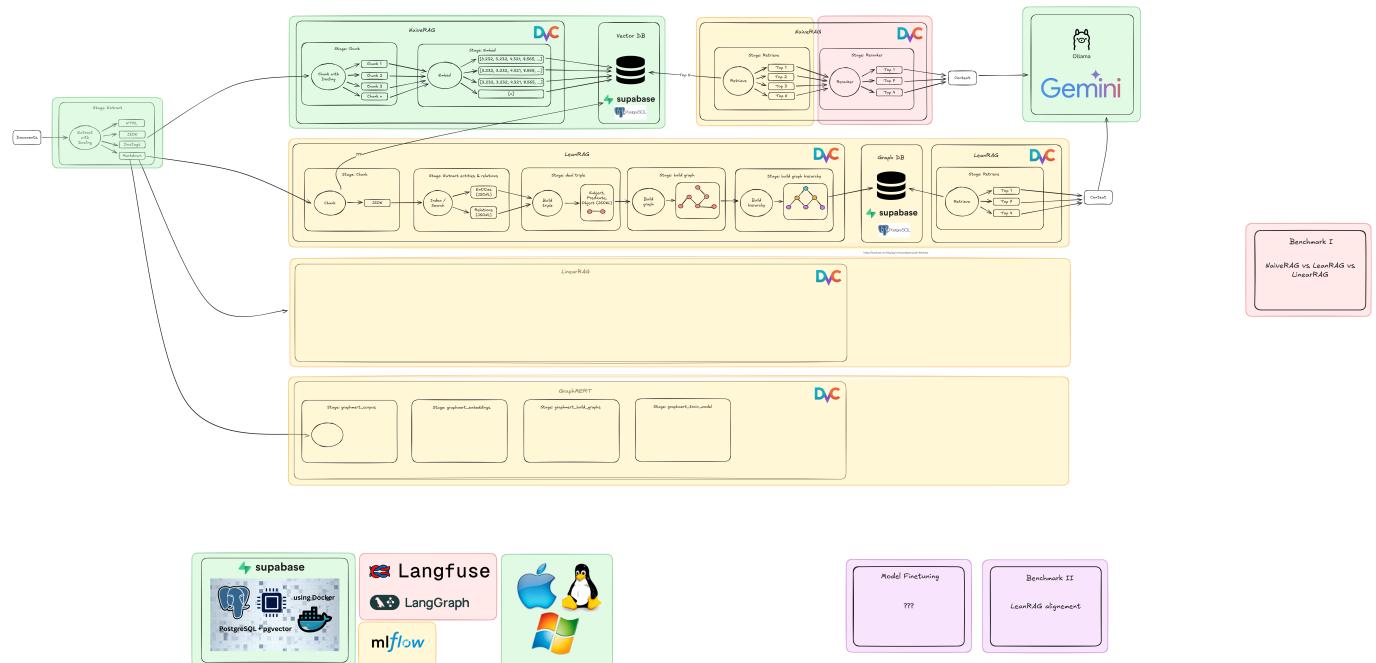
## 4.2 Fine-tuning Strategie

- Verwendung des **Unslot Frameworks** für ressourceneffizientes Training.
- Integration von Ansätzen wie **GraphRAFT** oder **GraphMERT** zur Distillation von Wissen in kleine, domänenspezifische Modelle.

## 5. Implementierung

### 5.1 Systemarchitektur

Beschreibung der Pipeline von der PDF-Eingabe bis zur Antwortgenerierung.



## 5.2 Verwendete Hardware

1 Lenovo Tower i9-14900, RAM 64.0 GB, GPU 4090 Desktop 16GB VRAM. 1 Generic Tower i9-14900, RAM 256.0 RAM, 3x RTX 6000 48.0GB VRAM => Total 144GB VRAM 1 HP EliteBook X G11 => Massenextraktion mit Docing Prozessor Intel 5U, RAM 32.0GB

1 Lenovo Notebooks Legion 9 16IRX8 Prozessor 13th Gen Intel(R) Core(TM) i9-13980HX (2.20 GHz) Installierter RAM 32.0 GB (31.7 GB verwendbar) GPU Nvidia RTX4090 Mobile mit 16GB VRAM 1 MacBook M3 Pro RAM 32.0 GB shared 1 Generic Tower RAM64.0 GB, GPU Nvidia RTX5060TI 16GB VRAM 1 Lenovo T14 RAM, OS Pop!\_OS 22.04 LTS, GPU embedded

## 6. Evaluation und Benchmarking

### 6.1 Benchmark-Design

- 

**Ansatz 1:** Generierung eines Testdatensatzes mittels Synthetic Data Generation (SDG) und Evaluierung durch ein "LLM als Judge".

- 

**Ansatz 2:** Nutzung publizierter Benchmarks wie dem Open RAG Benchmark.

### 6.2 Ergebnisse

Vergleich der Performance: Standard RAG vs. IMARA GraphRAG vs. Fine-tuned Model.

=====

### Docing Results

angetroffene Herausforderungen **Challenge:** Die Qualität der Ergebnisse liegt unter den Erwartungen.

