

Novelty Assessment Report

Paper: AtlasKV: Augmenting LLMs with Billion-Scale Knowledge Graphs in 20GB VRAM

PDF URL: <https://openreview.net/pdf?id=6i1jVAYbHs>

Venue: ICLR 2026 Conference Submission

Year: 2026

Report Generated: 2026-01-07

Abstract

Retrieval-augmented generation (RAG) has shown some success in augmenting large language models (LLMs) with external knowledge. However, as a non-parametric knowledge integration paradigm for LLMs, RAG methods heavily rely on external retrieval modules and the retrieved textual context prior. Especially for very large scale knowledge augmentation, they would introduce substantial inference latency due to expensive searches and much longer relevant context. In this paper, we propose a parametric knowledge integration method, called $\text{\textbf{AtlasKV}}$, a scalable, effective, and general way to augment LLMs with billion-scale knowledge graphs (KGs) (e.g. 1B triples) using very little GPU memory cost (e.g. less than 20GB VRAM). In AtlasKV, we introduce KG2KV and HiKVP to integrate KG triples into LLMs at scale with sub-linear time and memory complexity. It maintains strong knowledge grounding and generalization performance using the LLMs' inherent attention mechanism, and requires no external retrievers, long context priors, or retraining when adapting to new knowledge.

Disclaimer

This report is **AI-GENERATED** using Large Language Models and WisPaper (a scholar search engine). It analyzes academic papers' tasks and contributions against retrieved prior work. While this system identifies **POTENTIAL** overlaps and novel directions, **ITS COVERAGE IS NOT EXHAUSTIVE AND JUDGMENTS ARE APPROXIMATE**. These results are intended to assist human reviewers and **SHOULD NOT** be relied upon as a definitive verdict on novelty.

Note that some papers exist in multiple, slightly different versions (e.g., with different titles or URLs). The system may retrieve several versions of the same underlying work. The current automated pipeline does not reliably align or distinguish these cases, so human reviewers will need to disambiguate them manually.

If you have any questions, please contact: mingzhang23@m.fudan.edu.cn

Core Task Landscape

This paper addresses: **Augmenting Large Language Models with Billion-Scale Knowledge Graphs**

A total of **6 papers** were analyzed and organized into a taxonomy with **6 categories**.

Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Parametric Knowledge Integration into LLMs**
- **Text-Mediated Knowledge Integration**
- **Hybrid Knowledge Integration Architectures**
- **Knowledge Graph Construction and Schema Induction**
- **Enterprise Knowledge Graph Applications**

Complete Taxonomy Tree

- Augmenting Large Language Models with Billion-Scale Knowledge Graphs Survey Taxonomy
- Parametric Knowledge Integration into LLMs
 - Direct KG-to-Parameter Encoding ★ (1 papers)
 - [0] AtlasKV: Augmenting LLMs with Billion-Scale Knowledge Graphs in 20GB VRAM (Anon et al., 2026) [View paper](#)
 - Pre-trained Knowledge Graph Embeddings for LLMs (1 papers)
 - [6] PKGM: A Pre-trained Knowledge Graph Model for E-commerce Application (Zhang Wen, 2022) [View paper](#)
- Text-Mediated Knowledge Integration
 - KG Verbalization for Corpus Augmentation (1 papers)
 - [2] Knowledge Graph Based Synthetic Corpus Generation for Knowledge-Enhanced Language Model Pre-training (Agarwal, 2021) [View paper](#)
- Hybrid Knowledge Integration Architectures
 - Multi-Encoder Fusion for Structured and Textual Knowledge (2 papers)
 - [3] StATIK+: Structure and Text for Inductive Knowledge Graph Modeling and Paths towards Enterprise Implementations. (ES Markowitz, 2023) [View paper](#)
 - [4] Integrating Large Language Models and Knowledge Graphs to Improve Factuality and Reasoning (Unknown, 2025) [View paper](#)
- Knowledge Graph Construction and Schema Induction
 - Autonomous Schema-Free KG Construction (1 papers)
 - [1] AutoSchemaKG: Autonomous Knowledge Graph Construction through Dynamic Schema Induction from Web-Scale Corpora (Bai Jiaxin, 2025) [View paper](#)
- Enterprise Knowledge Graph Applications
 - Entity Resolution and Disambiguation (1 papers)
 - [5] Hybrid Neural Ensemble Architectures for Investment Entity Disambiguation in Enterprise Graphs (C Bishop, 2025) [View paper](#)

