# A New Frontier in Language Model Optimization: A Detailed Research Proposal for Knowledge-Guided Dynamic Sparsity

## Executive Summary

The exponential growth in the scale of Large Language Models (LLMs) has undeniably advanced artificial intelligence capabilities across a myriad of tasks, yet it simultaneously introduces formidable challenges related to computational demands and memory footprints. These challenges impede widespread real-world deployment and raise significant concerns regarding environmental sustainability. Weight pruning has emerged as a particularly promising optimization technique, offering a pathway to substantial model compression without compromising performance.

This proposal undertakes a critical evaluation of "Knowledge-Guided Dynamic Sparsity" (KG-DS), a novel, state-of-the-art methodology for LLM weight pruning. KG-DS is conceptualized to dynamically adapt LLM sparsity based on input, building upon a protected "knowledge core." The methodology aims for significant efficiency gains, encompassing reduced memory usage and accelerated inference speed, while aspiring to maintain or even surpass the quality of the original, dense model.

The analysis presented herein, grounded in a thorough review of contemporary research, highlights KG-DS's alignment with cutting-edge advancements such as gradient-based saliency estimation and hardware-friendly structured pruning. Its innovative embrace of dynamic sparsity paradigms represents a significant conceptual leap. However, a detailed feasibility assessment reveals several critical challenges. These include the potential computational overhead associated with real-time dynamic mask application, the imperative for robust and accurate instruction classification, and the practical complexities of achieving hardware-aware dynamic sparsity.

This document outlines a comprehensive research agenda designed to refine and validate KG-DS. The proposed research focuses on optimizing its core mechanisms, rigorously addressing the identified technical and practical challenges through advanced algorithmic and architectural designs, and establishing a robust, multi-faceted evaluation framework. The overarching objective is to pave the way for

the development and deployment of truly efficient, adaptable, and high-quality LLMs.

# 1. Introduction: The Imperative of LLM Optimization

The advent of Large Language Models, characterized by their immense parameter counts, has ushered in a new era of artificial intelligence, enabling unprecedented performance across diverse applications, from complex reasoning to sophisticated content generation. Models ranging from billions to trillions of parameters have demonstrated remarkable emergent abilities, pushing the boundaries of what is computationally feasible.[1] However, this impressive capability is inextricably linked to substantial resource consumption. The sheer scale of these models translates into astronomical computational requirements for both their initial training and subsequent inference, alongside immense memory footprints.[3] These resource demands present significant barriers to their widespread accessibility, practical real-world deployment, and long-term environmental sustainability. The challenges manifest as increased latency in real-time applications and high energy consumption, underscoring the urgent need for effective optimization strategies.[4]

Historically, the primary concern in optimizing large neural networks revolved around reducing the raw computational operations (FLOPs). However, a deeper understanding of LLM inference reveals a more nuanced bottleneck. During inference, the dominant constraint is often not the raw computation but rather the memory bandwidth required to load model weights to the processing units.[3] This fundamental shift in the performance bottleneck implies that optimization strategies must extend beyond merely minimizing FLOPs. They must critically consider how memory access patterns are affected and strive for memory efficiency. This understanding underscores why techniques that reduce the active memory footprint, such as weight pruning, are particularly relevant. Structured sparsity, which involves removing groups of parameters rather than individual weights, is generally more amenable to existing hardware architectures and can lead to tangible latency reductions on general-purpose devices, unlike unstructured sparsity which often necessitates specialized hardware for practical speed-ups.[1]

To address these multifaceted challenges, a variety of model compression techniques have emerged. These include quantization, which reduces the precision of model weights; knowledge distillation, which transfers learned knowledge from a large

teacher model to a smaller student model; and various forms of pruning.[3] Among these, weight pruning stands out as a prominent and effective approach. It involves the selective removal of redundant parameters, effectively setting a proportion of weights to zero, thereby reducing the overall model size and computational load.[2]

A paramount constraint for any pruning methodology is the imperative to achieve significant compression *without* compromising the model's performance. This "quality preservation" objective becomes increasingly challenging as higher sparsity levels are pursued, often leading to a delicate trade-off between model size reduction and accuracy.[2]

This report introduces and critically evaluates "Knowledge-Guided Dynamic Sparsity" (KG-DS), a novel, state-of-the-art methodology for LLM weight pruning. KG-DS is founded on the principle of creating a dynamic, input-aware sparse version of the model. This is achieved by identifying and protecting a crucial "knowledge core" while dynamically activating other task-relevant weights during inference. The proposed methodology aims to deliver substantial efficiency gains, encompassing both reduced memory usage and accelerated inference speed, while maintaining or potentially even exceeding the performance of the original, dense model.

The objectives of this research proposal are therefore multi-fold:

- To provide a detailed technical exposition of the KG-DS methodology.
- To contextualize KG-DS within the current landscape of LLM pruning research, drawing parallels and distinctions with existing state-of-the-art approaches.
- To conduct a thorough feasibility analysis of KG-DS, identifying its inherent strengths and critical challenges based on empirical evidence and theoretical considerations from the literature.
- To propose concrete research directions and enhancements aimed at advancing KG-DS or fostering the development of related dynamic sparsity paradigms.
- To outline a robust and comprehensive evaluation methodology to rigorously validate the claims of quality preservation and efficiency gains for KG-DS.

## 2. Understanding Knowledge-Guided Dynamic Sparsity (KG-DS)

The Knowledge-Guided Dynamic Sparsity (KG-DS) methodology represents a significant conceptual departure from conventional static pruning techniques. Its core

philosophy centers on creating a dynamically reconfigurable LLM where the active parameters are adapted in real-time based on the input prompt. This adaptive mechanism is built upon a carefully identified and protected "knowledge core," ensuring foundational capabilities are preserved while task-specific resources are dynamically allocated.

## Phase 1: Knowledge Core Identification (Offline)

This initial phase is a one-time, pre-computation process designed to identify and safeguard the most critical weights that encapsulate the fundamental knowledge and core capabilities of the LLM.