**Massnahme 1:** Optimierung der Parameter. Die optimierte Version der Parameter ist massiv schneller und viel genauer.

```

310
311 ## 6 CONCLUSION
312
313 In this paper, we presented the DocLayNet dataset. It provides the document conversion and layout anal
314
315 From the dataset, we have derived on the one hand reference metrics for human performance on document
316
317 To date, there is still a significant gap between human and ML accuracy on the layout interpretation t
318
319 ## REFERENCES
320
321 [1] Max Göbel, Tamer Hassan, Erminda Orse, and Giorgio Orsi. Idar 2013 table competition. In 2013 I
322 [2] Chiaki Cleaver, Alessandro Sartori, and Stefan Jeschke. Doclens: A dataset for document layout analysis
323 [3] Giorgio Orsi, Tamer Hassanein, and Giorgio Orsi. Doclens: A dataset for document layout analysis
324 [4] Antonio Billone, Giorgio Orsi, and Douglas S. Hahn. Scientific literature in the International Conference
325 [5] Logan Merkewich, Hao Tong, Haoyue Xie, Xing Bao, Xu Li, Sheng, and Jian Li. Data-Sharpening for Doc
326 [6] Xu Li, Zhang Bing, and Jian Li. Data-Sharpening for Image Annotation and Document Analysis. In
327 [7] Xiao Xiang, Zhang Jian, and Jian Li. Data-Sharpening for Image Annotation and Document Analysis. In
328
329 [8] Ross B. Girshick, Andrew Donahue, Trevor Darrell, and J. J. Hindradrik. Technical Journal. Rich feature hierarc
330 [9] Ross B. Girshick. Fast R-CNN, and J. J. Hindradrik. Technical Journal. International Conference on Com
331 [10] Shaogang Gong, Ross B. Girshick, and J. J. Hindradrik. Technical Journal. Fast-r-cnn-towards-the-world
332 [11] Kaiming He, Georgios Gkioxari, Piotr Dollar, and J. J. Hindradrik. Technical Journal. IEEE International
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
599
600
601
602
603
604
605
606
607
608
609
609
610
611
612
613
614
615
616
617
618
619
619
620
621
622
623
624
625
626
627
628
629
629
630
631
632
633
634
635
636
637
638
639
639
640
641
642
643
644
645
646
647
648
649
649
650
651
652
653
654
655
656
657
658
659
659
660
661
662
663
664
665
666
667
668
669
669
670
671
672
673
674
675
676
677
678
679
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
698
699
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887

```

```

310 ## 6 CONCLUSION
311
312 #In this paper, we presented the DocLayout dataset. It provides the document conversion and layout analysis
313 #From the dataset, we have derived on the one hand reference metrics for human performance on document
314 #interpretation, and on the other hand, we have provided a baseline for state-of-the-art models.
315 #To date, there is still a significant gap between human and ML accuracy on the layout interpretation
316 #task.
317
318 ## REFERENCES
319
320 # [1] Max Gobel, Tamir Hassan, Ermelinda Oro, and Giorgio Orsi. Icdar 2013 table competition. In 2013
321 # International Conference on Document Analysis and Recognition, pages 1–6. IEEE, 2013.
322 # [2] Mingxing Hu, Djenji Ng, Junjie Yan, and Lingqiao Gao. and Stefan Fleuret. Icdar2015 competition. In 2015
323 # International Conference on Document Analysis and Recognition, pages 1–6. IEEE, 2015.
324 # [3] Antonio J. Moreno, Yves G. Leblebici, Peter Zhang, and Douglas Bartzakos. A survey on scientific literature pa-
325 #persing. In 2013 International Conference on Document Analysis and Recognition (ICDAR), pages 1–6. IEEE, 2013.
326 # [4] Logan Merkiewich, Hao Zhang, and Vivian Xing. Naval: LabelBinarizer-Shizrid, Jiang Shexin, Roy Lee, Zhen Li, Li
327 #Xu, and Ming Shou. LabelBinarizer-Shizrid: A fast and accurate document layout analysis system. In 2018 Interna-
328 #tional Conference on Document Analysis and Recognition (ICDAR), pages 1–6. IEEE, 2018.
329 # [5] Xu Chen, Ming Shou, Li Li, Zheng Yu, Lei Cui, Shuchen Huang, Furu Wei, Shuchun Li, and Ming Shou. BookBank: A
330 #Chinese book dataset. In 2018 International Conference on Document Analysis and Recognition (ICDAR), pages 1–6.
331 # [6] Riaz Ahmad, Muhammad Tariq Afzal, and M. Qadir. Information extraction from pdf sources based on
332 #hierarchical features. In 2018 International Conference on Document Analysis and Recognition (ICDAR), pages 1–6.
333 # [7] Ross B. Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accu-
334 #rate object detection and semantic segmentation. In 2014 IEEE International Conference on Computer Vision (ICCV),
335 # pages 580–587.
336 # [8] Shaqin Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object
337 #detection with region proposal networks. In 2015 IEEE International Conference on Computer Vision (ICCV),
338 # pages 919–926.
339 # [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In 2017 IEEE Interna-
340 #tional Conference on Computer Vision (ICCV), pages 743–751.
341 # [10] Glenn Goher, Alex Stokos, Aysha Chauhan, Jitka Borovcik, Nancodellis, ZhaoLi, Yonghe Feng,
342 # and Ming Shou. LayoutNet: Pre-training
343 # [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Se-
344 # [12] Mingxing Hu, Ruinan Feng, and Quoc V. Le. Efficientdet: Scalable and efficient object detect-
345 # [13] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross S. Girshick, James H.
346 # [14] Tsuxin Lu, Alexander Kirillov, Francisco Massa, Wan-You Lo, and Ross Girshick. Detectron2. In 2020
347 # [15] Nikolaos Litsakhinos, Cesare Baraclo, Mihaly Lysik, Viktor Korogodetsky, Ahmed Hassas, Andre Car-
348 # [16] Yilong Xu, Minghao Li, Lei Cui, Shuchen Huang, Furu Wei, and Ming Shou. LayoutNet: Pre-training
349 # [17] Shoubin Liu, Xuyan Wu, Shuaiqin Han, Chenjun Shi, and Qing Wang. Vlayout: Fusion of visual
350 # [18] Peng Zhang, Kai Li, Liang Qiao, Fanben Zhang, Shiliang Qi, and Qiang Wang. Vl layout: A unified
351 # [19] Peter W. Stach, Michele Dolfi, Christoph Auer, and Costas Bekas. Corpus conversion services A
352 # [20] Praveen Sharpen, and Tachi M. Khoshnoudfar. A survey on image data augmentation for deep learning