Narrative

Core task: Augmenting large language models with billion-scale knowledge graphs. The field organizes around several complementary strategies for integrating structured knowledge into LLMs. Parametric Knowledge Integration into LLMs explores methods that encode graph facts directly into model parameters, often through continued pretraining or specialized embedding techniques. Text-Mediated Knowledge Integration converts structured triples into natural language passages that LLMs can consume more naturally, exemplified by approaches like Synthetic Corpus Generation[2]. Hybrid Knowledge Integration Architectures combine parametric and retrieval-based mechanisms, as seen in systems like Hybrid Neural Ensemble[5], to balance memorization with dynamic access. Knowledge Graph Construction and Schema Induction focuses on building and refining the graphs themselves, with works such as AutoSchemaKG[1]

addressing schema-level challenges at scale. Enterprise Knowledge Graph Applications targets domain-specific deployments, including e-commerce scenarios explored by PKGM Ecommerce[6], where billion-scale graphs must support real-world query loads.

A central tension across these branches involves the trade-off between embedding knowledge within model weights versus retrieving it on demand. Parametric approaches promise faster inference and tighter integration but face scalability limits and update challenges, while text-mediated and hybrid methods offer flexibility at the cost of additional retrieval overhead. Within the parametric branch, AtlasKV[0] pursues direct KG-to-parameter encoding, aiming to compress vast relational structures into the model itself. This contrasts with hybrid strategies like Hybrid Neural Ensemble[5], which maintain separate retrieval modules, and with text-mediated techniques such as Synthetic Corpus Generation[2], which rely on verbalized triples. Ongoing questions include how to maintain factual consistency, as highlighted by LLMs Knowledge Graphs Factuality[4], and how to efficiently update billion-scale embeddings when graphs evolve. AtlasKV[0] sits squarely in the parametric camp, emphasizing direct encoding over retrieval augmentation.

Related Works in Same Category

No sibling papers were found in the same taxonomy leaf. A taxonomy-subtopic-level comparison will be produced instead.

Sibling Subtopics

- **Pre-trained Knowledge Graph Embeddings for LLMs** (leaves: 1, papers: 1)
- Scope: Approaches that learn unified knowledge graph representations during pre-training for direct integration into downstream language model tasks.
- Exclude: Excludes runtime retrieval-based methods and schema-free construction; see Non-Parametric Integration and Autonomous KG Construction.

Contributions Analysis

Overall novelty summary. The paper proposes AtlasKV, a parametric method for integrating billion-scale knowledge graphs directly into LLM parameters via key-value mappings. It occupies the 'Direct KG-to-Parameter Encoding' leaf within the 'Parametric Knowledge Integration into LLMs' branch. Notably, this leaf contains only the original paper itself—no sibling papers were identified in the taxonomy. This suggests the specific approach of converting KG triples to attention-compatible key-value structures at billion-scale represents a relatively sparse research direction within the broader parametric integration landscape.

The taxonomy reveals neighboring approaches in adjacent branches. 'Pre-trained Knowledge Graph Embeddings for LLMs' explores unified representations learned during pretraining, while 'KG Verbalization for Corpus Augmentation' converts triples to natural language text. 'Multi-Encoder Fusion' architectures combine graph neural networks with language models using separate encoders. AtlasKV diverges by avoiding external retrievers, text conversion, or separate graph encoders, instead leveraging the LLM's native attention mechanism. The taxonomy's scope notes explicitly exclude retrieval-based and text-mediated methods from this parametric branch, positioning AtlasKV as pursuing tighter integration than hybrid alternatives.

Among twenty-four candidates examined across three contributions, none were identified as clearly refuting the work. The AtlasKV framework examined ten candidates with zero refutable matches; KG2KV pipeline examined five with none refutable; HiKVP algorithm examined nine with none refutable. This suggests that within the limited search scope, the specific combination of billion-scale direct encoding, sub-linear complexity guarantees, and hierarchical pruning mechanisms appears relatively unexplored. However, the search examined only top-K semantic matches plus citations, not an exhaustive survey of parametric knowledge integration literature.

Given the sparse taxonomy leaf and absence of refuting candidates among twenty-four examined papers, the work appears to occupy a distinct position within parametric integration approaches. The limited search scope means potentially relevant work in adjacent areas—such as efficient attention mechanisms or knowledge distillation—may not have been captured. The analysis reflects what was found in targeted semantic search, not a comprehensive field review.

