

Tangram: Accelerating Serverless LLM Loading through GPU Memory Reuse and Affinity

Wenbin Zhu
Shandong University
China
wenbinzhu@mail.sdu.edu.cn

Zhaoyan Shen
Shandong University
China
shenzhaoyan@sdu.edu.cn

Zili Shao
The Chinese University of Hong Kong
China
shao@cse.cuhk.edu.hk

Hongjun Dai
Shandong University
China
dahogn@sdu.edu.cn

Feng Chen
Indiana University Bloomington
USA
fchen25@iu.edu

Abstract

Serverless Large Language Models (LLMs) have emerged as a cost-effective solution for deploying AI services by enabling a “pay-as-you-go” pricing model through GPU resource sharing. However, cold-start latency, especially the model loading phase, has become a critical performance bottleneck, as it scales linearly with model size and severely limits the practical deployment of large-scale LLM services. This paper presents *Tangram*, a novel system that accelerates Serverless LLM loading through efficient GPU memory reuse. By leveraging the unused GPU memory to retain model parameters, Tangram significantly reduces model transfer time and cold-start latency. Its design includes three key components: unified GPU memory pool for tensor-level parameter sharing across models, on-demand KV cache allocation for dynamic memory management, and GPU-affinity-aware scheduling for maximizing resource utilization. These techniques collectively address the critical challenges of inefficient memory usage and the cold-start problem in Serverless LLM platforms. We have implemented a fully functional prototype, and experiments show that Tangram achieves up to 6.2× faster loading and reduces Time-To-First-Token (TTFT) during cold-start by 23–55% over state-of-the-art methods.

1 Introduction

In recent years, the demand for Large Language Model (LLM) services has rapidly surged across nearly all sectors. A critical challenge for many users and organizations is the prohibitively high deployment costs. A promising emerging solution is *Serverless LLM* [14], represented by Google Vertex AI [17], Amazon Bedrock [4], Alibaba PAI-EAS [3], and Hyperbolic [19]. As an online service platform, Serverless LLM provides on-demand services, allowing users to share GPU resources and deploy their LLM models. By multiplexing the expensive hardware infrastructure, the system dynamically switches between LLM models and reallocates GPU resources among users, enabling a flexible,

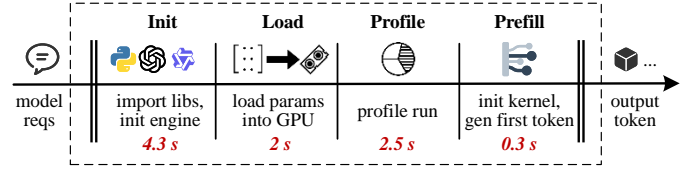


Figure 1. Multi-phase initialization of Serverless LLM and Time-To-First-Token (TTFT) composition for GPT-20B.

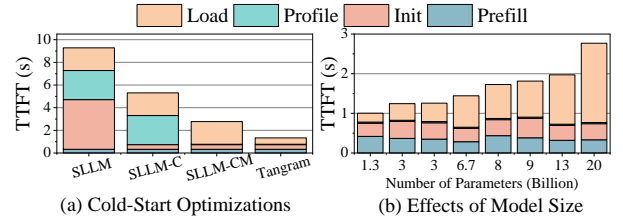


Figure 2. (a) TTFT breakdown for GPT-20B under different cold-start optimizations. (b) TTFT breakdown of SLLM-CM across models of varying parameter sizes.

efficient “pay-as-you-go” pricing model that significantly lowers the operation costs for users [2, 14, 18].

However, due to the ever-growing parameter size and the complexity of multi-phase initialization, an increasingly prominent challenge is that Serverless LLMs tend to suffer from significant model-switching latency, referred to as the *cold-start problem*, which severely impacts latency-sensitive applications [5, 13, 14, 47–49].

In a typical Serverless LLM deployment, initiating a new model instance involves four phases: *Init*, which configures the runtime environment, including Python libraries and inference engines; *Load*, which transfers the model parameters into GPU memory; *Profile*, which performs a profiling run to determine the maximum available KV cache size; and *Prefill*, which initializes the CUDA kernels and executes the initial forward pass. After these steps, the activated model generates its first inference token. Consequently, the *Time-to-First-Token* (TTFT) is primarily

determined by the cumulative latency of these phases, significantly affecting the responsiveness of Serverless LLM services and user experience.

Prior studies have made efforts to reduce TTFT. SLLM [14] improves *Load* efficiency through CPU-based model caching and parallel parameter loading; Several works [10, 12, 35, 44] exploit Checkpoint/Restore In Userspace (CRIU) to bypass redundant context initialization in *Init*; Medusa [49] uses offline materialization to pre-compute the information required for *Profile*; Tidal [13] creates adaptive CUDA context templates to reduce the latency of kernel loading in *Prefill*. Despite these notable advances, cold-start latency remains high, especially for large models.

To investigate the cold-start bottleneck of Serverless LLMs, we incrementally extended the widely used SLLM platform with state-of-the-art optimizations and evaluated models of varying sizes (ranging from OPT-1.3B to GPT-20B). As shown in Figure 2a, with CRIU-enabled checkpointing (*SLLM-C*) and further with offline materialization (*SLLM-CM*), the latencies of *Init* and *Profile* are significantly reduced. However, the latency of *Load* remains high. Figure 2b shows that as model size increases, the latencies of other phases remain relatively stable, whereas the *Load* latency scales almost linearly. This indicates that the *Load* latency is becoming the major performance bottleneck, which must be addressed.

The key to reducing *Load* latency is to minimize the amount of data transferred, which can be achieved by retaining model parameters as much as possible in GPU memory. This idea is motivated by two observations: First, in Serverless LLM, model requests exhibit a strong locality. Our analysis on Serverless LLM workloads reveals a repeated loading pattern. Certain models are predictably loaded and evicted. By retaining such models completely or partially in GPU memory, the volume of data transfer across PCIe can be reduced, thereby reducing the data loading time. Second, both model parameter sizes and KV cache requirements vary substantially across inference requests. In many cases, GPU memory is not fully used during a run. Exploiting this idle memory capacity opens an opportunity to leverage the unused GPU memory for retaining data.

Unfortunately, current Serverless LLM deployments operate under a restrictive assumption—once a model instance is scheduled to a GPU, it exclusively occupies the GPU memory for its entire service lifetime. During initialization, after the model parameters are loaded, most of the remaining GPU memory is reserved for the KV cache to reduce computation overhead. However, the actually required KV cache size is highly dependent on the input sequence length and inference batch size, which vary significantly across queries and over time. As a result, a large portion of the reserved GPU memory for KV cache remains unused. When a model’s lifecycle ends, all GPU-resident data, including reusable model parameters, are discarded, even

if subsequent requests require the same model. This over-conservative design, inherited from traditional single-model serving architectures where it works well, unfortunately results in severe resource underutilization and eliminates reuse opportunities in the Serverless LLM environment. To fully exploit GPU memory and reduce model-switching latency, *the deployment paradigm must evolve from an exclusive, single-model approach to a shared-memory, multi-model architecture.*

In this work, we present **Tangram**, a Serverless LLM framework that accelerates model loading through GPU memory reuse and affinity. Unlike the traditional exclusive resource-binding approach, Tangram enables parameters from multiple models to share GPU memory. When loading a new model, it selectively evicts resident parameters based on size and access frequency. Frequently accessed (*hot*) parameters are retained. These preserved parameters are reused to serve subsequent requests for the same model, thereby reducing data transfer and model-switching latency.

To minimize model parameter loading overhead, Tangram employs a *tensor-level reuse strategy* along with a *unified memory pool* for memory allocation and reclamation. For efficient memory management, Tangram formulates tensor reuse as the Multi-Choice Multi-Dimensional Knapsack Problem and applies a two-stage heuristic algorithm to guide allocation decisions. Tangram also supports on-demand KV cache allocation by integrating the Unified Memory Pool with the ElasticKV engine and E-Attention kernel, enabling KV blocks to be allocated dynamically during runtime decoding. To exploit GPU affinity, Tangram employs a *GPU-affinity-aware scheduler* that considers reusable parameters on GPUs when making scheduling decisions and selects the GPU with the highest expected loading efficiency.

We have implemented a prototype of Tangram based on SLLM [14] and VLLM inference engine [23]. The implementation is open-sourced [37]. We evaluate Tangram on a Serverless LLM platform using real-world workloads. Experimental results show that Tangram accelerates model loading by up to 6.2× and reduces the overall TTFT by 23%–55% compared to state-of-the-art approaches.