```

Die Unterschiede sind z.T. ganze Tabellen.

91     >We did not control the document selection with regard to language. The vast majority of documents come  
92     >To ensure that future benchmarks in the document-layout analysis can be easily compared, we split up D  
93     >E" (2)Se, e.g. AAFL from <https://www.semilegible.com/>  
94  
95     Table 1 shows the overall frequency of documents among the different sets. We ensure that subsets are  
96     In order to accommodate the different types of models currently in use by the community, we provide Doc  
97     Despite being cost-inducing and far less scalable than other languages, human annotation has several b  
98  
99  
100     The annotation campaign was carried out in four phases. In phase one, we identified and prepared the d  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
  
the textual content of an element, which goes beyond visual layout recognition, in particular outside  
At first sight, the task of visual document-layout interpretation appears intuitive enough to obtain p  
Obviously, this is an issue with plausible annotations. For example, examples of plausible but inconsi  
- (1) Every list-item is an individual object with class label list-item . This definition is different  
- (2) A List-item is a paragraph with hanging indentation. Single elements can be manipulated as List-  
- (3) For every Caption , there must be exactly one corresponding Picture or Table.  
- (4) Connected sub-pictures are grouped together in one Picture object.

```

91 *We did not control the document selection with regard to language. The vast majority of documents come
92 *To ensure that future benchmarks in the document-layout analysis community can be easily compared, we
93 *e.g. AAPL from https://www.annotationsreports.com/.
94 *
95 *Table 1 shows the overall frequency and distribution of the labels among the different sets. Importantly,
96 *In order to accommodate the different types of models currently in use by the community, we provide both
97 *descriptions in JSON and XML format.
98 *Despite being cost-intense and far less scalable than automation, human annotation has several benefits
99 *such as the ability to handle complex visual structures and the capacity to detect subtle visual
100 *variations in the data.
101 *# 4 ANNOTATION CAMPAIGN
102 *
103 *The annotation campaign was carried out in four phases. In phase one, we identified and prepared the
104 *relevant documents.
105 *Table 1: DocLayoutNet dataset overview. Along with the frequency of each class label, we present the relative
106 *distribution of the data across three splits: Train, Test, and Val.
107 *
108 *

| class label | Count  | % of Total | % of Total | % of Total | triple inter-annotator mAP |       |
|-------------|--------|------------|------------|------------|----------------------------|-------|
| Caption     | 22524  | 2.04       | 1.77       | 2.32       | All                        | 84-89 |
| Footnote    | 6318   | 0.59       | 0.50       | 0.58       |                            | 85-91 |
| Form        | 11937  | 2.25       | 1.80       | 2.36       |                            | 85-95 |
| List-item   | 185660 | 17.19      | 13.34      | 15.82      |                            | 87-88 |


222 *Figure 3: Corpus Conversion Service annotation user interface. The PDF page is shown in the background
223 *and the annotation tool is overlaid on top of it.
224 *
225 *
226 *distributed the annotation workload and performed continuous quality controls. Phase one and two ran
227 *from January to April 2019. Phase three ran from May to June 2019. Phase four ran from July to August 2019.
228 *
229 *Phase 1: Data selection and preparation. Our inclusion criteria for documents were described in Section
230 *3.1.
231 *Preparation work included uploading and parsing the sourced PDF documents in the Corpus Conversion
232 *System.
233 *
234 *Phase 2: Label selection and guideline. We reviewed the collected documents and identified the most common
235 *label classes.
236 *https://arxiv.org/
237 *
238 *the textual content of an element, which goes beyond visual layout recognition, in particular outside
239 *of tables.
240 *At first sight, the task of visual document-layout interpretation appears intuitive enough to obtain
241 *high-quality annotations.
242 *Obviously, this inconsistency in annotations is not desirable for datasets which are intended to be used
243 *in machine learning models.
244 *


245 *- (1) Every list-item is an individual object instance with class label List-item. This definition is
246 *widely accepted.

247 *- (2) A List-item is a paragraph with hanging indentation. Singling out elements can qualify as List-item.

248 *- (3) For every Caption, there must be exactly one corresponding Picture or Table.

249 *- (4) Connected sub-pictures are grouped together in one Picture object.

250 *
```