This paper presents **3 main contributions**, each analyzed against relevant prior work:

Contribution 1: AtlasKV framework for billion-scale KG augmentation

Description: AtlasKV is a parametric knowledge integration framework that enables LLMs to incorporate billion-scale knowledge graphs with minimal GPU memory requirements. It maintains strong knowledge grounding and generalization without requiring external retrievers, long context priors, or retraining when adapting to new knowledge.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Structured Knowledge Integration and Memory Modeling in Large Language Systems

URL: [View paper](#)

Brief Assessment

Structured Knowledge Integration[16] focuses on integrating memory networks and perception graphs for multi-hop reasoning tasks, not on billion-scale KG augmentation with minimal GPU memory. The technical approaches and objectives differ fundamentally from AtlasKV's parametric knowledge integration paradigm.

2. Paths-over-graph: Knowledge graph empowered large language model reasoning

URL: [View paper](#)

Brief Assessment

Paths-over-graph[8] focuses on multi-hop reasoning path exploration from KGs to enhance LLM reasoning accuracy, not on parametric knowledge integration with minimal GPU memory requirements. The technical approaches are fundamentally different: Paths-over-graph[8] uses path pruning and exploration for reasoning tasks, while the original work addresses memory-efficient KG integration into LLM parameters.

3. Unlock the power of frozen llms in knowledge graph completion

URL: [View paper](#)

Brief Assessment

Frozen LLMs Completion[12] focuses on knowledge graph completion (predicting missing triples) using frozen LLMs' intermediate representations, not on augmenting LLMs with billion-scale KGs for generation tasks. The technical approaches differ fundamentally: [12] extracts hidden states for classification, while AtlasKV injects KG triples into attention layers.

4. A few-shot learning method based on knowledge graph in large language models

URL: [View paper](#)

Brief Assessment

Few Shot KG[9] focuses on few-shot learning methods using knowledge graphs in LLMs, not on billion-scale parametric KG integration with minimal GPU memory requirements like AtlasKV.

5. Zep: a temporal knowledge graph architecture for agent memory

URL: [View paper](#)

Brief Assessment

Zep Temporal[7] focuses on temporal knowledge graphs for agent memory in conversational contexts, not on parametric integration of billion-scale KGs into LLMs with minimal GPU memory requirements. The systems address fundamentally different problems and use cases.

6. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph

URL: [View paper](#)

Brief Assessment

Think on Graph[10] focuses on interactive LLM-KG reasoning via beam search exploration, not parametric knowledge integration. It does not address billion-scale KG storage in GPU memory or the parametric embedding approach central to AtlasKV.

7. MemVerse: Multimodal Memory for Lifelong Learning Agents

URL: [View paper](#)

Brief Assessment

MemVerse Multimodal[14] focuses on multimodal memory for lifelong learning agents using hierarchical knowledge graphs and parametric memory distillation. It does not address billion-scale KG augmentation with minimal GPU memory requirements, which is the core technical contribution of the original paper.

8. Network for knowledge Organization (NEKO): An AI knowledge mining workflow for synthetic biology research

URL: [View paper](#)

Brief Assessment

NEKO Workflow[15] focuses on knowledge graph construction and organization for synthetic biology research using RAG approaches, not on parametric integration of billion-scale KGs into LLM attention mechanisms with minimal GPU memory as AtlasKV does.

9. Memory Matters: The Need to Improve Long-Term Memory in LLM-Agents

URL: [View paper](#)

Brief Assessment

Memory Matters Agents[11] focuses on long-term memory management in LLM agents using vector databases for episodic/semantic memory, not on parametric knowledge integration frameworks for billion-scale knowledge graphs with minimal GPU memory requirements.

10. Memory-augmented query reconstruction for llm-based knowledge graph reasoning

URL: [View paper](#)

Brief Assessment

Memory Augmented Query[13] focuses on LLM-based knowledge graph question answering through query reconstruction and memory-augmented reasoning, not on parametric knowledge integration of billion-scale KGs into LLM attention mechanisms with minimal GPU memory.

Contribution 2: KG2KV pipeline for converting KG triples to Q-K-V data

Description: KG2KV is a pipeline that naturally transforms knowledge graph triples into high-quality query-key-value data by leveraging the structural similarity between KG triples and self-attention Q-K-V vectors. This design enhances generalization performance by ensuring diverse enquiry attributes from massive KG relations.