### Knowledge-Rich Probing

The process commences by exposing the full LLM to a diverse, curated dataset specifically designed to probe its core knowledge domains. This includes tasks requiring factual recall, commonsense reasoning, and deep linguistic understanding. The selection and curation of this "knowledge-rich" dataset are paramount. The quality and diversity of this dataset directly influence the robustness and generalizability of the "knowledge core" that is subsequently identified. If the dataset used for probing is incomplete or contains biases, the protected core might not truly encapsulate all fundamental knowledge. This could lead to a degradation in performance when the model encounters unseen or out-of-distribution tasks, even with the dynamic masks applied later. Therefore, the integrity of this initial probing phase is a critical determinant for KG-DS's ability to preserve quality across diverse applications. This step aligns with existing practices in pruning research that utilize small calibration datasets for importance estimation.[7] Future research should focus on developing robust methodologies for defining and curating optimal "knowledge-rich" datasets, potentially leveraging techniques such as active learning or advanced data diversity metrics to ensure comprehensive coverage.

### Gradient-Activation Saliency Mapping

Following the probing, a composite saliency score is computed for each weight within the LLM. This score is a sophisticated metric that integrates multiple factors: the weight's magnitude, the activation of its connected neurons, and, crucially, the gradient flow. The inclusion of gradient information is a key differentiator, aiming to identify weights that are not merely highly active but are also critical for the model's learning and generalization capabilities. This approach represents a significant advancement over simpler metrics like magnitude-only pruning, which often lead to substantial performance loss at higher sparsity levels.[3]

This saliency mapping aligns with recent state-of-the-art pruning methods. For instance, Wanda++ explicitly leverages "decoder-block-level regional gradients" and has demonstrated superior performance compared to Wanda, which relies solely on layer-wise weight and activation information, achieving up to a 32% perplexity improvement.[10] Similarly, GBLM-Pruner incorporates gradient information into its pruning criteria, showing performance gains over Wanda while maintaining comparable computational efficiency.[12] The feasibility of incorporating gradient information, despite the prohibitive cost of full-model backpropagation for LLMs [12], is supported by the success of methods like Wanda++. These methods demonstrate that localized or "regional" gradient computations, performed for a single decoder block at a time, can be both efficient and effective.[13] Moreover, gradients can be computed once and reused, further optimizing the process.[12] Therefore, the success of KG-DS's saliency mapping depends on adopting these efficient, localized gradient computation strategies, which are essential to uphold its claim of being "retraining-free" by avoiding computationally intensive full-model backpropagation.

**Structural Knowledge Preservation**

The calculated saliency scores are then utilized to identify and preserve entire structural components of the LLM. This includes elements such as attention heads and feed-forward network (FFN) neurons that are consistently activated or deemed critical during knowledge-intensive tasks. This process culminates in the creation of a "knowledge-core" mask, which designates these identified components as protected from aggressive pruning. This hybrid approach, combining different granularities of pruning, is hypothesized to yield superior results. This strategy is consistent with the current trend towards structured pruning, which is inherently more hardware-friendly

on general devices and can lead to tangible latency reductions.[1] Research, such as that on CFSP and LLM-Pruner, has demonstrated the efficacy of structured pruning, particularly highlighting that FFNs often exhibit higher structural sparsity and are more amenable to pruning than Multi-Head Attention (MHA) modules.[1]

While KG-DS defines a static "knowledge-core" mask based on this offline analysis, the broader paradigm of dynamic sparsity and non-uniform sparsity allocation suggests an even more advanced potential. Instead of a rigidly fixed structural core, the granularity and composition of the protected core itself could potentially be dynamically adapted or learned over time, or across broader task families. For example, Týr-the-Pruner employs an evolutionary search to identify optimal *non-uniform* sparsity distributions across different layers and components of an LLM.[15] This suggests a future direction for KG-DS where the "structural knowledge preservation" is not entirely static but can evolve or be fine-tuned based on evolving task requirements or model usage patterns, leading to an "adaptive structural knowledge preservation" mechanism.

### Phase 2: Dynamic Sparsity via Input-Modulated Pruning (Online)

This phase represents the truly innovative aspect of KG-DS, where the LLM's sparsity pattern is dynamically reconfigured in real-time during inference based on the user's input prompt.

### Instruction-Type Classification

Upon receiving a prompt, a lightweight classifier, trained on top of the LLM's own embeddings, rapidly categorizes the input into a predefined set of task types (e.g., "creative writing," "code generation," "question answering"). The concept of input-dependent model activation is supported by existing research, which indicates that dynamically adjusting model components based on input can be more effective than relying on a unified, static model.[16] However, the accuracy and robustness of this "lightweight classifier" are critically important. LLMs are known to be highly sensitive to prompt variations; even minor alterations in phrasing, including typos or semantic nuances, can lead to substantial performance degradation.[17] Furthermore, LLMs can

struggle with understanding complex, multi-step instructions, interpreting subtle nuances, or handling inconsistent information.[20]

The entire dynamic sparsity mechanism of KG-DS fundamentally hinges on the flawless operation of this classifier. If the classifier miscategorizes a prompt due to subtle phrasing, unexpected input, or even adversarial perturbations, the incorrect task-specific mask will be applied. This would directly lead to a degradation in the model's performance, contradicting the core objective of "quality preservation." Therefore, the "lightweight" nature of this classifier might represent a significant trade-off for its robustness in real-world scenarios. Extensive research is required to ensure this classifier is exceptionally robust to the diverse and often unpredictable nature of user prompts. This might necessitate more sophisticated training methodologies, such as self-denoising or representation alignment [17], or even a small, dedicated fine-tuning step, which would then challenge the "retraining-free" claim of KG-DS.

**Task-Specific Sparsity Masks**

For each identified task category, a corresponding pre-computed "sparsity mask" is readily available. These masks are generated offline and specifically designate which weights, beyond the protected "knowledge core," are most relevant and thus should remain active for that particular task.

**Dynamic Mask Application**

The appropriate task-specific mask is then applied "on the fly" to the LLM. This means that for every given input, a unique, sparse version of the LLM is instantiated and utilized in real-time. This approach is conceptually aligned with the emerging paradigm of dynamic sparsity, which emphasizes activating different neural pathways based on input, moving beyond static feature selection.[21] This paradigm underpins advanced architectures like sparse Mixtures of Experts (MoE) and enables conditional computation, where the effective model size adapts to the complexity of the input.[21]

While the conceptual power of dynamic sparsity is undeniable, its "on-the-fly"

implementation for LLM weights presents a significant practical challenge. Current hardware architectures are predominantly optimized for dense matrix multiplications or fixed sparsity patterns, such as N:M sparsity, which offers a balance between precision and hardware efficiency.[2] Applying a truly dynamic, input-dependent mask in real-time could introduce substantial overhead dueating to inefficient memory access patterns, increased cache misses, and the potential need for specialized sparse computation kernels that may not be universally available or optimally integrated. For instance, even dynamic data masking, a related concept, has been shown to add latency and hinder workflows in real-world applications.[23] While some dynamic data masking approaches claim no performance penalty

*during* database operations once data is already masked [24], the act of dynamically

*applying* a weight mask within an LLM inference pipeline is a distinct challenge. Unstructured pruning, which results in irregular patterns, typically requires specialized hardware support to achieve practical speed-ups.[1] KG-DS's dynamic masks, even if structured, introduce a level of variability that standard hardware might not efficiently exploit. Therefore, the "significant inference speedup" advantage claimed by KG-DS is heavily contingent on the efficiency of this dynamic mask application and the availability of hardware specifically optimized for such dynamic sparse computations. This necessitates a strong emphasis on hardware-software co-design or the development of innovative software optimizations to minimize runtime overhead.