problematische Parameter:

```
226     params = [
227         "pipeline": "vlm",
228         "from_formats": ["docx", "pptx", "html", "image", "pdf", "asciidoc", "md", "xlsx"],
229         "to_formats": [ "md", "json", "html", "text", "doctags"], # Option "html_split_page"
230         "image_export_mode": "placeholder", # Allowed values: "placeholder", "embedded", "referenced". Optional, defaults to eml
231         "do_ocr": True,
232         "force_ocr": False,
233         "ocr_engine": "easyocr",
234         "ocr_lang": [ "en"], # en, fr, de, es
235         "pdf_backend": "diparse_v4",
236         "table_mode": "accurate",
237         "abort_on_error": False,
238
239         "do_table_structure": True, # default is True
240         "include_images": True, # default is True
241         # "do_code_enrichment": True, # default is False
242         # "do_formula_enrichment": True, # default is False
243         # "do_picture_classification": True, # default is False
244         "do_picture_description": True, # default is False
245         "picture_description_api": None, # "http://localhost:11435/v1/",
246         # "vlm_pipeline_model": "granite3.2-vision:2b",
247         # "vlm_pipeline_model_api": "http://localhost:11434/v1/chat/completions", # vlm_pipeline_model_api,
```

## erfolgreiche Parameter:

```

73 parameters = {
74     "from_formats": ["docx", "pptx", "html", "image", "pdf", "asciidoc", "md", "xlsx"],
75     "to_formats": ["md", "json", "html", "text", "doctags"], # Option "html_split_page"
76     "image_export_mode": "placeholder", # Allowed values: placeholder, embedded, referenced. Optional, defaults to embedded.
77     "do_ocr": True,
78     "force_ocr": False,
79     "ocr_engine": "easyocr",
80     "ocr_lang": ["en"],
81     "pdf_backend": "diparse_v4",
82     "table_mode": "accurate",
83     "abort_on_error": False,
84     # "do_table_structure": True, # default is True
85     # "include_images": True, # default is True
86     # "do_code_enrichment": True, # default is False
87     # "do_formula_enrichment": True, # default is False
88     # "do_picture_classification": True, # default is False
89     # "do_picture_description": True, # default is False
90     # "picture_description_api": "http://localhost:11434/v1/chat/completions",
91     # "vlm_pipeline_model": "granite3.2-vision:2b",
92     # "vlm_pipeline_model_api": vlm_pipeline_model_api,
93
94 }
95     # "target": "zip",

```

**Challenge:** Die 16GB VRAM waren nicht genug, um alle features von docling zu unterstützen. Das verursachte periodische Endless-loop's in Docling serve.

**Massnahme 1:** Der Verzicht auf die Container-Version "Docling serve" und die Verwendung direkt in Python.

**Massnahme 2:** Die Ausführung von Docling auf der CPU, um das VRAM-Limit zu umgehen

**Challenge:** Die cloudcode\_cli.exe in der VSCode-Umgebung hat durch einen etremen RAM-Verbrauch im Hintergrund die Ausführung von docling verhindert. freeze, not started, ... <https://forum.cursor.com/t/high-memory-consumption-on-cloudcode-cli/106122>

**Massnahme 1:** Ein Uninstall von cloudcode\_cli.exe war unumgänglich.

**Challenge:** Das parsen von Formeln in Docling mit CPU oder GPU ist sehr langsam. Den Verzicht auf die Extraktion der Formeln war keine Option, da eine maximale Qualität des Extrakts abgestrebt wurde, um die over-all Performance nicht zu beeinträchtigen.

Docling Log Ausschnitt:

```

[WindowsPath('C:/Users/ML4SE/Desktop/openspec_demo/configs/data/OpenRAGBench/pdfs/2411.02951v2.pdf')]
2025-12-17 19:08:35,249 - INFO - detected formats: [<InputFormat.PDF: 'pdf'>]
2025-12-17 19:08:35,259 - INFO - Going to convert document batch...
2025-12-17 19:08:35,260 - INFO - Processing document 2411.02951v2.pdf
2025-12-18 01:37:07,514 - INFO - Finished converting document 2411.02951v2.pdf in 23312.29 sec.
mpve the source file to the target directory
2025-12-18 01:37:07,940 - INFO - Processed 1 docs, of which 0 failed and 0 were partially converted.
2025-12-18 01:37:07,948 - INFO - Document conversion complete in 203589.20 seconds. it successfully completed 1 out of 287
[WindowsPath('C:/Users/ML4SE/Desktop/openspec_demo/configs/data/OpenRAGBench/pdfs/2411.03001v2.pdf')]
2025-12-18 01:37:07,968 - INFO - detected formats: [<InputFormat.PDF: 'pdf'>]
2025-12-18 01:37:07,972 - INFO - Going to convert document batch...

