

NEUSTREAM: Bridging Deep Learning Serving and Stream Processing

Haochen Yuan^{*}
Peking University

Yu Cheng[†]
Peking University, Microsoft Research

Yuanqing Wang^{†*}
Peking University, Microsoft Research

Jilong Xue
Microsoft Research

Ziming Miao
Microsoft Research

Wenhai Xie
Peking University

Lingxiao Ma
Microsoft Research

Zhi Yang[‡]
Peking University

Abstract

Modern Deep Neural Network (DNN) exhibits a pattern where multiple sub-models are executed, guided by control flows such as loops and switch/merge operations. This dynamic nature introduces complexities in batching the requests of such DNNs for efficient execution on GPUs. In this paper, we present NEUSTREAM, a programming model and runtime system for serving deep learning workloads using stream processing. NEUSTREAM decomposes the inference workflow into modules and forms them into a streaming processing system where a request flows through. Based on such abstraction, NEUSTREAM is able to batch requests at fine-grained module granularity. To maximize serving goodput, NEUSTREAM exploits a two-level scheduling approach to decide the best batching requests and resource allocation for each module while satisfying service level objectives (SLOs). Our evaluation of NEUSTREAM on a set of modern DNNs like Large Language Models (LLM) and diffusion models, etc., shows that NEUSTREAM significantly improves goodput compared to state-of-the-art DNN serving systems.

CCS Concepts: • Computing methodologies → Machine learning; • Computer systems organization → Real-time systems.

Keywords: Deep Learning Serving, dynamic Neural Networks, Streaming Processing System

1 Introduction

The modern DNN applications often exhibit a dynamic execution pattern where each input might follow different execution paths guarded by control flows like loops or branches. For example, an LLM generates tokens iteratively through a decoding model, with different inputs requiring varying numbers of decoding steps to produce sentences of different lengths. Similarly, in a diffusion-based image generation model, a random input tensor may need to undergo a varying

number of diffusion steps to generate high-quality images. Their dynamic execution patterns are inherently inefficient when run on GPUs. Typically, to fully utilize a GPU, a DNN model needs to process multiple requests in a batch to reuse the weight data. However, if each input sample follows different execution paths, it becomes difficult to batch the inputs together, necessitating sequential execution. This can easily cause the GPU to become bottlenecked by memory access, thereby reducing overall efficiency.

To address such inefficiency, we observe that the dynamic behavior in these applications is fundamentally different from traditional dynamic neural networks [30]. First, the bodies of these control flows are often much larger in granularity rather than consisting of just a few operators. For example, each decoding step in an LLM involves all the layers of a decoding model. Second, although different inputs follow different execution paths, they often repeatedly go through the same modules at different times. For example, different inputs to a diffusion model need to call the same UNet module multiple times. These patterns motivate us to rethink model serving from the perspective of these repeatedly executed control-flow "bodies": each body model could continuously process inputs and forward its outputs to downstream models based on control flow guidance until the full path is completed. This essentially forms a streaming-based computation flow, where each application is decomposed into modules (i.e., the control flow bodies) and connects these modules in a data-flow manner guided by control flow logic. Such an execution model could effectively address the inefficiency caused by the inability to batch many inputs together. Specifically, each module of the application could batch multiple inputs at a fine-grained granularity for execution and then dispatch their results to downstream modules.

However, executing a DNN model in such a streaming-based manner is non-trivial. First, existing DNN frameworks like PyTorch often represent a DNN model as a single Python application with control flows to manage all the different execution paths. It is challenging to decouple such a model into streaming modules and data-flow connections easily. Second,

^{*}Both authors contributed equally to the paper

[†]The work is done when Yuanqing Wang and Yu Cheng are interns at Microsoft Research.

[‡]Zhi Yang is the corresponding author (yangzhi@pku.edu.cn).

existing GPUs are designed for batch processing, which cannot efficiently support streaming execution, where multiple heterogeneous modules are continuously and concurrently running. Third, DNN serving often has strict service-level objective (SLO) requirements for latency. Decoupling an application into multiple modules and scheduling them independently makes it more challenging to optimize for the best overall goodput, i.e., the throughput of requests that satisfy SLOs, given that each module could have very different trade-offs between latency and batch size (§2).

In this paper, we present NEUSTREAM, a deep learning serving system that aims to maximize GPU utilization by exploiting the ability of GPUs to batch multiple requests while meeting certain latency SLOs. NEUSTREAM leverages the following key designs to address the aforementioned challenges: First, to program a DNN application in a streaming manner, NEUSTREAM introduces first-class *stream* and *stream-module* interfaces to explicitly express a continuous data flow execution. The *stream-module* concept inherits from the existing DNN module concept in frameworks like PyTorch, making it natural to adapt existing models. These interfaces implicitly require converting control flows in the native language into data flow connections, i.e., streams, which allows NEUSTREAM to schedule the requests in each stream to optimize GPU utilization. Second, to execute all the heterogeneous stream-modules on GPUs, NEUSTREAM abstracts the GPU resource into fine-grained computation units called *Streaming Processing Units (SPUs)*, where stream-modules can concurrently execute on different subsets of SPUs. Third, given the explicit stream-modules and fine-grained SPUs, NEUSTREAM decomposes the global SLO of each request into partial SLOs for each individual module. NEUSTREAM then decouples the whole scheduling space into two levels: an intra-module level scheduling decision to determine how many requests to execute in a batch for a specific module, and an inter-module level scheduling decision to determine how many resources to allocate across modules. Through such scheduling decoupling, NEUSTREAM can easily scale out to multiple GPUs, i.e., through inter-module resource allocation, and maximize serving throughput, i.e., through intra-module batch execution.