**The KG-DS Advantage (as proposed by the user's AI)**

The conceptualized KG-DS methodology posits several key advantages over existing pruning techniques:

- **Quality Preservation:** By safeguarding a "knowledge core" and dynamically allocating resources to the most relevant parts of the network for a given task, KG-DS is expected to maintain, and in some specialized domains, potentially exceed the performance of the original, dense model.
- **Significant Efficiency Gains:** At any given time, a large portion of the model's weights are inactive, leading to a substantial reduction in memory usage and a potential for significant inference speedup, particularly with hardware designed for sparse computations.
- **Versatility and Adaptability:** KG-DS aims to create a highly adaptable model

capable of effectively handling a wide range of tasks by dynamically reconfiguring its active parameters.

- **Retraining-Free:** The core methodology is stated to not require expensive retraining of the LLM. The lightweight classifier for instruction-type identification can be trained efficiently with minimal data.

The claim of being "retraining-free" requires careful consideration. While KG-DS asserts that it avoids expensive retraining, many state-of-the-art pruning methods, even those described as "one-shot" like SparseGPT [9], often implicitly or explicitly rely on some form of "recovery fine-tuning" or "weight updates" to compensate for performance loss after pruning.[3] For instance, LLM-Pruner leverages Low-Rank Adaptation (LoRA) for efficient performance recovery.[14] The KG-DS description itself acknowledges that "the high cost of retraining these massive models has rendered such approaches impractical," implying that some methods

*do* necessitate retraining. If the "lightweight classifier" or the generation of task-specific masks involves any form of learning or parameter updates, it technically constitutes a form of "training," even if minimal. Therefore, the "retraining-free" assertion likely refers specifically to the avoidance of *full model fine-tuning*. If any component, such as the instruction classifier, requires training, its computational cost and impact on the overall efficiency of the "retraining-free" paradigm must be thoroughly analyzed and minimized. This also raises the question of whether KG-DS would benefit from sparsity-aware fine-tuning techniques like Sparsity Evolution Fine-Tuning (SEFT) or the Axolotl/LLM Compressor framework, which are designed to preserve sparsity during fine-tuning, addressing the issue of traditional fine-tuning reverting models to dense states.[5]

---

## 3. Current State-of-the-Art in LLM Weight Pruning

The field of LLM weight pruning has rapidly evolved, driven by the need to make increasingly large models more deployable and efficient. Current methodologies can be broadly categorized based on the granularity of pruning and the criteria used for identifying redundant parameters.

## Categorization of Pruning Methods

- **Unstructured Pruning:** This approach involves removing individual weights from the neural network. While it can achieve very high levels of sparsity and often maintain performance effectively at these high ratios, unstructured sparsity typically results in highly irregular patterns.[2] These irregular patterns are challenging to optimize on standard hardware architectures, meaning that practical speed-ups often necessitate specialized hardware support.[1]
- **Structured Pruning:** In contrast to unstructured methods, structured pruning removes entire groups of parameters, such as full attention heads, feed-forward network (FFN) neurons, or even entire layers.[4] This approach imposes structured sparsity, which is significantly more hardware-friendly on general-purpose devices and can directly lead to reduced latency.[1] Research indicates that FFN modules, due to their inherent characteristics, generally exhibit higher structural sparsity and are more amenable to structured pruning compared to Multi-Head Attention (MHA) layers.[1]
- **Semi-Structured Pruning (e.g., N:M sparsity):** This category represents a hybrid approach, retaining a fixed number of non-zero elements (N) within every group of M elements. N:M sparsity strikes an optimal balance between maintaining model precision and achieving hardware efficiency. This specific pattern can accelerate both matrix multiplication and memory access, thereby enhancing the performance of both pre-filling and decoding processes on off-the-shelf GPUs.[2]

## Key Pruning Criteria for Importance Evaluation

The effectiveness of any pruning method hinges on its ability to accurately identify and remove redundant weights while preserving critical ones. Various criteria have been developed for this "importance estimation" phase:

- **Magnitude-based Pruning:** This is the simplest and earliest approach, where weights with the smallest absolute magnitudes are pruned.[3] While straightforward to implement, magnitude-based pruning often leads to significant performance degradation, particularly at higher sparsity levels, as it does not account for the functional importance of weights in the context of network activations or gradients.[9]

- **Activation-based Pruning (e.g., Wanda):** More advanced methods consider the contribution of weights in relation to their input activations. Wanda, for instance, prunes weights based on the smallest product of the weight's magnitude and the norm of the corresponding input activation.[3] This approach has proven more effective than magnitude-only pruning.[3] Frameworks like CFSP also leverage activations as an importance criterion, citing their low computational overhead as they can be obtained with a single forward pass.[1]
- **Gradient-based Pruning (e.g., SparseGPT, LLM-Pruner, Wanda++, GBLM-Pruner):** The most sophisticated pruning criteria incorporate gradient information, which directly reflects a weight's contribution to the model's loss function and thus its importance for learning and generalization.
  - **SparseGPT:** This is a post-training, one-shot, layer-wise pruning method that utilizes approximations of the Hessian matrix to minimize reconstruction error. It is particularly effective for very large models and can achieve significant sparsity (e.g., 50% for OPT-175B) with minimal accuracy loss and without requiring extensive fine-tuning.[3] SparseGPT prunes weights incrementally and performs local updates to maintain the input-output relationship of each layer, avoiding the need for global gradient information.[25]
  - **LLM-Pruner:** This method focuses on structural pruning, identifying and removing non-critical coupled structures based on gradient information.[14] It aims for task-agnostic compression and can efficiently recover performance post-pruning using Low-Rank Adaptation (LoRA) with minimal data and training time.[14]
  - **Wanda++:** This novel framework significantly advances gradient-based pruning by leveraging "decoder-block-level regional gradients." It has demonstrated substantial performance improvements over Wanda, which relies solely on layer-wise weight and activation data, while maintaining high efficiency. The use of regional gradients makes gradient calculation feasible even for very large LLMs by localizing the backpropagation process to individual decoder blocks.[10]
  - **GBLM-Pruner:** This approach also integrates gradient information into its pruning criteria, showing performance gains compared to Wanda while maintaining a comparable computational footprint. A key efficiency aspect is that gradients only need to be computed once and can be reused across iterative pruning steps.[12]