```

```
2025-12-18 01:37:07,973 - INFO - Processing document 2411.03001v2.pdf
2025-12-18 14:22:26,866 - INFO - Finished converting document 2411.03001v2.pdf in
45918.92 sec.
mpve the source file to the target directory
2025-12-18 14:22:27,152 - INFO - Processed 1 docs, of which 0 failed and 0 were
partially converted.
2025-12-18 14:22:27,160 - INFO - Document conversion complete in 249508.41
seconds. it successfully completed 1 out of 286
[WindowsPath('C:/Users/ML4SE/Desktop/openspec_demo/configs/data/OpenRAGBench/pdfs/
2411.03166v3.pdf')]
2025-12-18 14:22:27,193 - INFO - detected formats: [<InputFormat.PDF: 'pdf'>]
2025-12-18 14:22:27,201 - INFO - Going to convert document batch...
2025-12-18 14:22:27,202 - INFO - Processing document 2411.03166v3.pdf
2025-12-19 03:50:46,515 - INFO - Finished converting document 2411.03166v3.pdf in
48499.35 sec.
mpve the source file to the target directory
2025-12-19 03:50:47,201 - INFO - Processed 1 docs, of which 0 failed and 0 were
partially converted.
2025-12-19 03:50:47,229 - INFO - Document conversion complete in 298008.48
seconds. it successfully completed 1 out of 285
[WindowsPath('C:/Users/ML4SE/Desktop/openspec_demo/configs/data/OpenRAGBench/pdfs/
2411.03257v3.pdf')]
2025-12-19 03:50:47,249 - INFO - detected formats: [<InputFormat.PDF: 'pdf'>]
2025-12-19 03:50:47,257 - INFO - Going to convert document batch...
2025-12-19 03:50:47,259 - INFO - Processing document 2411.03257v3.pdf
2025-12-19 23:49:15,094 - INFO - Finished converting document 2411.03257v3.pdf in
71907.86 sec.
mpve the source file to the target directory
2025-12-19 23:49:17,939 - INFO - Processed 1 docs, of which 0 failed and 0 were
partially converted.
2025-12-19 23:49:18,034 - INFO - Document conversion complete in 369919.29
seconds. it successfully completed 1 out of 284
```

**Massnahme 1:** Einen zweiten Rechner 100% dafür einsetzen.

=====

leanRAG Results

Ressourcenbedarf nach den Refactoring (Schritt tripple extraction):

Prozesse				 Neuen Task ausführen
Name	Status	CPU	Arbeitsspeicher	Uhrzeit
>  Visual Studio Code (21)	0%	7%	1'663.6 MB	10:45:00
>  Microsoft Edge (17)	0%	49%	689.9 MB	10:45:00
>  LM Studio (10)	3.0%	49%	540.4 MB	10:45:00

=====

## linearRAG Results

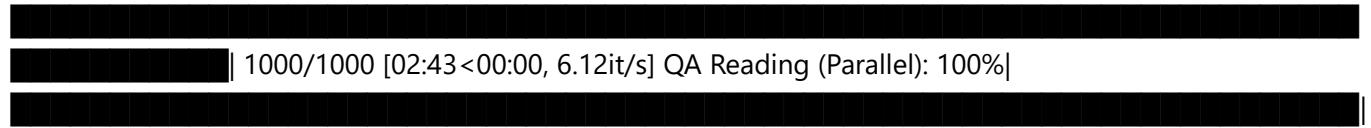
LinearRAG, Dataset: 2wikimultihop, Results with local GPT-OSS-20b Model

```
[passage] Loaded 658 records from ./import\2wikimultihop\passage_embedding.parquet [entity] Loaded  
40320 records from ./import\2wikimultihop\entity_embedding.parquet [sentence] Loaded 21206 records from  
./import\2wikimultihop\sentence_embedding.parquet
```

2025-12-09 12:16:23,189 - INFO - Evaluation Results: 2025-12-09 12:16:23,191 - INFO - LLM Accuracy: 0.7350  
(735.0/1000) 2025-12-09 12:16:23,191 - INFO - Contain Accuracy: 0.7210 (721/1000)

LinearRAG, Dataset: 2wikimultihop, Results with online gpt-4o-mini Model

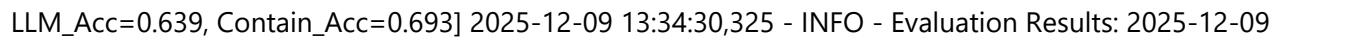
```
[passage] Loaded 658 records from ./import\2wikimultihop\passage_embedding.parquet [entity] Loaded  
40320 records from ./import\2wikimultihop\entity_embedding.parquet [sentence] Loaded 21206 records from  
./import\2wikimultihop\sentence_embedding.parquet Retrieving: 100%
```



1000/1000 [02:43<00:00, 6.12it/s] QA Reading (Parallel): 100%



1000/1000 [03:48<00:00, 4.37it/s] Evaluating samples: 100%



LLM\_Acc=0.639, Contain\_Acc=0.693] 2025-12-09 13:34:30,325 - INFO - Evaluation Results: 2025-12-09  
13:34:30,325 - INFO - LLM Accuracy: 0.6390 (639.0/1000) 2025-12-09 13:34:30,325 - INFO - Contain Accuracy:  
0.6930 (693/1000)

LinearRAG, Dataset: 2wikimultihop, Results with remote gemma3:17b Model

```
[passage] Loaded 658 records from ./import\2wikimultihop\passage_embedding.parquet [entity] Loaded  
40320 records from ./import\2wikimultihop\entity_embedding.parquet [sentence] Loaded 21206 records from  
./import\2wikimultihop\sentence_embedding.parquet Retrieving: 100%
```