The remainder of this paper is organized as follows. Section 2 introduces the background and motivation. Section 3 and 4 describe the design and implementation. Section 5 presents the performance evaluation. Section 6 reviews the related work. Section 7 concludes the paper.

2 Background and Motivation

2.1 Serverless LLM

Serverless LLMs adopt a distributed architecture, where *central Controller* nodes manages the scheduling of model requests, and *Worker* nodes execute model inference, as illustrated in Figure 3. Controller nodes maintain a *resource map* recording the real-time status of worker nodes,

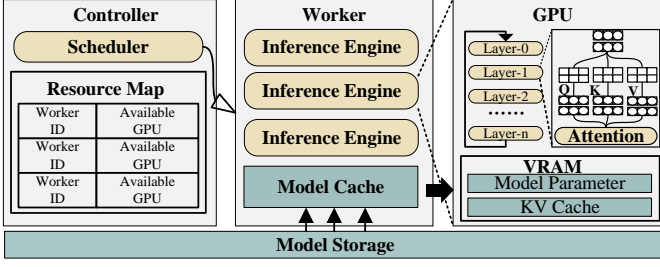


Figure 3. Serverless LLM Architecture.

including available GPUs and the execution state of each model instance.

A model must be registered before handling inference requests. During *registration*, the user uploads the model parameters to the platform, which are then stored in persistent storage and assigned a unique *model ID*. When an inference request for a model arrives, Controller nodes first check whether the model is already loaded in an inference engine. If so, the request is dispatched to that engine and queued for execution. Otherwise, the Controller invokes its scheduling algorithm to locate idle GPU resources. If no resources are available, the request must wait until the running tasks complete and GPU resources are released.

When a model instance is allocated to a worker node, it runs exclusively on GPUs, depending on the size of the model. After completing its inference task, the instance enters a brief idle period to reduce cold-start frequency by keeping the model temporarily active. When the idle period ends, the instance is terminated and its parameters are evicted to release GPU resources.

2.2 LLM Inference

Once a model instance is scheduled to a worker node, it goes through a multi-phase initialization process to set up the inference environment and generate the first token. The initialization process mainly includes four phases:

Init establishes the runtime context for the target model, including Python library initialization, inference engine startup, and model architecture setup. Recent studies leverage CRIU checkpointing to preserve initialized contexts and restore them for subsequent deployments [35, 44], thereby eliminating redundant overhead and reducing the latency of Init.

Load transfers model parameters to GPU memory, reconstructing them into tensors. SLLM [14] optimizes this process by partitioning parameters into manageable chunks and caching them in CPU memory. These chunks are transferred to the GPU in parallel, improving Load efficiency by maximizing PCIe bandwidth utilization.

Profile conducts a preliminary inference run of the target model to determine the optimal memory allocation for various data components, such as the KV cache, activation

data, and other intermediate tensors, which are required during model execution. Medusa [49] improves this process by pre-computing model-specific profiles offline, thereby removing profiling overhead from the critical path and further accelerating startup.

Prefill runs the entire input prompt in parallel to compute the KV pairs and generate the first token. Due to current GPU’s lazy kernel loading approach, *Prefill* typically includes kernel launch overhead. Tidal [13] proposes using an adaptive kernel template to reduce the launch overhead.

After *Prefill*, the model enters the *Decode* stage, during which it generates tokens sequentially using all previously computed KV pairs. In each decoding step, the input passes through the stack of Transformer layers, where the attention mechanism computes intermediate representations until the final layer projects them into logits for the next token [6, 41]. At the same time, KV pairs are also generated and stored for each new token.

During inference, a model maintains a KV cache to avoid recomputing KV pairs for previously generated tokens. However, storing this cache incurs substantial memory overhead, as its size depends on both *sequence length* and *batch size*. Sequence length is the total number of tokens generated during Prefill and Decode for a single inference request; the batch size is the number of inference requests processed simultaneously. Both factors vary dynamically, making it challenging to predict memory requirements. To ensure stable performance, existing systems conservatively pre-allocate GPU memory for the maximum sequence length for each request.

2.3 Performance Bottleneck and Opportunities

To investigate cold-start bottlenecks, we have developed a prototype of Serverless LLM that integrates three state-of-the-art optimizations: SLLM parallel loading [14], CRIU-based checkpointing [12, 35, 44], and Medusa’s offline profiling [49], which we refer to as *SLLM-CM*.

We measure the TTFT breakdown of SLLM-CM across different model sizes. As shown in Figure 2, the integrated optimizations significantly reduce overall TTFT, but *Load* latency remains dominant and grows linearly with model size, revealing parameter loading as the primary bottleneck in cold starts. This paper explores opportunities for reducing Load latency through GPU memory reuse and GPU affinity, guided by two key observations:

Observation 1: Model accesses exhibit locality. We find that in Serverless LLM platforms, the incoming inference tasks exhibit temporal locality. To investigate this, we use a real-world serverless trace [31] to generate a sequence of access requests across eight different models. Figure 4a illustrates the access interval of each model, where the interval is defined as the number of intervening models between two consecutive requests to the same model. The results show that most requests are consecutive, meaning

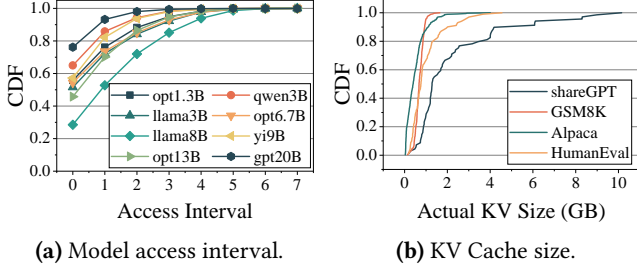


Figure 4. (a) Model access intervals in Serverless LLM trace. (b) Measured KV cache size during inference (Batch Size=16).

that if the model parameters remain in GPU memory, subsequent requests for the same model could avoid redundant data transfers.

Traditional serverless deployments exploit this locality by configuring the `keep_alive` parameter to retain active model instances for a specified period, e.g., 4 minutes [7, 15, 21]. However, the request interval for the same model is highly variable in real-world scenarios, making it challenging to select a static, pre-configured value that minimizes cold starts and prevents unnecessary GPU occupancy.

Observation 2: Current KV cache allocation is overly conservative. LLM inference engines typically pre-allocate the *maximum* GPU memory for the KV cache, without considering the dynamic variability in sequence lengths observed in real-world inference. As shown in Figure 4b, empirical measurements on the Llama8B model (batch size 16) across four representative datasets—ShareGPT [33], GSM8K [26], Alpaca [38], and HumanEval [27]—reveal that actual KV cache consumption can vary by a factor of 25–137 across datasets. While the KV cache may occupy nearly all available memory in worst-case scenarios, it usually consumes far less than the maximum allocation.

Based on the above observations, we have identified two important opportunities for optimizations: (1) Leveraging the temporal locality of inference tasks, retaining parameters in GPU memory even *after* the model instance is released can avoid the latency of reloading data in the next run. However, simply caching the last model’s parameters is insufficient, as access intervals vary significantly across models (Figure 4a), particularly at interval 0. If only retaining the last used model’s parameters in GPU memory, models that have a relatively long reuse interval (e.g., llama8B in this example) would still be frequently loaded and evicted. (2) Allocating KV cache based on *actual demand* frees unnecessarily reserved GPU memory, which can then be used to retain model parameters. This allows parameters of inactive model instances to be preserved in GPU memory.

3 Design of Tangram

In this section, we present the design details of Tangram. Our design aims to achieve three important goals.

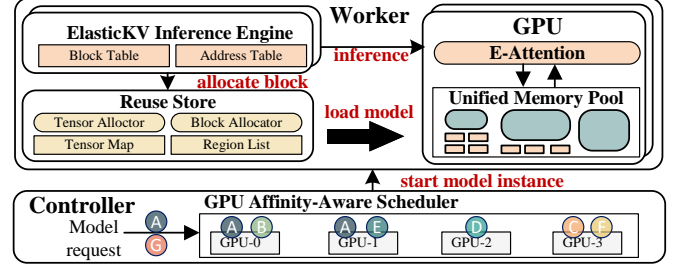


Figure 5. Tangram Architecture.

Aim #1: Minimize model parameter loading. To accelerate model loading, it is crucial to eliminate redundant parameter transfers by exploring model temporal locality. Tangram seeks to minimize parameter loading by (i) selecting an appropriate reuse granularity, and (ii) identifying and utilizing the reusable parameters during model loading.

Aim #2: Enable efficient memory management. To maximize reuse efficiency, Tangram aims to incorporate an efficient memory management mechanism that (i) prioritizes retaining more high-value model parameters, (ii) mitigates memory fragmentation caused by frequent model switching, and (iii) supports on-demand KV cache allocation.