This contribution was assessed against **5 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. A Knowledge Augmented Framework for Multimodal News{Object-Entity Relation Extraction

URL: [View paper](#)

Brief Assessment

Multimodal News Extraction[29] focuses on multimodal relation extraction from news headlines and images using knowledge graphs for semantic enhancement, not on converting KG triples into query-key-value representations for self-attention mechanisms in LLMs.

2. Named entity recognition for Chinese marine text with knowledge-based self-attention

URL: [View paper](#)

Brief Assessment

Marine Text NER[28] focuses on named entity recognition for Chinese marine text using knowledge-based self-attention, not on converting knowledge graph triples into query-key-value representations for general LLM augmentation.

3. Knowledge graph embeddings based on 2d convolution and self-attention mechanisms for link prediction.

URL: [View paper](#)

Brief Assessment

KG Embeddings Convolution[27] focuses on converting triples into Q-K-V for 2D convolution-based embeddings, not for augmenting LLMs with parametric knowledge integration as in the original paper's KG2KV pipeline.

4. A personalized paper recommendation method based on knowledge graph and transformer encoder with a self-attention mechanism

URL: [View paper](#)

Brief Assessment

Personalized Paper Recommendation[26] uses knowledge graphs with query-key-value vectors for paper recommendation, but does not describe a pipeline for converting KG triples into Q-K-V data. The candidate focuses on recommendation systems rather than general KG-to-attention integration methods.

5. Query-Guided Graph Neural Networks for Knowledge Graph Reasoning

URL: [View paper](#)

Brief Assessment

Query Guided GNN[30] focuses on query-key-value formulations for graph neural network message passing in knowledge graph reasoning tasks, not on converting KG triples into training data for LLM attention mechanisms.

Contribution 3: HiKVP algorithm for hierarchical key-value pruning

Description: HiKVP is a hierarchical pruning algorithm that organizes knowledge keys into a three-layer structure and progressively selects relevant key-value pairs. It achieves sub-linear time and memory complexity, enabling scalable integration of billion-scale KGs while preserving high accuracy.

This contribution was assessed against **9 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Learning Logical Rules from Knowledge Graphs

URL: [View paper](#)

Brief Assessment

Learning Logical Rules[25] focuses on learning logical rules from knowledge graphs for reasoning tasks, not on hierarchical pruning algorithms for key-value pairs in LLM attention mechanisms or knowledge integration systems.

2. Lhrs-bot: Empowering remote sensing with vgi-enhanced large multimodal language model

URL: [View paper](#)

Brief Assessment

LHRS Bot[17] focuses on remote sensing image understanding using vision-language alignment with VGI data, not on hierarchical pruning algorithms for knowledge graph integration or key-value pair management in LLMs.

3. Scalable billion-point approximate nearest neighbor search using

URL: [View paper](#)

Brief Assessment

Billion Point Search[19] focuses on hierarchical indexing for approximate nearest neighbor search in vector databases, not on hierarchical pruning of key-value pairs for knowledge graph integration into LLMs. The technical domains and applications are fundamentally different.

4. TALE: Token-Adaptive Low-Rank KVCache Approximation with Reconstruction Elimination

URL: [View paper](#)

Brief Assessment

TALE KVCache[24] focuses on token-adaptive low-rank approximation for KVCache compression in LLM inference, not hierarchical pruning of knowledge graph key-value pairs for knowledge integration.

5. On-Device Large Language Models: A Survey of Model Compression and System Optimization

URL: [View paper](#)

Brief Assessment

On Device LLMs[18] is a survey on model compression and system optimization for on-device deployment. It does not propose hierarchical pruning algorithms for key-value pairs in knowledge integration; instead, it reviews existing compression techniques like quantization, pruning, and low-rank methods for general LLM deployment.

6. Hierarchical Memory Networks for Multi-Hop Reasoning in Large-Scale Knowledge Bases

URL: [View paper](#)

Brief Assessment

Hierarchical Memory Networks[21] focuses on hierarchical memory organization for multi-hop reasoning in knowledge graphs, not on hierarchical pruning of key-value pairs for scalable KG integration into LLMs as in the original paper.