1000/1000 [03:10<00:00, 5.24it/s] QA Reading (Parallel): 100%

```
[ 1000/1000 [1:22:15<00:00, 4.94s/it] Evaluating samples: 100% | 1000/1000 [03:24<00:00, 4.88sample/s, LLM_Acc=0.240, Contain_Acc=0.351] 2025-12-09 19:02:34,979 - INFO - Evaluation Results: 2025-12-09 19:02:34,980 - INFO - LLM Accuracy: 0.2400 (240.0/1000) 2025-12-09 19:02:34,981 - INFO - Contain Accuracy: 0.3510 (351/1000)
```

LinearRAG, Dataset: 2wikimultihop, Results with online gpt-4o Model

```
[passage] Loaded 658 records from ./import\2wikimultihop\passage_embedding.parquet [entity] Loaded 40320 records from ./import\2wikimultihop\entity_embedding.parquet [sentence] Loaded 21206 records from ./import\2wikimultihop\sentence_embedding.parquet Retrieving: 100%
```

```
[ 1000/1000 [03:00<00:00, 5.55it/s] QA Reading (Parallel): 100% | 1000/1000 [03:29<00:00, 4.78it/s] Evaluating samples: 100% | 1000/1000 [00:40<00:00, 24.96sample/s, LLM_Acc=0.590, Contain_Acc=0.755] 2025-12-09 19:32:14,264 - INFO - Evaluation Results: 2025-12-09 19:32:14,264 - INFO - LLM Accuracy: 0.5900 (590.0/1000) 2025-12-09 19:32:14,265 - INFO - Contain Accuracy: 0.7550 (755/1000)
```

LinearRAG, Dataset: hotpotqa, Results with local GPT-OSS-20b Model

```
[passage] Loaded 1311 records from ./import\hotpotqa\passage_embedding.parquet [entity] Loaded 66846 records from ./import\hotpotqa\entity_embedding.parquet [sentence] Loaded 38455 records from ./import\hotpotqa\sentence_embedding.parquet Retrieving: 100%
```

```
[ 1000/1000 [03:46<00:00, 4.42it/s] QA Reading (Parallel): 100% | 1000/1000 [1:51:26<00:00, 6.69s/it] Evaluating samples: 100% |
```

```
1000/1000 [24:59<00:00, 1.50s/sample, LLM_Acc=0.771, Contain_Acc=0.662] 2025-12-10 20:59:41,463 - INFO - Evaluation Results: 2025-12-10 20:59:41,463 - INFO - LLM Accuracy: 0.7710 (771.0/1000) 2025-12-10 20:59:41,463 - INFO - Contain Accuracy: 0.6620 (662/1000)
```

LinearRAG, Dataset: musique, Results with local GPT-OSS-20b Model

```
[passage] Loaded 1354 records from ./import\musique\passage_embedding.parquet [entity] Loaded 67532 records from ./import\musique\entity_embedding.parquet [sentence] Loaded 39110 records from ./import\musique\sentence_embedding.parquet Retrieving: 100%
```

```
[ 1000/1000 [03:15<00:00, 5.13it/s] QA Reading (Parallel): 100% | 1000/1000 [3:51:21<00:00, 13.88s/it] Evaluating samples: 100% | 1000/1000 [17:39<00:00,
```

```
1.06s/sample, LLM_Acc=0.642, Contain_Acc=0.317] 2025-12-11 02:00:28,341 - INFO - Evaluation Results: 2025-12-11 02:00:28,342 - INFO - LLM Accuracy: 0.6420 (642.0/1000) 2025-12-11 02:00:28,342 - INFO - Contain Accuracy: 0.3170 (317/1000)
```

LinearRAG, Dataset: medical, Results with local GPT-OSS-20b Model

[passage] Loaded 225 records from ./import\medical\passage\_embedding.parquet [entity] Loaded 9033 records from ./import\medical\entity\_embedding.parquet [sentence] Loaded 8985 records from ./import\medical\sentence\_embedding.parquet Retrieving: 100%

[redacted] | 2062/2062 [06:03<00:00, 5.67it/s] QA Reading (Parallel): 100%  
[redacted] | 2062/2062 [10:51<00:00, 3.17it/s] Evaluating samples: 100%| 2062/2062 [01:26<00:00,  
23.72sample/s, LLM\_Acc=0.694, Contain\_Acc=0.032] 2025-12-11 09:33:43,939 - INFO - Evaluation Results:  
2025-12-11 09:33:43,939 - INFO - LLM Accuracy: 0.6940 (1431.0/2062) 2025-12-11 09:33:43,939 - INFO -  
Contain Accuracy: 0.0320 (66/2062)

## 7. Diskussion der Ergebnisse

- Qualität der generierten Graphen.
- Effektivität des Fine-tunings im Vergleich zu GPT-basierten Modellen.
- Ressourcenverbrauch und Skalierbarkeit.

## 8. Risikomanagement und Lessons Learned

Reflektion über die im Antrag identifizierten Risiken:

- Datenqualität und Graph-Dichte.
- Rechenintensität des Fine-tunings.
- Teamkoordination.
- Der Vorsatz Plattformunabhängig zu sein hatte sich im Laufe des Projekts als unnötige Herausforderung herausgestellt. Konkret Microsoft Windows hatte bei der Installation spezielle Anforderungen, Inkompatibilität mit MLFlow und letztlich erzwungene Reboots, die mehrfach lang laufende Prozesse abgeschlossen haben.