Aim #3: Exploit GPU affinity. With memory reuse enabled, GPUs differ in their affinity for specific models based on the parameters already resident in memory. Tangram requires a GPU affinity-aware scheduler that assigns requests to selected GPUs to maximize the benefits of memory reuse and minimize end-to-end latency.

Figure 5 illustrates the architecture of Tangram. Tangram employs a *Tensor-level Model Reuse and Loading* mechanism which retains the model tensors in reused GPU memory, and overrides the Load operation to only transfer the necessary tensors (Section 3.1). To enable efficient memory management, Tangram integrates the *Reusable Memory Management* (Section 3.2) to solve memory shortage problems at minimal cost within limited time, and leverages the *On-Demand KV cache Allocation* (Section 3.3) to further increase reusable GPU memory. Finally, Tangram also designs a *GPU Affinity-Aware Scheduler* (Section 3.4) to allocate models to their favored GPUs. In the following sections, we introduce the details of each component.

3.1 Tensor-level Model Reuse and Loading

In traditional Serverless LLMs, model parameters are loaded and released at the *model* level, dictated by the exclusive use of the GPU by each model. However, adopting model as the basic unit leads to overly coarse-grained memory management, causing memory waste and poor utilization in Serverless LLM scenarios. A naïve alternative is to partition GPU memory into equally sized *pages* and map parameters to them, allowing page-level loading and eviction during model switches. However, with a large number of

concurrent hardware threads, GPU execution depends on coalesced accesses to global memory. Mapping parameters to non-contiguous physical pages breaks coalescing, wastes bandwidth, and could degrade performance by up to an order of magnitude [25]. In addition, due to the overly small granularity, managing a large number of pages exacerbates the management cost.

We propose to realize parameter reuse at the *tensor* level. Tensors are the fundamental units for organizing and scheduling computation in LLMs. Reusing data in units of tensors can ensure that each tensor’s data remains contiguous, which helps preserve computational efficiency and avoid extensive modifications to existing GPU kernels. As a model typically comprises dozens of tensors, this approach is more fine-grained than model-level reuse, reaching a balance between fine-grained memory utilization and acceptable management complexity.

As illustrated in Figure 5, Tangram deploys a *Reuse Store* on each worker node to manage the reusable memory of each GPU (the Unified Memory Pool) and to support tensor-level model loading. For each GPU, the Reuse Store generates a unique fingerprint as the index for every tensor. A hash table-based *Tensor Mapping* is maintained to keep track of each tensor’s index (fingerprint) and its corresponding address in GPU memory.

During a model switch, once a model is assigned to a specific GPU, the Reuse Store first searches in the Tensor Mapping table for the model’s tensors. For tensors not found, memory is allocated through the *Tensor Allocator* in the Unified Memory Pool and the tensors are loaded from either the CPU-side Model Cache or the persistent Model Store. If the GPU does not have sufficient free memory space, the system evicts selected tensors from other “inactive but resident” models to free space. After all required tensors are loaded, the Reuse Store updates the tensor mapping and returns the actual GPU memory addresses of all tensors to the Inference Engine for computation.

3.2 Reusable Memory Management

To enable tensor-level model reuse, the Reuse Store logically divides the GPU memory space into a sequence of *regions*. Each region is a contiguous memory space in the Unified Memory Pool; and the regions are chained together following their physical layout. A region is either *free* or *allocated*.

During the model loading process, if a tensor is absent from the Tensor Map, the Tensor Allocator needs to select a free region to allocate memory for it. A commonly used strategy is the *best-fit* approach, which assigns the smallest free region that can accommodate the tensor, and if necessary, evicts existing tensors using an LRU policy. Each time when a tensor is allocated from a free region, a new, smaller free region could also be produced (for the remaining unused space). Many small, non-contiguous free regions are scattered in the memory space, none of which is large enough

to hold a tensor’s contiguous memory requirements, even though the total free space is sufficient. In such cases, model loading could fail due to memory fragmentation.

The fragmented memory space must be reorganized, but this incurs non-trivial memory-copy overhead [28, 34, 42]. Tangram addresses this by modeling tensor reuse allocation as a Multi-Choice Multi-Dimensional Knapsack Problem (MCMCKP) [16, 20] and applying a two-stage heuristic loading strategy to minimize overhead.

3.2.1 Multi-Choice Multi-Dimension Knapsack Problem. When allocating space for new tensors, we can either (1) *evict* existing tensors to free regions or (2) *move* a tensor to merge adjacent free regions into a larger one. Eviction may cause additional latency if the evicted tensors are later accessed, while movement incurs immediate merge overhead, increasing loading latency. We formalize this process as a *Multi-Choice Multi-Dimensional Knapsack Problem* (MCMCKP), formulated in Equation 1.

$$\begin{aligned}
 &\textbf{Given: } T = \{t_1, \dots, t_N\}, T' = \{t'_1, \dots, t'_M\}, S, \\
 &\quad s_i, s'_j, c_j, m_j. \\
 &\textbf{Variables: } x_i, y_j, z_j \in \{0, 1\}. \\
 &\textbf{Constraints: } \sum_{i=1}^N s_i x_i \leq S - \sum_{j=1}^M s'_j (1 - y_j), \\
 &\quad y_j + z_j \leq 1, \quad \forall j, \\
 &\quad \sum_{i=1}^N x_i = N. \\
 &\textbf{Objective: } \min \sum_{j=1}^M (c_j y_j + m_j z_j)
 \end{aligned} \tag{1}$$

In the equation above, $T = \{t_1, \dots, t_N\}$ denotes the set of new tensors, where the total number of tensors is N , and each tensor t_i of which needs to be allocated with space of size s_i ; $T' = \{t'_1, \dots, t'_M\}$ denotes the set of existing tensors (total number is M) currently residing in GPU memory, each tensor t'_j of which has size s'_j , eviction cost c_j , and merge cost m_j . The total size of the Unified Memory Pool is denoted by S . We define three binary decision variables: $x_i = 1$ if a new tensor t_i is successfully allocated, $y_j = 1$ if an existing tensor t'_j is evicted, and $z_j = 1$ if an existing tensor t'_j is merged.

The constraints in the equation guarantees that: (1) the memory consumed by new tensors does not exceed the available GPU memory after counting in evicted tensors; (2) an existing tensor cannot be both evicted and merged simultaneously; (3) all new tensors must be allocated with space. The objective is to minimize the overall allocation cost, defined as the weighted sum of eviction and merge costs across all existing tensors.

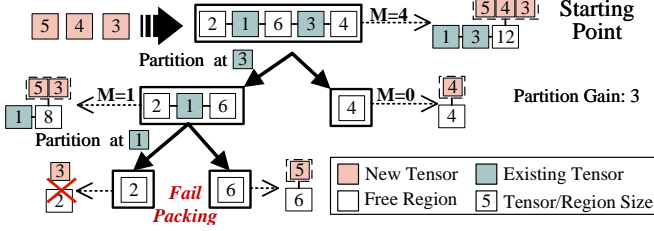


Figure 6. Partitioned-Gain Packing algorithm.

3.2.2 Two-Stage Loading Heuristic. The MCMDKP is NP-hard, as it can be reduced from the classical 0-1 knapsack problem. Tangram solves it with a two-stage heuristic algorithm that decomposes the multi-dimensional, multi-choice problem into stepwise subproblems.

Stage 1: Minimal-Cost Eviction. The goal of the first stage is to identify a set of tensors for eviction such that the total size of free regions is sufficient for the new tensors. Tangram defines the eviction cost c_j for tensor t'_j as:

$$c_j = p_m \cdot \left(\frac{s'_j}{b_m} \right) \cdot \alpha_m, \quad (2)$$

where model m is the model to which tensor t'_j belongs. Here, p_m denotes the tensor miss probability from model m , s'_j the size of tensor t'_j , b_m the model loading bandwidth, and α_m the model's loading-latency sensitivity. The value of α_m is specified by user during registration, since different models may have distinct responsiveness requirements [8, 22]; It is set to 1 by default; for less latency-sensitive models, α_m can be set to a smaller value.

Tangram adopts a greedy strategy, which evicts tensors based on ascending eviction cost until the released size meets the requirement for loading the new tensors. In other words, we reclaim memory with the lowest-cost tensor.

Stage 2: Partitioned-Gain Packing. Due to fragmentation, the total free space may be sufficient yet still unable to accommodate all new tensors. A simple solution is to merge all free regions together to one end, forming a single range of contiguous free space, but this incurs significant merge cost. Tangram implements a *Partitioned-Gain Packing* algorithm to solve this problem [9]. The objective is to fit new tensors directly in available free regions and perform merging *only* when necessary.