We implemented NEUSTREAM in 5,000 lines of Python code. NEUSTREAM's executors for diffusion models and LLM models are built on top of PyTorch and vLLM [10], respectively. We evaluated NEUSTREAM for serving mainstream DNN applications, such as LLM-based chat service, image generation, and multi-agent tasks, on various NVIDIA GPUs, including the GeForce RTX 4090, Hopper H100, and RTX A6000. Our evaluation demonstrates that NEUSTREAM achieves significantly higher serving goodput compared to state-of-the-art baseline systems, e.g., outperforming Clockwork [8] by up to 5.26× for serving Diffusion models and vLLM by 1.53×–69.31× for serving OPT models at various request rates, respectively. Looking forward, the rapid advancements in DNN

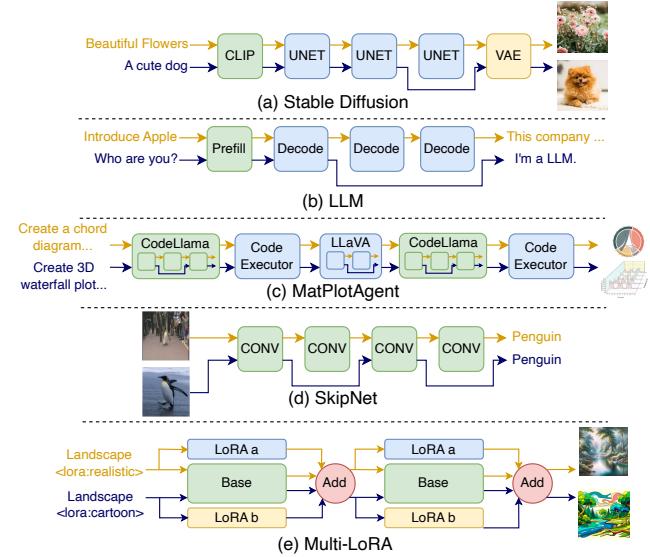


Figure 1. Examples of dynamic neural networks. They dynamically adjust computation paths on per-input basis.

model capabilities could enable more complex execution flows in future applications. We believe NEUSTREAM provides a scalable, efficient, and hardware-friendly system design to accommodate these ever-increasing applications.

2 Motivation

Existing DNN serving systems face the following challenges due to treating the entire model as a monolithic component:

Challenge I: Dynamic execution paths. Most of the existing serving systems consider *static* DNN models whose execution paths are fixed, i.e., independent of the input. *Dynamic* DNN models, in contrast, can adjust their computational path with respect to input, and such models are now becoming prevalent. As illustrated in Figure 1, diffusion [19] models could work at different numbers of denoising steps, depending on preferences on generation quality and inference speed of each input. The LLM [36] generates output tokens autoregressively on a per-input basis, where each request may require a different number of decode iterations. The LLM-based multi-agents [31] further exacerbate this dynamics properties. Skipping and early stopping models such as Skipnet [29] dynamically decide whether to skip certain layers based on the input. Multiple Low-Rank Adaptation (LoRA [9]) requests share the same base model layer but enter into different LoRA adapter layers. Dynamic per-input execution paths in these models pose challenges for the optimizations that have been key to obtaining high GPU efficiency. For example, batching cannot be applied if the execution paths of inputs are not identical. Also, one cannot *proactively* perform scheduling actions (e.g., when or whether to execute a request), because dynamic execution paths have unpredictable latency performance.

166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220

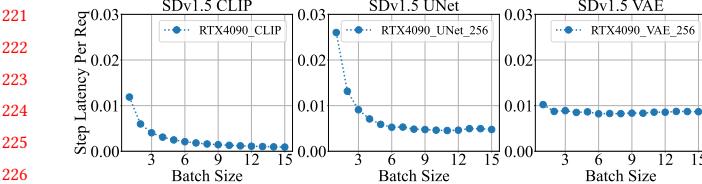


Figure 2. Diffusion phase execution latency on RTX 4090.

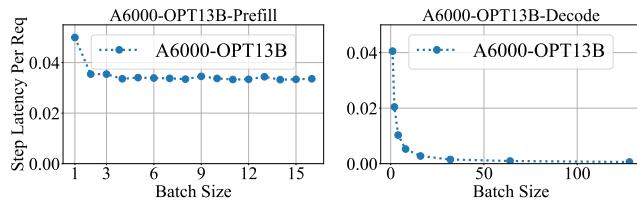


Figure 3. LLM phase execution latency on RTX A6000.