## 9. Fazit und Ausblick

Zusammenfassung, ob ein 80M domänenspezifisches Modell tatsächlich grössere Modelle übertreffen konnte, und mögliche nächste Schritte.

Ausblick:

Aus den Ergebnissen konnten folgende Ansätze für die weitere Entwicklung abgeleitet werden:

- Die Qualität einer Knowledge Graphen wird hauptsächlich durch die Qualität der Entities beeinflusst. Ein Ansatz, um das Problem der sprachlichen Mehrdeutigkeit im Label der Entities ist, diese durch Attribute, abgeleitet aus dem Kontext, zu differenzieren. Ein Beispiel ist: "Der Müller hat dem Beruf eines Maurers"
  - Die Entity "Müller" ist folglich eine Maurer mit dem Familiennamen "Müller" und nicht eine Person mit dem Beruf Müller.

- Für eine produktive Lösung, sollten möglichst viele Verarbeitungsschritte im Scope eines einzelnen Dokuments (vor-)verarbeitet werden, bis und mit entity-relation Triples. Dies bringt folgende Vorteile mit sich:
  - Eine kontinuierliche Erweiterung des Graphen durch Vorverarbeitete Datensätze.
  - Parallelisierung
  - Die Möglichkeit, zu Entfernen, wenn Datensätze ungültig werden vereinfacht. Mögliche Gründe sind Fehler in den Daten oder Zeitbasierte Daten würde durch aktuellere ersetzt.
  - Mehrere Graphen können mit minimalem Offset für verschiedene Berechtigungsstufen erzeugt werden.
- Der Einfluss von Raum und Zeit muss systematisch im Graph-Modell berücksichtigt werden. z.B. Schwierigkeiten mit der Atmung werden auf Meereshöhe anders interpretiert wie auf dem Everest. Aktienkurse sind abhängig von der Zeit oder auch sich mit 100km/h zu bewegen war um 1900 rasend schnell und heute eher Durchschnitt.
- Für eine Knowledge Base mit verschiedenen Sprachen, können entities nur mit einer semantisch korrekten Übersetzung zusammengeführt werden. Um ständige Übersetzungen zwischen den Sprachen zu verhindern, könnte eine höhere Hierarchie mit einem Konzept-Graph repräsentiert werden. das heisst einzelne Fakten werden als Knowledge Graph dargestellt und darüber auf Konzepte abgebildet.
- Das Clustering identischer Relationen zu einem Hypergraph ist ein weiterer Ansatz, Teilgraphen zusammen zu führen, ohne sich die Möglichkeit zu verbauen Teile wieder zu entfernen. Ebenso können wahrscheinliche Relationen abgeleitet werden. (vom Hypergraph zurück zum Knowledge Graph)

## 10. Referenzen

- [1] Docling: An Efficient Open-Source Toolkit. <https://arxiv.org/abs/2501.17887> Docling: An Efficient Open-Source Toolkit for AI-driven Document Conversion <https://www.docling.ai/> <https://docling-project.github.io/docling/>
- [2] LeanRAG: Knowledge-Graph-Based Generation. <https://arxiv.org/abs/2508.10391> Knowledge-Graph-Based Generation with Semantic Aggregation and Hierarchical Retrieval <https://github.com/KnowledgeXLab/LeanRAG>
- [3] LinearRAG: A relation-free graph construction method for efficient GraphRAG. <https://arxiv.org/abs/2510.10114> LinearRAG: Linear Graph Retrieval Augmented Generation on Large-scale Corpora <https://github.com/DEEP-PolyU/LinearRAG>
- [4] GraphMERT: Efficient Distillation of Reliable KGs. <https://arxiv.org/abs/2510.09580> GraphMERT: Efficient and Scalable Distillation of Reliable Knowledge Graphs from Unstructured Data <https://github.com/creativeautomaton/graphMERT-python>
- [5] Open RAG Bench Dataset <https://github.com/vectara/open-rag-bench> Open RAG Benchmark (1000 PDFs, 3000 Queries): A Multimodal PDF Dataset for Comprehensive RAG Evaluation
- ... (Weitere Quellen gemäss Antrag).

## 11. Glossar

A - C

- **AI-Native GraphRAG:** Ein weiterentwickeltes Paradigma von GraphRAG, das den gesamten Workflow von unstrukturierten Daten bis zur Antwortgenerierung automatisiert und dabei die Komplexität von Graphentheorie und Datenbankmanagement abstrahiert.
- **Chunking:** Der Prozess des Zerlegens von Texten in kleinere Abschnitte (Chunks). Im Bericht wird dies als kritischer Faktor für *naives RAG* identifiziert, da suboptimale Chunk-Größen (zu gross oder zu klein) zu Kontextverlust oder Rauschen führen können.
- **CommonKG:** Eine im Kontext von *LeanRAG* erwähnte Methode zur Erstellung von Wissensgraphen, bei der Entitäten und Relationen (Triples) aus Text-Chunks extrahiert und dedupliziert werden.