The trajectory of research in saliency metrics for pruning demonstrates a clear progression. Early methods were simplistic, relying on static properties like magnitude. Subsequent advancements incorporated dynamic, input-dependent

behaviors through activation-aware criteria. The current state-of-the-art increasingly integrates gradient information, which captures a weight's direct influence on the model's learning objective. KG-DS's "Gradient-Activation Saliency Mapping," which combines magnitude, activation, and gradient flow, directly reflects this convergence towards a more holistic understanding of weight importance. This indicates that effective pruning necessitates considering a weight's static value, its dynamic response to inputs, and its fundamental contribution to the model's overall function and generalization capabilities.

**Post-training vs. In-training Pruning and the Role of Fine-tuning**

Pruning methods can also be classified by when the sparsity is introduced relative to the model's training:

- **One-shot Pruning:** This highly desirable paradigm aims to prune models without requiring extensive, computationally prohibitive retraining.[2] SparseGPT and Wanda are prominent examples of successful one-shot post-training pruning techniques for LLMs.[2]
- **Challenges of Retraining:** Full retraining of pruned LLMs is a major computational and logistical burden. It is exceptionally expensive and time-consuming, often consuming extensive compute and memory resources. Furthermore, retraining can inadvertently compromise the LLM's essential language understanding and reasoning abilities, which are significantly more difficult to restore than simpler metrics like perplexity.[2]
- **Recovery Fine-tuning (e.g., LoRA):** To mitigate the performance degradation that can occur after pruning, many methods employ a recovery step. Low-Rank Adaptation (LoRA) is a widely adopted parameter-efficient fine-tuning (PEFT) technique. It modifies only a small number of additional parameters, allowing for efficient recovery of performance with minimal training overhead compared to full fine-tuning.[1]
- **Sparsity-aware Fine-tuning (e.g., SEFT, Axolotl/LLM Compressor):** A more recent development addresses the challenge that traditional fine-tuning methods (even PEFTs like LoRA) can inadvertently revert pruned models to a dense state, undermining the efficiency gains of sparsity.[5] Sparsity Evolution Fine-Tuning (SEFT) is a novel method specifically designed for sparse LLMs. Inspired by Dynamic Sparse Training (DST), SEFT enables dynamic sparsity evolution during fine-tuning, aiming to adapt the sparse connectivity to downstream tasks while

maintaining the desired sparsity level.[5] Similarly, open-source solutions like Axolotl and LLM Compressor offer workflows that combine pruning, sparse fine-tuning, and quantization, explicitly preserving the sparsity structure during the fine-tuning phase to recover accuracy without compromising efficiency.[26]

**Dynamic Sparsity and Conditional Computation**

Beyond static pruning, an emerging and highly promising paradigm is dynamic sparsity. This approach shifts from traditional static feature selection in neural representations to a dynamic model where different neural pathways are activated depending on the input.[21] This concept is a driving force behind new architectures for foundation models, notably sparse Mixtures of Experts (MoE), which allow for conditional computation. This enables the model to adapt its effective architecture or representation size based on the complexity or type of the input, offering advantages in incorporating structural constraints and potentially achieving the performance of dense models with significantly reduced active parameters.[21]

KG-DS fundamentally operates as a hybrid of pruning and conditional computation. By dynamically activating different parts of the model (via task-specific masks) based on the input type, it effectively creates a form of "mixture of experts," where each "expert" corresponds to a specialized sparse sub-network for a particular task. This represents a significant conceptual advancement beyond conventional static pruning techniques. Consequently, the evaluation and optimization of KG-DS must extend beyond traditional pruning metrics to also consider challenges and performance indicators typically associated with conditional computation and MoE architectures, such as routing efficiency, expert specialization, and load balancing across the dynamically activated components.

---

**Table 1: Comparison of State-of-the-Art LLM Pruning Methods**

| Method Name | Pruning Type | Saliency Criterion | Retraining/Fine-tuning Required | Key Advantage | Key Challenge | Hardware Friendliness | How KG-DS Relates |
|---|---|---|---|---|---|---|---|
| Magnitu | Unstruct | Weight | Often | Simplicit | Significa | Low | Baseline |

| de Pruning | ured | Magnitude | requires fine-tuning for high sparsity | y, ease of implementation | nt performance loss at high sparsity [9] | (irregular patterns) [4] | ; KG-DS improves upon by adding activation and gradient. |
|---|---|---|---|---|---|---|---|
| Wanda | Unstructured | Weight Magnitude * Input Activation Norm | No (one-shot) [3] | More effective than magnitude-only [3] | Performance degradation at higher sparsities [3] | Low (irregular patterns) [1] | KG-DS extends this with gradient information. |
| SparseGPT | Unstructured/Semi-structured | Hessian-based Reconstruction Error | No (one-shot) [9] | High sparsity with minimal accuracy loss for large models [9] | Computationally intensive for large models (though efficient for its class) [25] | Moderate (can support N:M) [2] | KG-DS shares "one-shot" goal; its saliency is different. |
| LLM-Pruner | Structural | Gradient Information | LoRA for recovery [14] | Task-agnostic structural compression [14] | Requires recovery fine-tuning [14] | High (structural) [1] | KG-DS also uses structural pruning and gradient. |
| Wanda++ | Structural (regional) | Regional Gradients + Weight/Activation | Optional LoRA compatibility [10] | Significant perplexity improvement over Wanda, efficient gradient | Still a post-training approach; regionality might miss global depende | High (structural) [1] | KG-DS's gradient inclusion aligns directly with Wanda++. |