The algorithm starts with the worst-case “merge-all” scenario, then explores feasible assignment schemes, evaluating their reductions in merge cost to identify the most cost-effective option. Thus, the problem of minimizing merge cost is reformulated as the problem of maximizing the assignment gain. Since enumerating all allocation schemes is computationally intractable, we adopt a partitioning and bin-packing strategy to obtain a near-optimal solution.

Partitioning subspaces. A *subspace* (P) is defined as a consecutive set of regions that begins and ends with a

Algorithm 1: Partitioned-Gain Packing Algorithm

Input: New tensors T (sorted by size in descending order), Region List R

Output: Allocation set \mathcal{A} with subspaces and tensor subsets, and overall merge cost \mathcal{M}

```

1  $\mathcal{P} \leftarrow \{(R, T)\}; \mathcal{A} \leftarrow \emptyset;$ 
    $\mathcal{M} \leftarrow$  size of all allocated regions in  $R$ ;
2 while  $\mathcal{P} \neq \emptyset$  do
3   foreach  $(P, \mathcal{T})$  in  $\mathcal{P}$  do
4      $split\_occurred \leftarrow \text{false};$ 
5     foreach  $t'_p \in P$  do
6       Split  $P$  at  $t'_p$  into  $(P_1, P_2)$ ;
7        $\mathcal{T}_1, \mathcal{T}_2 \leftarrow \text{TryPacking}(\mathcal{T}, P_1, P_2)$ ;
8       if  $(\mathcal{T}_1, \mathcal{T}_2) \neq (\emptyset, \emptyset)$  then
9          $split\_occurred \leftarrow \text{true};$ 
10        Replace  $(P, \mathcal{T})$  with  $(P_1, \mathcal{T}_1)$  and
            $(P_2, \mathcal{T}_2)$  in  $\mathcal{P}$ ;
11         $\mathcal{M} \leftarrow \mathcal{M} - t'_p.size;$ 
12        break;
13   if  $split\_occurred = \text{false}$  then
14     Remove  $(P, \mathcal{T})$  from  $\mathcal{P}$ ;
15     Add  $(P, \mathcal{T})$  to  $\mathcal{A}$ ;
16 return  $\mathcal{A}, \mathcal{M};$ 
17 Function  $\text{TryPacking}(\mathcal{T}, P_1, P_2)$ 
18    $\mathcal{T}_1, \mathcal{T}_2 \leftarrow \emptyset;$ 
19    $C_1 \leftarrow P_1.capacity; C_2 \leftarrow P_2.capacity;$ 
20   foreach  $tensor\ t' \in \mathcal{T}$  do
21     if  $t'.size \geq \min(C_1, C_2)$  then
22       return  $\emptyset, \emptyset;$ 
23     if  $C_1 \geq C_2$  then
24       Add  $t'$  to  $\mathcal{T}_1$ ;  $C_1 \leftarrow C_1 - t'.size;$ 
25     else
26       Add  $t'$  to  $\mathcal{T}_2$ ;  $C_2 \leftarrow C_2 - t'.size;$ 
27 return  $\mathcal{T}_1, \mathcal{T}_2;$ 

```

free region, characterized by its *capacity* C (the total size of free and allocated regions) and *maximum merge cost* M (the total size of allocated regions). For example, in Figure 6, the initial subspace has $M = 4$ and $C = 12$. A subspace can accommodate a set of new tensors (\mathcal{T}), if its capacity is at least their total size of \mathcal{T} , meaning that all tensors can be placed with no more than cost M . Our objective is to partition the subspace so that all new tensors are packed into free regions at the lowest possible merge cost.

Based on this idea, we employ Partitioned-Gain Packing (Algorithm 1) to recursively search for an allocation plan with minimal merge cost. The algorithm maintains a pending set \mathcal{P} of subspace–tensor pairs and an allocation set \mathcal{A} of

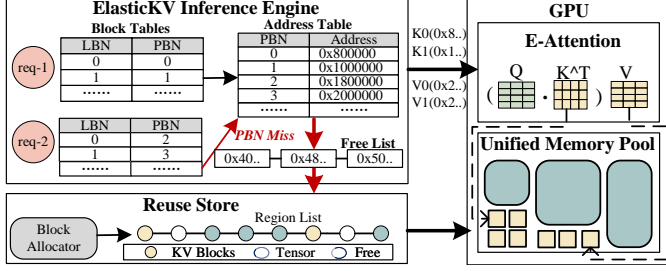


Figure 7. On-demand KV cache Allocation.

finalized allocations (line 1). For each pair (P, \mathcal{T}) in \mathcal{P} (line 3), it traverses the allocated regions t'_p in P as potential *partition points* (line 5), each splitting P into P_1 and P_2 and yielding a *partition gain* of $t'_p.size$ (line 7), since the tensor no longer contributes to the merge cost, thereby reducing \mathcal{M} by $t'_p.size$ (line 11). At each candidate point, the TRYPACKING function attempts to divide \mathcal{T} into \mathcal{T}_1 and \mathcal{T}_2 that fit into P_1 and P_2 . Successful partitioning generates two new subspace–tensor pairs added to \mathcal{P} for further recursion (line 10), while subspaces that cannot be partitioned are moved to \mathcal{A} (line 12). The process terminates when \mathcal{P} becomes empty.

Packing tensors in subspace. Determining whether a partitioning attempt succeeds can be viewed as a bin-packing problem with two bins. Tangram employs a variant of the *Best Fit Decreasing* (BFD) policy [39], which places the descending-sorted tensors into the subspace with the largest remaining space one by one (line-17 to line 27 in Algorithm 1). If any tensor cannot fit into the subspace, the process terminates immediately and returns failure, meaning that the partitioning attempt at this partition point is unsuccessful. Note that the tensors are only sorted once, and the packing process maintains this order across all generated subsets. Given the small number of tensors, the sorting cost is negligible.

Figure 6 illustrates an example, in which the initial subspace has two candidate partition points with partition gains of 1 and 3. It first selects the point with gain 3; if succeeds, the subspace splits into two with tensors divided accordingly, and the total merge cost $\mathcal{M} = 4$ is reduced to 1. However, the left subspace fails to pack its subset, and the right subspace cannot be further partitioned. The algorithm therefore terminates with a final merge cost of 1.

3.3 On-demand KV cache Allocation

In traditional Serverless LLM architectures, once the model is loaded, all remaining GPU memory is reserved for the KV cache. Although it guarantees inference performance in the worst-case scenario, such a conservative approach leads to significant memory waste due to workload variability, severely restricting the amount of memory available for Tangram to retain tensors for reuse.

Instead of pre-allocating KV cache space at initialization, Tangram utilizes the ElasticKV Inference Engine to allocate memory dynamically during inference. Since the KV size per token is fixed, the total cache requirement depends only on the current token count, which is determined by the input length in *Prefill* and incremented by one token per step in *Decode*. Before each phase, the ElasticKV Engine estimates the token count and allocates the exact required space.

In practice, multiple requests are often processed simultaneously [1, 45, 50]. Therefore, similar to conventional inference engines (e.g., vLLM [23]), the ElasticKV Inference Engine maintains a KV *Block Table* for each request (as shown in the Figure 7). Each block stores the KV pairs of several consecutive tokens (e.g., 16). The Block Table records the mapping from a request’s *Logical Block Numbers* (LBNs) to the globally unique *Physical Block Numbers* (PBNs). Each PBN corresponds to a physical address in the Unified Memory Pool, which is tracked in the *Address Table*. When a block becomes full, the ElasticKV Engine allocates a new physical block with the mapping updated in the Block Table.

The ElasticKV Inference Engine requests physical block space from the Unified Memory Pool by invoking the *Block Allocator* in the Reuse Store, which allocates and reclaims blocks using the same Region List as the Tensor Allocator. However, in urgent situations with no free space, the Block Allocator directly reclaims tensors from inactive models based on the cost-aware eviction policy in Section 3.2.2 to quickly reclaim memory and avoid impacting decode performance. The Block Allocator marks regions assigned to KV blocks in the Region List and reclaims them collectively once the current model instance completes.

Before each inference begins, the ElasticKV Engine provides the Attention Kernel with all KV blocks of the current request along with their physical addresses for computation. Tangram implements an *E-Attention Kernel* based on PagedAttention [23], which executes attention at block granularity and retrieves KV pairs dynamically from the Unified Memory Pool using physical addresses provided by the ElasticKV Engine.