7. ArceKV: Towards Workload-driven LSM-compactions for Key-Value Store Under Dynamic Workloads

URL: [View paper](#)

Brief Assessment

ArceKV[20] focuses on LSM-tree compaction strategies for key-value stores under dynamic workloads, not on hierarchical pruning algorithms for knowledge graph integration into LLMs. The technical domains are entirely different.

8. VAFTrack: asynchronous feature fusion via visual receptive weighted key-value perceptual for visual tracking

URL: [View paper](#)

Brief Assessment

VAFTrack[22] focuses on visual tracking with feature fusion and layer pruning in transformers for computer vision tasks, not on hierarchical key-value pruning for knowledge graph integration in LLMs.

9. Leveraging KV Similarity for Online Structured Pruning in LLMs

URL: [View paper](#)

Brief Assessment

KV Similarity Pruning[23] focuses on online token-level pruning during inference using key-value similarity to skip redundant attention computations. This differs from HiKVP's hierarchical clustering approach for organizing and selecting knowledge key-value pairs from billion-scale knowledge graphs for parametric knowledge integration.

Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

References

- [0] AtlasKV: Augmenting LLMs with Billion-Scale Knowledge Graphs in 20GB VRAM [View paper](#)
- [1] AutoSchemaKG: Autonomous Knowledge Graph Construction through Dynamic Schema Induction from Web-Scale Corpora [View paper](#)
- [2] Knowledge Graph Based Synthetic Corpus Generation for Knowledge-Enhanced Language Model Pre-training [View paper](#)
- [3] STATIK+: Structure and Text for Inductive Knowledge Graph Modeling and Paths towards Enterprise Implementations. [View paper](#)

- [4] Integrating Large Language Models and Knowledge Graphs to Improve Factuality and Reasoning [View paper](#)
- [5] Hybrid Neural Ensemble Architectures for Investment Entity Disambiguation in Enterprise Graphs [View paper](#)
- [6] PKGM: A Pre-trained Knowledge Graph Model for E-commerce Application [View paper](#)
- [7] Zep: a temporal knowledge graph architecture for agent memory [View paper](#)
- [8] Paths-over-graph: Knowledge graph empowered large language model reasoning [View paper](#)
- [9] A few-shot learning method based on knowledge graph in large language models [View paper](#)
- [10] Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph [View paper](#)
- [11] Memory Matters: The Need to Improve Long-Term Memory in LLM-Agents [View paper](#)
- [12] Unlock the power of frozen llms in knowledge graph completion [View paper](#)
- [13] Memory-augmented query reconstruction for llm-based knowledge graph reasoning [View paper](#)
- [14] MemVerse: Multimodal Memory for Lifelong Learning Agents [View paper](#)
- [15] Network for knowledge Organization (NEKO): An AI knowledge mining workflow for synthetic biology research [View paper](#)
- [16] Structured Knowledge Integration and Memory Modeling in Large Language Systems [View paper](#)
- [17] Lhrs-bot: Empowering remote sensing with vgi-enhanced large multimodal language model [View paper](#)
- [18] On-Device Large Language Models: A Survey of Model Compression and System Optimization [View paper](#)
- [19] Scalable billion-point approximate nearest neighbor search using {SmartSSDs} [View paper](#)
- [20] ArceKV: Towards Workload-driven LSM-compactions for Key-Value Store Under Dynamic Workloads [View paper](#)
- [21] Hierarchical Memory Networks for Multi-Hop Reasoning in Large-Scale Knowledge Bases [View paper](#)
- [22] VAFTrack: asynchronous feature fusion via visual receptive weighted key-value perceptual for visual tracking [View paper](#)
- [23] Leveraging KV Similarity for Online Structured Pruning in LLMs [View paper](#)
- [24] TALE: Token-Adaptive Low-Rank KVCache Approximation with Reconstruction Elimination [View paper](#)
- [25] Learning Logical Rules from Knowledge Graphs [View paper](#)
- [26] A personalized paper recommendation method based on knowledge graph and transformer encoder with a self-attention mechanism [View paper](#)
- [27] Knowledge graph embeddings based on 2d convolution and self-attention mechanisms for link prediction. [View paper](#)
- [28] Named entity recognition for Chinese marine text with knowledge-based self-attention [View paper](#)
- [29] A Knowledge Augmented Framework for Multimodal News{Object-Entity Relation Extraction [View paper](#)
- [30] Query-Guided Graph Neural Networks for Knowledge Graph Reasoning [View paper](#)