| | | | | use [10] | ncies [15] | | |
|---|---|---|---|---|---|---|---|
| CFSP | Structural | Coarse-to-Fine Activation Information | LoRA for recovery [1] | Efficient (single forward pass), hardware-friendly structured pruning [1] | Requires recovery fine-tuning [1] | High (structural) [1] | KG-DS's structural preservation and activation use align with CFSP. |
| K-prune | Structural | Knowledge Preservation (iterative) | No (retraining-free) [28] | Accurate retraining-free structured pruning [28] | Iterative process might incur overhead [28] | High (structural) [1] | KG-DS shares "knowledge preservation" and "retraining-free" goals. |
| OTO | Structured | Zero-Invariant Groups (ZIGs) | No (one-shot, no fine-tuning) [29] | Guarantees no output change post-pruning [29] | Requires specific ZIG identification; applicability to all LLM structures [29] | High (structural) [1] | KG-DS aims for similar "no fine-tuning" outcome but via different mechanism. |
| SEFT | Unstructured/Structured (Dynamic) | Sensitivity-based Pruning | Sparsity Evolution Fine-Tuning [5] | Preserves sparsity during fine-tuning, dynamic adaptation [5] | Can lead to denser models if not carefully controlled [5] | Varies (dynamic) [5] | KG-DS's dynamic nature aligns; SEFT offers a sparsity-aware fine-tuning |

| | | | | | | | solution. |
|---|---|---|---|---|---|---|---|

**Value of Table 1:** This table provides a concise yet comprehensive overview of prominent LLM pruning methods, categorizing them by their core mechanisms, strengths, and limitations. For each method, its relationship to KG-DS is explicitly stated. This allows for a quick comparative analysis, highlighting where KG-DS aligns with state-of-the-art practices (e.g., gradient-based saliency, structured pruning, one-shot objective) and where it proposes novel advancements (e.g., dynamic, input-modulated sparsity). The table serves to contextualize KG-DS within the broader research landscape, demonstrating its theoretical grounding while also pinpointing areas where its claims (e.g., "retraining-free") might face practical challenges or require further refinement in comparison to established techniques.

---

# 4. Feasibility Analysis of Knowledge-Guided Dynamic Sparsity (KG-DS)

The Knowledge-Guided Dynamic Sparsity (KG-DS) methodology presents an innovative approach to LLM optimization, combining established pruning techniques with a novel dynamic activation paradigm. A thorough assessment of its feasibility requires a detailed examination of its strengths, alignment with current research, and the significant challenges it must overcome for real-world applicability.

## 4.1. Strengths and Alignment with Current Research

KG-DS leverages several concepts that are at the forefront of LLM compression research, indicating a strong theoretical foundation and potential for effectiveness:

- **Advanced Saliency Estimation:** The proposed "Gradient-Activation Saliency Mapping" is a robust approach to identifying critical weights. It moves beyond simpler magnitude-based pruning, which is known to be less effective at higher sparsity levels.[9] By incorporating both activation magnitudes (similar to Wanda) and gradient flow, KG-DS aligns with recent breakthroughs such as Wanda++ and GBLM-Pruner.[10] These methods have demonstrated that gradient information

provides crucial insights for pruning, leading to significantly better performance retention.[10] The ability to compute these gradients at a regional or layer-wise level, as shown by Wanda++, makes this computationally feasible even for very large LLMs, circumventing the prohibitive cost of full-model backpropagation.[12] This multi-faceted importance criterion is a key strength for preserving model quality.

- **Structured Pruning for Hardware Efficiency:** KG-DS's emphasis on "Structural Knowledge Preservation" aligns with the industry's shift towards structured pruning. Unlike unstructured pruning, which often results in irregular sparsity patterns that are difficult to accelerate on general-purpose hardware, structured pruning removes entire components (e.g., FFN neurons, attention heads).[4] This results in more hardware-friendly models that can achieve tangible latency reductions on standard devices.[1] The observation that FFNs are particularly amenable to structured pruning provides a clear target for efficient compression.[1]

- **Embracing Dynamic Sparsity and Conditional Computation:** The core innovation of KG-DS lies in its "Dynamic Sparsity via Input-Modulated Pruning." This concept is a recognized and rapidly evolving paradigm in machine learning, driving the development of architectures like sparse Mixtures of Experts (MoE).[21] By dynamically activating different parts of the model based on input, KG-DS aims to achieve conditional computation, adapting the model's active size and complexity to the specific task at hand. This offers the potential for significant efficiency gains without sacrificing the performance breadth of a dense model.[22] The idea of input-dependent model activation is supported by research showing its effectiveness compared to static, unified approaches.[16]

- **"Retraining-Free" Objective:** KG-DS's stated goal of being "retraining-free" is highly attractive for LLMs, given the immense computational cost and time associated with full model fine-tuning after pruning.[2] This aligns with the objectives of one-shot pruning methods like SparseGPT and K-prune, which aim to achieve high sparsity with minimal or no post-pruning recovery steps.[9]

## 4.2. Critical Challenges and Feasibility Concerns

Despite its promising aspects, KG-DS faces several significant challenges that must be rigorously addressed for its practical viability and to fully realize its claimed advantages:

- **Robustness of Instruction-Type Classification:** The entire dynamic sparsity

mechanism of KG-DS is critically dependent on the accuracy and robustness of the "lightweight classifier" that categorizes input prompts. LLMs are notoriously sensitive to prompt variations; even minor alterations in phrasing, syntax, or the presence of noise can lead to substantial performance degradation or misinterpretation of intent.[17] Furthermore, LLMs can struggle with complex, multi-step instructions, subtle nuances, or inconsistent information.[20] If this classifier miscategorizes a prompt, the wrong task-specific sparsity mask will be applied, directly leading to suboptimal performance or outright failure, thereby undermining the "quality preservation" claim. The "lightweight" nature of this classifier, while beneficial for inference speed, might compromise its ability to handle the full spectrum of real-world user inputs robustly. Ensuring its resilience against diverse, potentially ambiguous, or even adversarial prompts will require significant research, potentially involving advanced robustness techniques like self-denoising or representation alignment, which could add to its complexity and training requirements.[17]