Optimizations. On-demand KV cache allocation improves memory utilization but introduces runtime overhead. To mitigate this, completed request blocks are retained in a *Free List* rather than immediately returned to the Reuse Store, and new blocks are allocated from the Reuse Store only when the Free List is empty. Since multiple batched inference tasks often run concurrently, the ElasticKV Engine merges their block-allocation requests before dispatch, amortizing overhead through batch operations. In addition, using a larger block size further reduces allocation frequency; this granularity is defined at system initialization and applied consistently by ElasticKV and the Block Allocator for allocation and reclamation.

3.4 GPU Affinity-Aware Scheduler

In Serverless LLM, the *request scheduler* determines the GPU for each request, directly affecting cold-start latency. Conventional schedulers, e.g., SLLM [14], mainly consider CPU-side resources (Model Cache), leaving GPU-side resources overlooked. Tangram, however, retains inactive model tensors in GPU memory for reuse, so loading time depends on the amount of reusable tensors, making GPU memory state an essential factor in scheduling.

Tangram adopts *GPU affinity-aware scheduling* (Algorithm 2), which receives a sequence of model requests (each with a model ID) and produces a mapping of models to scheduled GPUs. It first enumerates available GPUs (line 1) and then traverses requests, selecting the best GPU for each model (lines 2–12). For a given request, it first checks GPU feasibility (e.g., sufficient available memory, lines 5–6), measures the size of reusable tensors (line 7), and estimates the loading time (line 8). The GPU with the lowest expected latency is chosen (lines 9–10). If no idle GPU can accommodate the model, the request remains in the queue for future scheduling; otherwise, the model is paired with its best candidate GPU, which is subsequently removed from the available pool (lines 11–12). By prioritizing GPUs with a higher proportion of reusable tensors, Tangram effectively reduces loading time by exploiting memory locality.

Loading Time Estimation. Tangram estimates the expected loading time for a model-GPU pair, which serves as the foundation for scheduling decisions:

$$t_{load} = \frac{S - S'}{B} \quad (3)$$

where S is the model size, S' is the size of reusable tensors on the GPU, and B is the measured model-loading bandwidth. The bandwidth depends on the model’s storage location: if it is in the Worker’s Model Store, the model must first be copied to the Model Cache before GPU loading. Tangram leverages the asynchronous loading mechanism of SLLM [14], overlapping the transfer from the Model Store to the Model Cache with the transfer from the Model Cache to GPU memory. Consequently, t_{load} is determined only by the slower transfer medium.

4 Implementation

We implemented a fully functional Tangram prototype based on SLLM [14] and VLLM [23] with 11,829 lines of CPU code and 725 lines of CUDA code, fully available as open source [37]. Specifically, Tangram allocates the Unified Memory Pool with `cuda_malloc`, transfers tensors via `cuda_mem_copy`, and shares GPU memory handles with ElasticKV through the CUDA IPC API. To improve copy efficiency, adjacent tensors are grouped until exceeding an empirically chosen 8 MB threshold. ElasticKV, built on VLLM, restores tensors with `Torch.from_blob` and passes KV block addresses to E-Attention, extending PagedAttention

Algorithm 2: GPU Affinity-Aware Scheduling

```

Input: Requests; // list of model IDs
Output: Schedules; // list of (model_id, gpu)
1 GPUs ← AvailGPUs(); Schedules ← [];
2 for  $m \in \text{Requests}$  do
3    $best\_gpu \leftarrow \emptyset$ ;  $best\_latency \leftarrow +\infty$ ;
4   for  $g \in \text{GPUs}$  do
5     if  $CanRun(m, g) = \text{False}$  then
6       continue;
7      $reuse\_size \leftarrow GetReuseSize(m, g)$ ;
8      $lat \leftarrow EstimateLoadTime(m, reuse\_size, g)$ ;
9     if  $lat < best\_latency$  then
10       $best\_gpu \leftarrow g$ ;  $best\_latency \leftarrow lat$ ;
11  if  $best\_gpu \neq \emptyset$  then
12    Schedules.push( $m, best\_gpu$ );
    GPUs.erase( $best\_gpu$ );
13 return Schedules;

```

with two CUDA kernels (`reshape_and_cache_segment` and `segmented_attention`) to enable physical-address-based KV cache access. Finally, the GPU Affinity-Aware Scheduler extends SLLM’s resource-aware scheduler by querying each worker’s Reuse Store for real-time reusable tensor sizes and schedules requests based on the estimated loading time.

5 Evaluation

5.1 Experimental Setup

Experimental Environment. We evaluate Tangram’s performance with two setups. We use a single server equipped with NVIDIA L40 GPU (45 GB VRAM), a 224-core Intel Xeon Platinum 8480+ CPU, 1 TB DDR4 memory, and 7 TB SSD storage for single-node studies. For scalability studies, we use a multi-GPU server with eight NVIDIA 4090 GPUs (24 GB VRAM each) in a Ray-based [30] distributed setup consisting of one controller node and eight worker nodes (one GPU per worker). The controller node has a 16-core Intel Xeon Platinum 8352V CPU and 16 GB DDR4 memory, while each worker node has one NVIDIA 4090 GPU and 64 GB DDR4 memory. All experiments are conducted with Python v3.10, CUDA v12.3, PyTorch v2.5.0, and Ray v2.10.0.

LLM Models. Following the model size distribution reported by multi-model service platforms [19] and prior work [48], we select popular models from the OPT [51], LLaMA [40], Qwen [43], Yi [46], and GPT [29] series. Among them, 30% have 1–3B parameters, 60% have 4–13B, and 10% have 14–30B.

Workloads and Datasets. We adopt a similar approach as SLLM [14] to generate workloads, combining the open-source serverless trace from Azura [31, 32] with the model

Table 1. Decode throughput (token/s) of different approaches.

| Model | SLLM | SLLM-C | SLLM-CM | Tangram |
|---------|--------|--------|---------|---------|
| gpt20B | 223.29 | 227.83 | 225.99 | 223.06 |
| opt13B | 303.03 | 299.54 | 297.33 | 293.22 |
| yi9B | 463.81 | 466.90 | 463.04 | 456.19 |
| llama8B | 509.02 | 499.07 | 507.82 | 497.84 |
| opt6.7B | 481.33 | 481.06 | 482.02 | 472.12 |
| llama3B | 620.24 | 618.30 | 620.40 | 611.80 |
| qwen3B | 552.30 | 554.35 | 551.43 | 538.11 |
| opt1.3B | 626.48 | 630.85 | 627.12 | 610.34 |

list to produce corresponding model functions [24]. Request arrivals follow a Gamma distribution. The original trace contains many consecutive model requests; to evaluate Tangram under varying locality conditions, we adjust the *Coefficient of Variation* (CV) of the Gamma distribution and reduce the proportion of consecutive requests when $CV < 1$ to create more challenging cases with lower locality. The rule is as follows: **L1-Locality**: $CV = 0.25$, no consecutive requests; **L2-Locality**: $CV = 0.5$, halved consecutive request length; **L3-Locality**: $CV = 1$, original consecutive request; **L4-Locality**: $CV = 2$, consecutive request. For each model request, we use real datasets as inputs: ShareGPT [33] (chat), GSM8K [26] (math), Alpaca [38] (daily conversations), and HumanEval [27] (programming).

Comparison Baselines. We compare Tangram with the latest Serverless LLM optimizations as follows. The baseline is **SLLM** [14] that caches model parameters in CPU memory and improves loading efficiency via chunking and parallel transfer. **SLLM-C** extends SLLM with CRIU checkpointing [12] to reduce the latency of Init. It creates and stores context-ready model images locally to resume initialization directly from these images. **SLLM-CM** is built on SLLM-C with further enhancement using Medusa’s offline profiling [49] to eliminate the Profile latency.

To isolate GPU loading performance, we assume CPU memory is sufficient, so all missing tensors are served from CPU memory in SLLM variants and Tangram.

5.2 Overall Performance

We evaluate TTFT across different approaches and models. Figure 8 shows the TTFT breakdown. With CRIU and offline profiling, SLLM-C and SLLM-CM significantly reduce the Init and Profile latency. As the model size increases, however, Load becomes the dominant bottleneck in SLLM-CM, accounting for up to 72% of TTFT. Tangram mitigates this by reusing model parameters to avoid redundant data transfers, achieving a $1.8\times$ – $6.2\times$ faster loading and reducing overall TTFT by 14%–60%.

Table 1 reports the decoding throughput of different approaches. SLLM, SLLM-C, and SLLM-CM achieve similar throughput since they share the same inference engine. Tangram’s on-demand KV cache allocation adds only minor overhead, with throughput loss under 3.2%, which is negligible in practice.