## D - G

- **Docling:** Ein Open-Source-Toolkit zur Dokumentenkonvertierung. Im Projekt wurde es genutzt, um komplexe PDFs in maschinenlesbare Formate (Markdown/JSON) zu wandeln. Es traten Herausforderungen bezüglich VRAM-Verbrauch und Performance auf.
- **Embeddings:** Vektorrepräsentationen von Texten (Sätze, Entitäten, Passagen). Sie dienen als Basis für die Ähnlichkeitssuche und das Clustering in den Graphen.
- **FactScore:** Eine Metrik zur Bewertung der faktischen Korrektheit eines Wissensgraphen oder einer generierten Antwort. Im Bericht erzielt *GraphMERT* hierbei deutlich höhere Werte als reine LLMs.
- **Fine-tuning:** Das nachtrainieren eines LLMs (z. B. Qwen) auf spezifischen, graphenbasierten Daten, um die Antwortqualität und Domänenexpertise zu erhöhen.
- **GraphMERT:** Ein kompaktes, rein grafisches Encoder-Modell (Neurosymbolische KI), das effizient zuverlässige und ontologiekonsistente Wissensgraphen aus unstrukturierten Texten generiert.
- **GraphRAG (Graph Retrieval-Augmented Generation):** Eine Erweiterung von RAG, die statt flacher Textlisten strukturierte Wissensgraphen nutzt. Dies ermöglicht das Erkennen komplexer Beziehungen und *Multi-Hop-Reasoning*.

## H - L

- **Hypergraph:** Eine im Ausblick erwähnte Graphenstruktur, bei der eine Kante (Edge) mehr als zwei Knoten verbinden kann. Dies wird als Ansatz vorgeschlagen, um identische Relationen zu clustern.
- **IMARA:** Der Name des Projekts. Es steht für die Entwicklung einer domänenspezifischen GraphRAG-Pipeline mit Model Fine-tuning.
- **Knowledge Graph (Wissensgraph):** Eine strukturierte Darstellung von Wissen in Form von Knoten (Entitäten) und Kanten (Beziehungen), die ein aktives, abfragefähiges Modell der Welt darstellt.
- **LeanRAG:** Ein GraphRAG-Ansatz, der auf semantische Aggregation und hierarchisches Retrieval setzt, um Redundanzen zu minimieren (ca. 46 % weniger Redundanz im Vergleich zu flachen Baselines).
- **LinearRAG:** Eine effiziente GraphRAG-Methode, die "relation-free" arbeitet. Sie nutzt leichtgewichtige Entity Recognition und semantische Verlinkung für schnelle Verarbeitung mit linearer Komplexität.
- **LLM (Large Language Model):** Große Sprachmodelle, die als generative Komponente im RAG-Prozess dienen (z. B. GPT-4o, Qwen, Gemma).

## M - O

- **Multi-Hop-Reasoning:** Die Fähigkeit, Informationen über mehrere Verbindungsschritte hinweg zu verknüpfen (z. B. A ist verbunden mit B, B ist verbunden mit C → Schlussfolgerung von A auf C). Eine Schwäche von *naivem RAG*, aber eine Stärke von *GraphRAG*.
- **Naives RAG:** Bezeichnet im Bericht konventionelle, vektorbasierte RAG-Architekturen, die Wissen als unzusammenhängende Fakten (Chunks) behandeln und oft an kontextueller Fragmentierung leiden.

- **Neurosymbolische KI:** Kombination aus neuronalen Netzwerken (Generalisierung, Lernen) und symbolischer KI (Abstraktion, Logik, Graphen), wie sie im *GraphMERT*-Ansatz verwendet wird.
- **OpenRAGBench:** Ein Referenzdatensatz (Benchmark), der im Projekt genutzt wurde, um die Messbarkeit und Vergleichbarkeit der Ergebnisse sicherzustellen.

## S - V

- **Semantic Aggregation:** Ein Feature von *LeanRAG*, bei dem Entitäten in semantisch kohärente Zusammenfassungen (Cluster) gruppiert werden, um die Navigation im Graphen zu verbessern.
- **Synthetic Data Generation (SDG):** Ein Ansatz zur Generierung von künstlichen Testdaten, um die Leistung des Systems zu evaluieren (z. B. mittels "LLM als Judge").
- **Triple:** Die grundlegende Dateneinheit eines Wissensgraphen, bestehend aus Subjekt, Prädikat (Relation) und Objekt (z. B. "Müller" -> "hat Beruf" -> "Maurer").
- **Unsloth:** Ein Framework, das im Projekt für das ressourceneffiziente *Fine-tuning* der Modelle verwendet wurde.
- **ValidityScore:** Eine Metrik zur Bewertung der Gültigkeit von Relationen (Ontologie-Konsistenz) innerhalb eines Wissensgraphen.
- **Vektorsimilaritätssuche:** Das Suchverfahren klassischer RAG-Systeme, das Textabschnitte basierend auf mathematischer Ähnlichkeit (Vektornähe) findet, aber explizite Beziehungen oft ignoriert.

---

## Tipps für die Ausarbeitung

- **Code-Beispiele:** Fügt kurze Snippets eurer Automatisierungslösung oder der Unsloth-Konfiguration in Kapitel 5 ein.
- **Metriken:** In Kapitel 6 solltet ihr Tabellen mit Latenzzeiten und Genauigkeitswerten (Accuracy/F1) eurer Benchmarks zeigen.