- **Real-time Dynamic Mask Application Overhead:** While conceptually powerful, the "on-the-fly" application of dynamic sparsity masks during inference presents a substantial practical hurdle. Current hardware architectures are primarily optimized for dense matrix operations or fixed, structured sparsity patterns (e.g., N:M sparsity).[2] Applying a truly dynamic, input-dependent mask in real-time could introduce significant runtime overheads. This stems from factors such as increased memory access latency due to irregular data loading patterns, cache misses when switching between different sparse sub-networks, and the computational cost of dynamically reconfiguring the computational graph or kernel execution. Analogous challenges are observed in dynamic data masking, which has been shown to add latency and hinder workflow efficiency in real-world applications.[23] Achieving the "significant inference speedup" promised by KG-DS will necessitate either highly optimized software implementations for dynamic sparse operations or the availability of specialized hardware (e.g., custom ASICs or reconfigurable FPGAs) explicitly designed to efficiently handle dynamic sparsity patterns.[1] Without such optimizations or hardware support, the theoretical gains from reduced parameter count might be negated by the overhead of dynamic mask application.
- **"Retraining-Free" Claim Nuance:** KG-DS states its core methodology is "retraining-free." While this likely refers to avoiding expensive *full model fine-tuning*, the process of generating "task-specific sparsity masks" and training the "lightweight classifier" still involves a form of learning or optimization. If these offline processes are not sufficiently efficient or require substantial data and compute, the overall "retraining-free" advantage diminishes. Furthermore, even in

"one-shot" pruning, some form of "recovery fine-tuning" (e.g., via LoRA) or "weight updates" is often necessary to regain lost performance.[3] The challenge of maintaining sparsity during any subsequent fine-tuning (e.g., for adapting to new tasks or domains) is also critical, as traditional fine-tuning methods can revert sparse models to dense states.[5] This implies that if any post-pruning adaptation is needed, it must employ sparsity-aware fine-tuning techniques like SEFT or Axolotl/LLM Compressor.[5]

- **Generalization of the "Knowledge Core":** The effectiveness of the static "knowledge core" relies heavily on the diversity and representativeness of the "knowledge-rich probing" dataset. If this dataset does not adequately cover the breadth of the LLM's intended applications or future emergent capabilities, the protected core might inadvertently omit crucial knowledge. This could lead to a degradation in performance on tasks not well-represented in the initial probing, even with dynamic task-specific masks. The risk of pruning removing connections crucial for specific downstream tasks is a known challenge in LLM compression.[5] This highlights a potential limitation of a fixed knowledge core in a truly versatile and adaptable LLM.

# 5. Research Proposal: Enhancing Knowledge-Guided Dynamic Sparsity

This section outlines a comprehensive research program aimed at addressing the identified challenges of KG-DS and advancing the state-of-the-art in dynamic LLM pruning. The proposed research focuses on refining the core mechanisms of KG-DS, exploring novel architectural and algorithmic solutions, and establishing a robust evaluation framework.

### 5.1. Refinement of Knowledge Core Identification

The efficacy of KG-DS hinges on the precise identification and preservation of the foundational "knowledge core." This research stream will focus on enhancing the robustness and generalizability of this critical component.

### 5.1.1. Advanced Knowledge-Rich Probing Strategies

Current approaches rely on a "diverse, curated dataset" for probing. This research will investigate methodologies for systematically constructing or selecting optimal "knowledge-rich" datasets. This includes:

- **Diversity Metrics and Active Learning:** Developing quantitative metrics to assess the diversity and coverage of probing datasets across various knowledge domains (e.g., factual, commonsense, linguistic, reasoning). This will involve exploring active learning strategies to iteratively select the most informative samples for probing, maximizing the impact of a limited calibration dataset. The goal is to ensure the identified knowledge core is truly comprehensive and not biased towards specific data distributions encountered during probing.
- **Adversarial Probing:** Introducing adversarial examples or out-of-distribution inputs during the probing phase to identify and protect weights crucial for model robustness and generalization under challenging conditions. This aims to safeguard the model against unexpected inputs that might otherwise lead to misclassification by the dynamic sparsity mechanism.

### 5.1.2. Refined Gradient-Activation Saliency Mapping

While KG-DS incorporates gradient flow, further optimization of the saliency metric is possible:

- **Dynamic Weight Update during Pruning:** Investigate the integration of efficient layer-wise weight updates during the pruning process, similar to ADMM-based algorithms [3] or SparseGPT's local updates.[25] This would aim to minimize reconstruction error at each pruning step, potentially leading to a more accurate "knowledge core" and reducing the need for any post-pruning recovery.
- **Second-Order Information beyond Gradients:** Explore the utility of second-order information (e.g., Hessian approximations or Fisher information) in a computationally feasible manner, building on methods like SparseGPT.[3] While full Hessian computation is prohibitive, localized or approximate second-order information could provide a more nuanced understanding of weight importance and interdependencies, further refining the saliency scores. This could lead to a

more effective identification of weights critical for complex emergent behaviors, which are often harder to restore after pruning.[2]

### 5.1.3. Adaptive Structural Knowledge Preservation

Moving beyond a static "knowledge-core" mask, this research will explore dynamic allocation of structural importance:

- **Non-Uniform Structural Allocation:** Inspired by methods like Týr-the-Pruner which use evolutionary search for optimal non-uniform sparsity distributions across layers and components [15], this research will investigate how the "knowledge core" itself can be dynamically composed or its granularity adapted. Instead of a single fixed core, different "sub-cores" could be activated based on broader task families or domain shifts, allowing for more flexible resource allocation. This would entail learning a meta-mask that governs the composition of the knowledge core.
- **Inter-Component Saliency:** Develop methods to assess the collective importance of structural components (e.g., groups of attention heads or FFN layers) and their interdependencies, rather than just individual component saliency. This would ensure that pruning decisions preserve critical functional pathways within the LLM, aligning with the insights from global pruning methods.[15]

### 5.2. Optimizing Dynamic Sparsity via Input-Modulated Pruning

This research stream will focus on mitigating the practical challenges associated with real-time dynamic mask application and ensuring the robustness of the instruction classification.

### 5.2.1. Robust and Efficient Instruction-Type Classification

The reliability of the lightweight classifier is paramount. This research will address its

critical vulnerabilities:

- **Robustness to Prompt Perturbations:** Implement and evaluate advanced techniques to enhance the classifier's robustness against diverse prompt variations, including subtle phrasing changes, typos, and adversarial inputs. This could involve training strategies such as self-denoising or representation alignment, which have shown promise in improving LLM robustness to perturbed instructions.[17]
- **Hierarchical Classification:** Instead of a flat classification of task types, explore a hierarchical approach. A coarse-grained classifier would first identify broad task categories, followed by finer-grained classifiers for sub-tasks. This could improve accuracy and allow for more nuanced mask selection.
- **Confidence-Aware Masking:** Implement a mechanism where the classifier provides a confidence score for its prediction. If confidence is low, a default, less aggressive sparsity mask (or even the full dense model) could be used as a fallback, preventing performance degradation due to misclassification.
- **Adaptive Classifier Training:** Investigate lightweight, sparsity-aware fine-tuning of the classifier itself, potentially using techniques like LoRA or SEFT, to adapt it to new task types or evolving user prompt patterns without incurring significant computational overhead.[5] This would address the nuance of the "retraining-free" claim.