5.3 Performance Breakdown

To analyze performance gains, Figure 9 shows Tangram’s load performance breakdown across batch sizes. “**+Reuse**” denotes SLLM with Tangram’s Reuse Store. At batch size 1, “+Reuse” achieves a $2.3\times$ – $7.6\times$ speedup across models. However, pre-allocating KV cache space reduces available GPU memory. In the worst case, the Memory Pool cannot fully accommodate the largest model, GPT-20B, forcing some tensors to be repeatedly loaded from CPU memory.

“**+ODKV**” augments “+Reuse” with Tangram’s on-demand KV cache mechanism, which allocates KV cache space from the Unified Memory Pool according to its actual needs, improving memory utilization. As batch size grows, “+ODKV” sustains strong performance, achieving $1.16\times$ – $2.6\times$ speedups over “+Reuse” and $1.9\times$ – $4\times$ over SLLM. These results highlight the importance of combining both optimizations to achieve the desired performance gains.

5.4 Tangram Performance Analysis

5.4.1 Tensor Allocation Policies. We evaluate how different tensor allocation strategies affect Tangram’s loading performance. “**Rand+GM**” evicts tensors randomly when space is insufficient and merges free space by moving all the existing tensors (*GlobalMerge*). “**MCE+GM**” applies Minimal-Cost Eviction in Section 3.2.2 but still uses *GlobalMerge*. “**MCE+PGP**” represents the Tangram strategy with Minimal-Cost Eviction and Partitioned-Greedy-Packing. For each strategy, we measure the latency overhead of *Compute*, *Merge*, and *Load*.

Figure 10 presents results for two representative models: GPT-20B (large) and OPT-1.3B (small), with other models showing similar trends. Compared with SLLM, *Rand+GM* reduces *Load* time by up to 68% through tensor reuse while Minimal-Cost Eviction further cuts *Load* time by retaining more valuable tensors. However, in both *GM* settings, merge overhead is nonnegligible, particularly for OPT-1.3B, as smaller models leave more existing tensors to be moved during Global Merge. This effect would even be amplified as Memory Pool becomes larger. Partitioned-Greedy-Packing alleviates this by minimizing unnecessary data movement, achieving a 93% reduction in merge overhead. Across all three Tangram strategies, the additional *Compute* latency remains under 1 ms, making it negligible.

5.4.2 On-demand KV cache Allocation. We evaluate the effect of On-Demand KV cache Allocation (ODKV) on tensor reuse space. Figure 11a shows how available

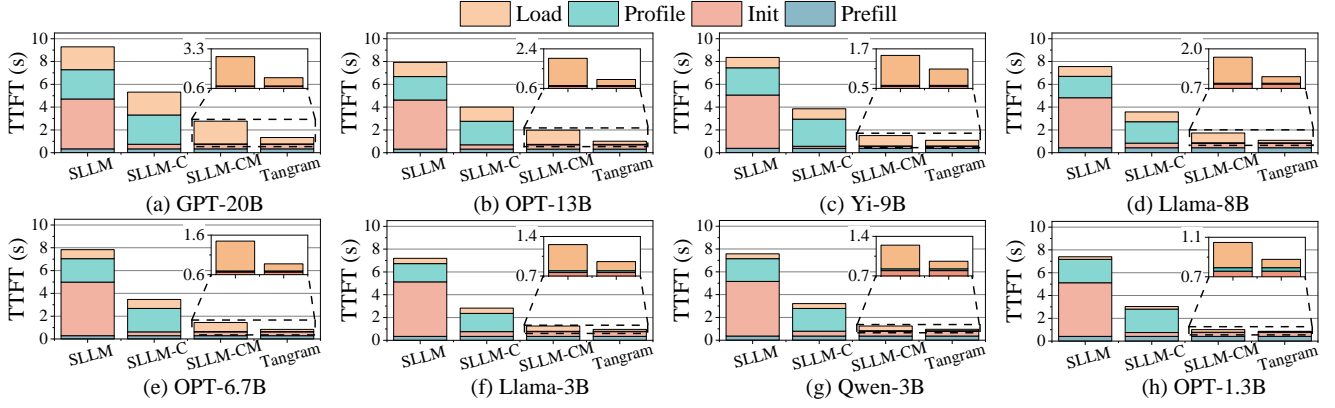


Figure 8. Overall performance: TTFT of different model under different Serverless LLM approaches.

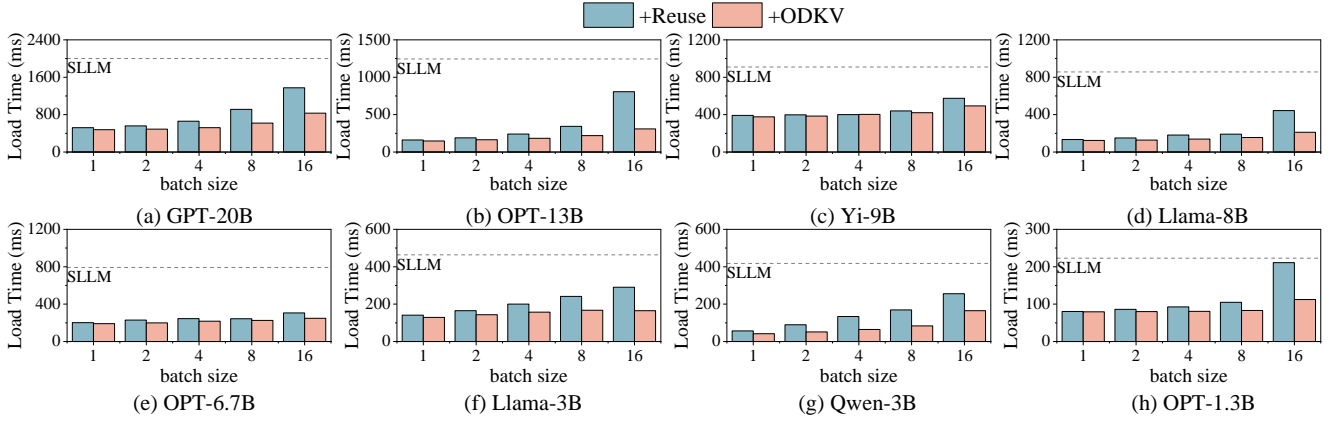


Figure 9. Performance breakdown: TTFT breakdown of Tangram under different batch size.

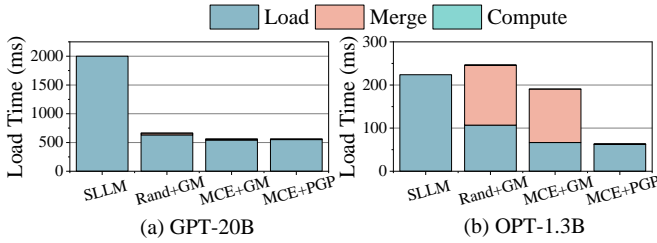


Figure 10. Effects of loading policies: Loading time breakdown for two representative model.

space changes with batch size. In the “w/o ODKV” setting, Tangram reserves KV cache space for each batch using the maximum sequence length. In contrast, “w/ ODKV” allocates space based on actual sequence length, leaving more space for tensors. As batch size grows, ODKV increases reusable space by 2.3%–41%.

Figure 11b reports the overhead of ODKV, where the average allocation overhead is normalized by decode time. Tangram reduces real-time allocation costs through delayed release, a dedicated block allocator, and batched allocation.

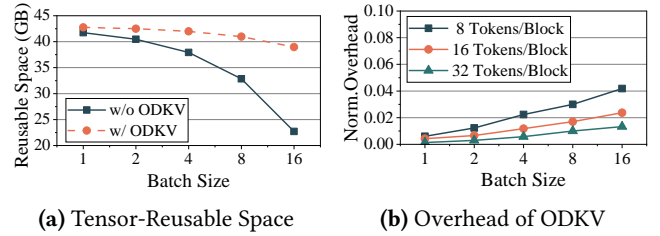


Figure 11. Effects of On-Demand KV cache Allocation: (a) Reusable space size under different batch sizes. (b) Overhead of ODKV under different block sizes.

As a result, ODKV overhead remains below 0.041 for a batch size of 16. Moreover, using coarser-grained block allocation further reduces overhead—increasing block size from 8 to 32 tokens lowers the overhead to 0.013.

5.5 Sensitivity Analysis

We conduct a sensitivity analysis of Tangram’s performance. Figure 12a shows the impact of workload locality and

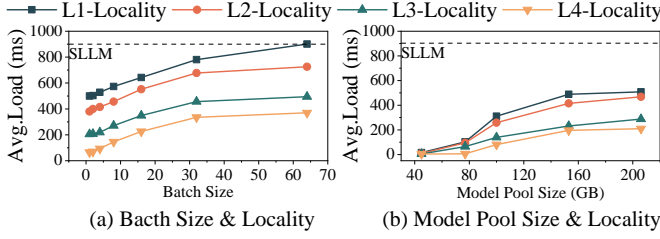


Figure 12. Sensitivity analysis: (a) Effects of workload locality and batch size (Model Pool Size = 120GB); (b) Effects of model pool size (Batch Size = 16).

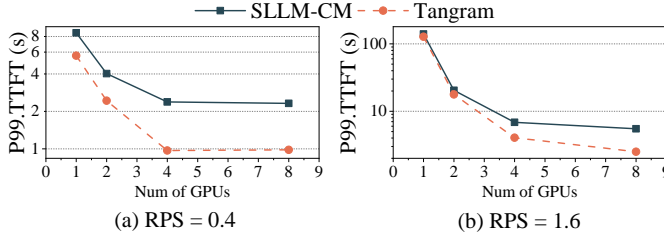


Figure 13. Multi-GPU Performance: 99th percentile tail latency of TTFT across different GPU configurations and request rates (Requests Per Second, RPS).

batch size on loading performance. As locality decreases, Tangram benefits less from directly reusing parameters of the previously loaded model. Similarly, larger batch sizes force the Reuse Store to evict more tensors from the Memory Pool. At a batch size of 64, Tangram can rarely obtain reusable tensors, reaching its worst case when locality is L1 and batch size is 64, where its performance matches the SLLM baseline.

Figure 12b shows the effect of model pool capacity on Tangram’s performance. In this experiment, we used a single GPU and gradually increased the total size of models scheduled to it. Tangram’s load latency initially increases but then stabilizes, and the higher locality is, the lower latency Tangram can achieve stably. Eventually, when the model pool exceeds 200 GB, Tangram’s average load latency remains only 23%–56% of SLLM’s, as it maximizes Memory Pool utilization across varying request lengths to reuse tensors.

5.6 Multi-GPU Performance

We evaluate the multi-GPU performance and scalability of Tangram under real workloads with various numbers of GPUs. Specifically, we control the request arrival rate, measured in requests per second (RPS), and scale the number of workers (each equipped with a single GPU) from 1 to 8. Figure 13 reports the 99th-percentile latency of SLLM-CM and Tangram under two workload intensities.

As shown in the figure, the tail latency of both SLLM-CM and Tangram decreases as the number of GPUs increases, with Tangram exhibiting a more pronounced reduction.

At the low request arrival rate (RPS=0.4), the latency for both methods plateaus after the number of GPUs reaches 4, since additional GPUs provide no further benefit due to low intensity of incoming requests. Nevertheless, under all GPU configurations, Tangram achieves noticeable latency reductions owing to its reuse mechanism, which directly lowers data transfer volume and thus shortens TTFT.

At high request densities (RPS=1.6), limited GPU resources cause Tangram and SLLM-CM to converge in tail latency with more GPUs, as queueing delays dominate. With more GPUs provisioned, Tangram achieves 8%–54% lower latency by leveraging GPU affinity to select the optimal GPU for each model, compared with SLLM-CM’s random GPU selection.

5.7 Overhead Analysis

We analyzed Tangram’s overhead on both Worker and Controller nodes. On Worker nodes, the primary storage cost arises from the CPU-side model cache, identical to SLLM, while the additional metadata overhead includes the *Tensor Map* (52 B per tensor) and the *Region List* (48 B per region), with typically fewer than 10^3 entries in total. On Controller nodes, processing a model instance request entails querying Worker nodes for the status of real-time reusable tensors. In our experiments with 8 Worker nodes, the RPC overhead of this query is 16.3 ms, which is negligible compared to the second-level latency of model initialization.

6 Other Related Work

Optimization of Serverless LLM Loading. SLLM [14] leverages locality to cache model parameters in CPU memory and uses block-wise parallel transfer to improve PCIe bandwidth utilization. Prism [48] coordinates multi-GPU transfers for further bandwidth gains, but both remain constrained by finite PCIe bandwidth versus growing model sizes. Sui et al. [36] and Tidal [13] aim to reduce loading latency through pre-loading, but they are limited by model size and stage-dependent constraints. Tangram instead directly reduces the transfer volume during the Load phase via GPU memory reuse. Tangram is generally orthogonal to these techniques and can be integrated with them for additional performance gains, e.g., purposefully evicting certain Prefill tensors to create better overlapping opportunities for Tidal.

GPU Memory Sharing. VLLM [23] shares GPU memory across KV caches from different sequences but assumes a single model per system, reserving maximal KV space and wasting memory in multi-model Serverless LLM settings. Choi et al. [11] first proposed spatially share a GPU to improve resource utilization for machine learning applications. Prism [48] focuses on enabling multiple LLMs to run concurrently on a single GPU rather than optimizing cold-start latency. It manages memory at the model granularity, requiring full parameter loads for each

launch and relying on CUDA for space management, which limits fragmentation control. Tangram introduces tensor-level GPU memory reuse, a unified pool for model tensors and KV cache, and fragmentation-aware management, which together enable efficient memory utilization and reduced model loading latency in Serverless LLMs.

7 Conclusion

Tangram is a system designed to accelerate Serverless LLM loading through GPU memory reuse. It leverages unused memory to retain model parameters, unifies tensor and KV cache management in a shared memory pool, employs cost-aware allocation with on-demand KV allocation, and performs GPU-affinity-aware scheduling. Our experiments demonstrate substantial improvements in reducing loading latency and TTFT, while scaling efficiently with multi-GPU resources at negligible overhead, which highlights GPU memory reuse as a promising direction for designing the next-generation Serverless LLM platforms.