### 5.2.2. Hardware-Aware Dynamic Mask Application

Addressing the real-time overhead of dynamic mask application is crucial for achieving practical speedups. This research will explore:

- **Software-Hardware Co-Design:** Collaborate with hardware architects to design or adapt existing hardware (e.g., specialized tensor cores, reconfigurable computing units) that can efficiently handle dynamic, irregular sparsity patterns. This includes optimizing memory access patterns and developing custom sparse computation kernels that minimize latency during mask application.[8]
- **Batching and Scheduling Optimizations:** Investigate dynamic batching and scheduling strategies for inference requests. By grouping similar instruction types, it might be possible to reduce the frequency of mask switching, thereby amortizing the overhead of dynamic mask application across multiple inferences.
- **Pre-fetching and Caching:** Explore intelligent pre-fetching and caching mechanisms for task-specific sparsity masks and their corresponding active

weights. Based on predicted upcoming tasks or user behavior, relevant sparse sub-networks could be loaded into faster memory proactively.

- **Progressive Sparsity Application:** Instead of a single "on-the-fly" mask application, explore a progressive approach where a baseline sparse model is always active, and additional task-specific weights are dynamically "activated" or "deactivated" in a more granular fashion, potentially reducing the overhead of full mask switching. This aligns with the concept of dynamic sparse attention.[30]

## 5.3. Comprehensive Evaluation Framework

Rigorous evaluation is essential to validate the claims of quality preservation and efficiency gains for KG-DS. The framework will extend beyond traditional perplexity metrics.

- **Perplexity and Downstream Task Performance:** Standard perplexity on language modeling benchmarks will be measured, but the primary focus will be on downstream task performance across a diverse set of benchmarks (e.g., MMLU, HELM, GLUE, SuperGLUE). This will include tasks requiring factual recall, reasoning, code generation, and creative writing, directly testing the versatility and quality preservation claims of KG-DS across its intended dynamic applications.
- **Human-in-the-Loop Evaluation:** Given that traditional metrics like BLEU/ROUGE often fail to capture the semantic nuance of LLM outputs [31], human evaluation and "LLM-as-a-judge" methodologies (e.g., G-Eval) will be critical. This will assess answer correctness, semantic similarity, relevance, and the absence of hallucinations, providing a more accurate reflection of real-world performance.[31]
- **Efficiency Metrics:**
  - **Memory Footprint:** Quantify the active memory usage of KG-DS during inference compared to the dense model and statically pruned models at various sparsity levels.
  - **Inference Latency/Throughput:** Measure end-to-end inference latency and throughput across different task types and input complexities, specifically quantifying the overhead introduced by the dynamic mask application and instruction classification. This will be measured on various hardware configurations (e.g., consumer GPUs, data center GPUs, edge devices) to assess real-world applicability.
  - **Energy Consumption:** Estimate the energy consumption per inference,

particularly important for sustainable AI deployment.[4]

- **Robustness Metrics:** Develop and apply specific metrics to evaluate the robustness of the instruction classifier and the overall KG-DS system to prompt perturbations, out-of-distribution inputs, and potential adversarial attacks.
- **Ablation Studies:** Conduct comprehensive ablation studies to understand the individual contributions of the "knowledge core," gradient-activation saliency, and dynamic mask application to overall performance and efficiency. This will also involve analyzing the impact of different "knowledge-rich" probing datasets and classifier architectures.

# 6. Conclusions and Recommendations

The proposed Knowledge-Guided Dynamic Sparsity (KG-DS) methodology represents a compelling and innovative direction for optimizing Large Language Models. Its core premise, combining a protected knowledge core with input-modulated dynamic sparsity, aligns with the cutting edge of research in conditional computation and advanced pruning techniques. The incorporation of gradient information for saliency mapping is a particularly strong aspect, building on the successes of methods like Wanda++ and GBLM-Pruner.

However, the analysis reveals that the practical feasibility and full realization of KG-DS's advantages are contingent upon addressing several critical challenges. The robustness of the "lightweight classifier" responsible for instruction-type classification is a primary concern; its susceptibility to prompt variations could undermine the very quality preservation KG-DS aims to achieve. Furthermore, the "on-the-fly" application of dynamic masks introduces significant computational overheads that current hardware architectures are not optimally designed to handle, potentially negating the promised inference speedups. The "retraining-free" claim, while desirable, requires a nuanced understanding, as certain components or post-pruning adaptations may still necessitate efficient, sparsity-aware learning.

To advance KG-DS from a conceptual framework to a practical, state-of-the-art solution, the following research recommendations are paramount:

1. **Prioritize Robust Instruction Classification:** Invest in dedicated research to enhance the robustness of the instruction-type classifier against diverse and challenging real-world prompts. This includes exploring advanced training

techniques (e.g., self-denoising, representation alignment) and potentially developing confidence-aware or hierarchical classification mechanisms to mitigate misclassification risks.

2. **Investigate Hardware-Software Co-Design for Dynamic Sparsity:** A significant portion of future work must focus on minimizing the real-time overhead of dynamic mask application. This involves both software optimizations (e.g., intelligent batching, caching, progressive activation) and, critically, collaborative efforts with hardware designers to develop architectures explicitly optimized for dynamic, irregular sparsity patterns. Without this, the memory and speed benefits may remain theoretical.

3. **Refine Knowledge Core Identification with Adaptive Mechanisms:** While a static knowledge core is a good starting point, explore how the core itself can adapt or be dynamically composed based on broader task families or evolving model usage. This could involve learning non-uniform structural allocations and considering inter-component saliency to ensure comprehensive knowledge preservation across diverse domains.

4. **Rigorously Define and Validate "Retraining-Free" Efficiency:** Clearly delineate the scope of "retraining-free" and quantify the computational cost of all offline processes (e.g., mask generation, classifier training). If any post-pruning adaptation or continuous learning is required, ensure it employs sparsity-aware fine-tuning techniques to preserve the efficiency gains.