References

- [1] Ahsan Ali, Riccardo Pincirol, Feng Yan, and Evgenia Smirni. 2020. Batch: Machine Learning Inference Serving on Serverless Platforms with Adaptive Batching. In *2020 International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*.
- [2] Ahsan Ali, Riccardo Pincirol, Feng Yan, and Evgenia Smirni. 2022. Optimizing Inference Serving on Serverless Platforms. *Proc. VLDB Endow.* 15, 10 (2022), 2071–2084.
- [3] Alibaba. 2025. *Alibaba Platform for AI: Elastic Algorithm Service*. <https://www.aliyun.com/product/bigdata/learn/eas> Accessed: 2025-07-07.
- [4] Amazon. 2025. *Amazon Bedrock: The easiest way to build and scale generative AI applications with foundation models*. https://aws.amazon.com/bedrock/?nc1=h_ls Accessed: 2025-07-07.
- [5] Lixiang Ao, George Porter, and Geoffrey M. Voelker. 2022. FaaSnap: FaaS Made Fast Using Snapshot-based VMs. In *2022 European Conference on Computer Systems (EuroSys)*. ACM, 730–746.
- [6] Muhammad Awais, Muzammal Naseer, Salman Khan, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, and Fahad Shahbaz Khan. 2025. Foundation Models Defining a New Era in Vision: A Survey and Outlook. *IEEE Trans. Pattern Anal. Mach. Intell.* 47, 4 (2025), 2245–2264.
- [7] Azure. 2025. *Configure TCP reset and idle timeout for Azure Load Balancer*. <https://learn.microsoft.com/en-us/azure/load-balancer/load-balancer-tcp-idle-timeout?tabs=tcp-reset-idle-portal> Accessed: 2025-07-07.
- [8] Aakash Bhattacharya and AWS Tian Wen. 2025. *Understanding and Remediating Cold Starts: An AWS Lambda Perspective*. <https://aws.amazon.com/cn/blogs/compute/understanding-and-remediating-cold-starts-an-aws-lambda-perspective/> Accessed: 2025-07-07.
- [9] Ernesto G. Birgin, Rafael D. Lobato, and Reinaldo Morabito. 2010. An Effective Recursive Partitioning Approach for the Packing of Identical Rectangles in a Rectangle. *J. Oper. Res. Soc.* 61, 2 (2010), 306–320.
- [10] James Cadden, Thomas Unger, Yara Awad, Han Dong, Orran Krieger, and Jonathan Appavoo. 2020. SEUSS: Skip Redundant Paths to Make Serverless Fast. In *2020 European Conference on Computer Systems (EuroSys)*. ACM, 32:1–32:15.
- [11] Seungbeom Choi, Sunho Lee, Yeonjae Kim, Jongse Park, Youngjin Kwon, and Jaehyuk Huh. 2022. Serving Heterogeneous Machine Learning Models on Multi-GPU Servers with Spatio-Temporal Sharing. In *2022 USENIX Annual Technical Conference (USENIX ATC)*.
- [12] CRIU. 2025. *CRIU: Checkpoint and Restore in Userspace*. https://criu.org/Main_Page Accessed: 2025-07-14.
- [13] Weihao Cui, Ziyi Xu, Han Zhao, Quan Chen, Zijun Li, Bingsheng He, and Minyi Guo. 2025. Efficient Function-as-a-Service for Large Language Models with TIDAL. *CoRR abs/2503.06421* (2025).
- [14] Yao Fu, Leyang Xue, Yeqi Huang, Andrei-Octavian Brabete, Dmitrii Ustiugov, Yuvraj Patel, and Luo Mai. 2024. ServerlessLLM: Low-Latency Serverless Inference for Large Language Models. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [15] Alexander Fuerst and Prateek Sharma. 2021. FaasCache: Keeping Serverless Computing Alive with Greedy-Dual Caching. In *2021 Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [16] Dylan Gaspar, Yun Lu, Myung Soon Song, and Francis J. Vasko. 2020. Simple Population-based Metaheuristics for the Multiple Demand Multiple-choice Multidimensional Knapsack Problem. *Int. J. Metaheuristics* 7, 4 (2020), 330–351.
- [17] Google. 2025. *Model Garden on Vertex AI: Jumpstart your ML project with a single place to discover, customize, and deploy a wide variety of models from Google and Google partners*. <https://cloud.google.com/model-garden?hl=en> Accessed: 2025-07-07.
- [18] Zicong Hong, Jian Lin, Song Guo, Sifu Luo, Wuhui Chen, Roger Wattenhofer, and Yue Yu. 2024. Optimus: Warming Serverless ML Inference via Inter-Function Model Transformation. In *2024 European Conference on Computer Systems (EuroSys)*.
- [19] Hyperbolic. 2025. *Hyperbolic: The Open Access AI Cloud. Hyperbolic Provides Affordable GPU Access and Inference Services for Those at the Edges of AI*. <https://www.hyperbolic.ai/> Accessed: 2025-07-07.
- [20] Shahrear Iqbal, Md. Faizul Bari, and M. Sohail Rahman. 2010. Solving the Multi-dimensional Multi-choice Knapsack Problem with the Help of Ants. In *Swarm Intelligence - 7th International Conference, ANTS 2010, Brussels, Belgium, September 8-10, 2010. Proceedings*, Vol. 6234. 312–323.
- [21] AWS James Beswick. 2025. *Operating Lambda: Performance optimization – Part 1*. <https://aws.amazon.com/cn/blogs/compute/operating-lambda-performance-optimization-part-1/> Accessed: 2025-07-07.
- [22] Artjom Joosen, Ahmed Hassan, Martin Asenov, Rajkarn Singh, Luke Nicholas Darlow, Jianfeng Wang, Qiwen Deng, and Adam Barker. 2025. Serverless Cold Starts and Where to Find Them. In *2025 European Conference on Computer Systems (EuroSys)*.
- [23] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *29th Symposium on Operating Systems Principles (SOSP)*.
- [24] Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [25] NVIDIA. 2021. *Improving GPU Memory Oversubscription Performance*. https://developer.nvidia.com/blog/improving-gpu-memory-oversubscription-performance/?utm_source=chatgpt.com Accessed: 2025-07-14.
- [26] openai. 2025. *gsm8k*. <https://huggingface.co/datasets/openai/gsm8k> Accessed: 2025-07-14.
- [27] openai. 2025. *openai humaneval*. https://huggingface.co/datasets/openai/openai_humaneval Accessed: 2025-07-14.
- [28] Anuj Pathania, Vanchinathan Venkataramani, Muhammad Shafique, Tulika Mitra, and Jörg Henkel. 2017. Defragmentation of Tasks in

- Many-Core Architecture. *ACM Trans. Archit. Code Optim.* 14, 1 (2017), 2:1–2:21.
- [29] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners. In *OpenAI blog*.
- [30] Ray. 2025. *Ray: A Distributed Framework for AI Workloads*. <https://github.com/ray-project/ray> Accessed: 2025-07-14.
- [31] Mohammad Shahradd, Rodrigo Fonseca, Iñigo Goiri, Gohar Irfan Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini. 2020. Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider. In *2020 USENIX Annual Technical Conference (USENIX ATC)*.
- [32] Mohammad Shahradd, Rodrigo Fonseca, Iñigo Goiri, Gohar Irfan Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini. 2020. Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider. In *2020 USENIX Annual Technical Conference (USENIX ATC)*.
- [33] shibing624. 2025. *sharegpt-gpt4*. https://huggingface.co/datasets/shibing624/sharegpt_gpt4 Accessed: 2025-07-14.
- [34] Matthias Springer and Hidehiko Masuhara. 2019. Massively Parallel GPU Memory Compaction. In *2019 ACM SIGPLAN International Symposium on Memory Management (ISMM)*.
- [35] Radostin Stoyanov, Viktória Spisaková, Jesus Ramos, Steven Gurfinkel, Andrei Vagin, Adrian Reber, Wesley Armour, and Rodrigo Bruno. 2025. CRUgpt: Transparent Checkpointing of GPU-Accelerated Workloads. *CoRR* abs/2502.16631 (2025).
- [36] Yifan Sui, Hanfei Yu, Yitao Hu, Jianxun Li, and Hao Wang. 2024. Pre-Warming is Not Enough: Accelerating Serverless Inference With Opportunistic Pre-Loading. In *2024 ACM Symposium on Cloud Computing (SoCC)*.
- [37] TangramGroup. 2025. *Tangram: A GPU Affinity-Aware Serverless LLM Inference Framework*. <https://anonymous.4open.science/r/ServerlessLLM-56B8> Accessed: 2025-08-20.
- [38] tatsulab. 2025. *alpaca*. <https://huggingface.co/datasets/tatsu-lab/alpaca> Accessed: 2025-07-14.
- [39] Takwa Tlili and Saoussen Krichen. 2023. Best Fit Decreasing Algorithm for Virtual Machine Placement Modeled as a Bin Packing Problem. In *9th International Conference on Control, Decision and Information Technologies, CoDIT 2023, Rome, Italy, July 3-6, 2023*. IEEE, 1261–1266.
- [40] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *CoRR* abs/2307.09288 (2023).
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*.
- [42] Ronald Veldema and Michael Philippsen. 2012. Parallel Memory Defragmentation on a GPU. In *2012 ACM SIGPLAN workshop on Memory Systems Performance and Correctness*.
- [43] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yeqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 Technical Report. *CoRR* abs/2407.10671 (2024).
- [44] Yanning Yang, Dong Du, Haitao Song, and Yubin Xia. 2024. On-demand and Parallel Checkpoint/Restore for GPU Applications. In *2024 ACM Symposium on Cloud Computing (SoCC)*.
- [45] Yanan Yang, Laiping Zhao, Yiming Li, Huanyu Zhang, Jie Li, Mingyang Zhao, Xingzhen Chen, and Keqiu Li. 2022. INFless: A Native Serverless System for Low-latency, High-throughput Inference. In *2022 Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [46] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. Yi: Open Foundation Models by 01.AI. *CoRR* abs/2403.04652 (2024).
- [47] Hanfei Yu, Rohan Basu Roy, Christian Fontenot, Devsh Tiwari, Jian Li, Hong Zhang, Hao Wang, and Seung-Jong Park. 2024. RainbowCake: Mitigating Cold-starts in Serverless with Layer-wise Container Caching and Sharing. In *2024 Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [48] Shan Yu, Jiarong Xing, Yifan Qiao, Mingyuan Ma, Yangmin Li, Yang Wang, Shuo Yang, Zhiqiang Xie, Shiyi Cao, Ke Bao, Ion Stoica, Harry Xu, and Ying Sheng. 2025. Prism: Unleashing GPU Sharing for Cost-Efficient Multi-LLM Serving. *CoRR* abs/2505.04021 (2025).
- [49] Shaoxun Zeng, Minhui Xie, Shiwei Gao, Youmin Chen, and Youyou Lu. 2025. Medusa: Accelerating Serverless LLM Inference with Materialization. In *30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [50] Chengliang Zhang, Minchen Yu, Wei Wang, and Feng Yan. 2019. MARK: Exploiting Cloud Services for Cost-Effective, SLO-Aware Machine Learning Inference Serving. In *2019 USENIX Annual Technical Conference (USENIX ATC)*.
- [51] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: Open Pre-trained Transformer Language Models. *CoRR* abs/2205.01068 (2022).