5. **Adopt a Holistic Evaluation Framework:** Move beyond traditional metrics to include extensive human-in-the-loop evaluations, LLM-as-a-judge methodologies, and comprehensive real-world efficiency measurements (latency, throughput, energy consumption) across diverse tasks and hardware. This will provide a more accurate and reliable assessment of KG-DS's true performance and practical utility.

By systematically addressing these challenges and pursuing the outlined research directions, Knowledge-Guided Dynamic Sparsity has the potential to become a pivotal methodology in the quest for more efficient, accessible, and sustainable Large Language Models, heralding a new era of optimized AI.

## Works cited

1. CFSP: An Efficient Structured Pruning Framework for LLMs with Coarse-to-Fine Activation Information - ACL Anthology, accessed June 20, 2025, https://aclanthology.org/2025.coling-main.626.pdf
2. Pruning Large Language Models with Semi-Structural Adaptive Sparse Training, accessed June 20, 2025,

https://ojs.aaai.org/index.php/AAAI/article/view/34592/36747
3. Fast and Optimal Weight Update for Pruned Large Language Models - arXiv, accessed June 20, 2025, https://arxiv.org/html/2401.02938v1
4. LLM Pruning for Enhancing Model Performance - Incubity by Ambilio, accessed June 20, 2025, https://incubity.ambilio.com/llm-pruning-for-enhancing-model-performance/
5. Leave it to the Specialist: Repair Sparse LLMs with Sparse Fine-Tuning via Sparsity Evolution - arXiv, accessed June 20, 2025, https://arxiv.org/html/2505.24037v1
6. The Efficiency Spectrum of Large Language Models: An Algorithmic Survey - arXiv, accessed June 20, 2025, https://arxiv.org/html/2312.00678v2
7. LLM pruning & distillation: Minitron approach - SuperAnnotate, accessed June 20, 2025, https://www.superannotate.com/blog/llm-pruning-distillation-minitron-approach
8. A Review on Edge Large Language Models: Design, Execution, and Applications - arXiv, accessed June 20, 2025, https://arxiv.org/html/2410.11845v2
9. SparseGPT: Massive Language Models Can be Accurately Pruned in One-Shot - arXiv, accessed June 20, 2025, https://arxiv.org/pdf/2301.00774
10. Wanda++: Pruning Large Language Models via Regional Gradients - arXiv, accessed June 20, 2025, https://arxiv.org/html/2503.04992v1
11. Wanda++: Pruning Large Language Models via Regional Gradients - OpenReview, accessed June 20, 2025, https://openreview.net/forum?id=WjnJf5ft0B&referrer=%5Bthe%20profile%20of%20Athanasios%20Mouchtaris%5D(%2Fprofile%3Fid%3D~Athanasios_Mouchtaris1)
12. Beyond Size: How Gradients Shape Pruning Decisions in Large ..., accessed June 20, 2025, https://openreview.net/forum?id=5BoXZXTJvL
13. WANDA++: PRUNING LARGE LANGUAGE ... - Amazon Science, accessed June 20, 2025, https://assets.amazon.science/47/f7/5021859e4427a31db7d0f2b6a75b/wanda-pruning-large-language-models-via-regional-gradients.pdf
14. [2305.11627] LLM-Pruner: On the Structural Pruning of Large Language Models - arXiv, accessed June 20, 2025, https://arxiv.org/abs/2305.11627
15. Týr-the-Pruner: Unlocking Accurate 50% Structural Pruning for LLMs via Global Sparsity Distribution Optimization - arXiv, accessed June 20, 2025, https://arxiv.org/pdf/2503.09657
16. Leveraging Self-Attention for Input-Dependent Soft Prompting in LLMs - arXiv, accessed June 20, 2025, https://arxiv.org/pdf/2506.05629
17. Robustness of Learning from Task Instructions - ResearchGate, accessed June 20, 2025, https://www.researchgate.net/publication/372917465_Robustness_of_Learning_from_Task_Instructions
18. Robustness of large language models in moral judgements - PMC, accessed June 20, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC12015570/
19. Enhancing LLM Robustness to Perturbed Instructions: An Empirical Study - arXiv, accessed June 20, 2025, https://arxiv.org/pdf/2504.02733

20. arXiv:2502.09101v1 [cs.HC] 13 Feb 2025, accessed June 20, 2025, https://arxiv.org/pdf/2502.09101
21. Dynamic Sparsity in Machine Learning | NeurIPS 2024 Tutorial, accessed June 20, 2025, https://dynamic-sparsity.github.io/
22. Dynamic Sparsity in Machine Learning: Routing Information through Neural Pathways, accessed June 20, 2025, https://neurips.cc/virtual/2024/tutorial/99527
23. LLM Data Masking: Silver Bullet or Double-Edged Sword? - Salesforce, accessed June 20, 2025, https://www.salesforce.com/blog/llm-data-masking/
24. Static vs Dynamic Data Masking - Tonic.ai, accessed June 20, 2025, https://www.tonic.ai/guides/static-vs-dynamic-masking
25. SparseGPT: Remove 100 billion parameters for free | Red Hat Developer, accessed June 20, 2025, https://developers.redhat.com/articles/2023/03/21/sparsegpt-remove-100-billion-parameters-free
26. Axolotl meets LLM Compressor: Fast, sparse, open | Red Hat Developer, accessed June 20, 2025, https://developers.redhat.com/articles/2025/06/17/axolotl-meets-llm-compressor-fast-sparse-open
27. Efficient Compression of Large Language Models using LLM-Pruner - YouTube, accessed June 20, 2025, https://www.youtube.com/watch?v=NcsiWxt55Xk
28. Accurate Retraining-free Pruning for Pretrained Encoder-based Language Models - arXiv, accessed June 20, 2025, https://arxiv.org/html/2308.03449v2
29. Only Train Once: A One-Shot Neural Network Training And Pruning Framework - NeurIPS, accessed June 20, 2025, https://proceedings.neurips.cc/paper_files/paper/2021/file/a376033f78e144f494bfc743c0be3330-Paper.pdf
30. [2506.11104] DAM: Dynamic Attention Mask for Long-Context Large Language Model Inference Acceleration - arXiv, accessed June 20, 2025, https://arxiv.org/abs/2506.11104
31. LLM Evaluation Metrics: The Ultimate LLM Evaluation Guide - Confident AI, accessed June 20, 2025, https://www.confident-ai.com/blog/llm-evaluation-metrics-everything-you-need-for-llm-evaluation