Challenge II: Distinct functional phases. There are different main functional phases during a dynamic model inference. For example, stable diffusion inference involves text encoding, diffusion and decoding phases, and LLM inference contains prefilling and decoding phases. We find that different phases have distinct batching characteristics. Figure 2 and Figure 3 show the request execution latency on each phase, which is batch’s latency / batch size (i.e., excluding queuing time). On a certain phase, as the batch size increases, the decreasing slope of a curve implies that batching leads to only a marginal gain in improving the execution latency (and thus the throughput). Once the gain becomes insignificant, adding more requests to the batch no longer improves the throughput. Instead, it proportionally extends the total processing time for the batch, inadvertently delaying all included requests. Thus, it is necessary to identify a critical threshold for batch size in each phase, beyond which the marginal gain becomes insignificant. Given that the slopes of different phases are distinct, we should identify the threshold specific to each phase. Batching more requests on a phase should only be considered when the batch size of the scheduled request is below its phase-specific threshold.

Insight I: decompose a DNN model into modules. The above challenges motivate us to *decompose* a dynamic DNN into multiple functional modules (e.g., phases), and schedule executions at the granularity of modules. This enables to batch the arrived request on a specific module, irrespective of their dynamic execution paths. For example, diffusion model requests with different numbers of iterations could batch at the granularity of iteration by decomposing the UNet from the computation. Serving multiple LoRA[9] adapters enables batching in the shared base model by decomposing the computation between the base model and the LoRA adapters.

```

1  class StreamModule:
2      def __init__(s_in:List[stream], s_out:List[stream]):
3          # gather message from multiple input stream
4          def gather() -> List[message]:
5              # kernel function
6              def compute(List[message]) -> List[message]:
7                  # assign batches to out streams
8                  def scatter(List[message]):
```

Figure 4. The stream module abstraction.

Also, the decomposition enables to use different batching designs and better-provisioned device resources for each phase, considering their distinct batching property and SLOs.

Insight II: Integrating stream programming model. The above decomposed execution can be essentially mapped to the stream programming model, which adopts module-level programming that decomposes applications into a set of modules operating on data streams (i.e., channels). However, existing DNN frameworks adopt tensor-level programming that models a program as a dataflow graph, where the nodes are operators and the edges are the tensors. The misaligned programming model is the major obstacle to bringing the benefits of modularity to deep learning in existing frameworks. We therefore advocate combining both programming models to enable a new DNN framework NEUSTREAM for efficient dynamic model inference.

3 NEUSTREAM Abstraction

In this section, we introduce the stream programming and execution models of NEUSTREAM.

3.1 Programming Model

The programmer writes a module-level program that manipulates named, first-class *streams* and stream *modules*: function units that take messages from and produce messages into streams. The computation logic within each stream module in turn, can be defined with a tensor-level program.

Stream module. A stream module provides the basic functionality for stream processing of deep learning inference requests. As shown in Figure 4, it is a schedulable entity consisting of the following blocks:

stream. A stream is a data channel that passes data between stream modules. Logically, a module can be regarded as a process in the concept of operating systems and a stream can be conceived as a pipe that connects different processes. Each `message=[header, data]` in the stream queue includes (1) a header containing identifying and routing information, and (2) data tensors. The header consists of a set of fields, which are set in the course of producing and delivering a message. Messages are ordered into streams according to priorities associated. Each stream module could contain (multiple) input streams `s_in` and (multiple) output streams `s_out`, specified in the `_init_` function.

```

331 1   class Stable-Diffusion:
332 2     #stream module definitions
333 3     class CLIP(StreamModule):
334 4       def compute(m:List[message])->List[message]
335 5         ... # processing code
336 6     class UNet(StreamModule):
337 7       def compute(m:List[message])->List[message]
338 8         ... # processing code
339 9       def scatter(m:List[message]):
34010         for msg in m:
34111           if msg.header.step > 0:
34212             self.s_out[0].push(msg)
34313           else:
34414             self.s_out[1].push(msg)
34515     class VAE(StreamModule):
34616       def compute(m:List[message])->List[message]
34717         ... # processing code
34818     #stream initialization
34919     stream text, latent, output, image;
35020     #define model pipeline via stream connection
35121     clip = CLIP([text], [latent])
35222     unet = UNet([latent], [latent, output])
35323     vae = VAE([output], [image])

```

348

Figure 5. Program for stable diffusion model inference.

350

gather. The gather function is used to get a *batch* from each input stream, which is a list of messages `List[message]`. If a module contains multiple input streams, the gather could combine the multiple messages of the same request from different input streams into one merging message using aggregation functions like sum, max, or concatenation.

compute. A stream module's computation is specified by a user-defined function `compute`, which takes *gathered* batches as input, performs processing, and generates the output batches. The user could directly define operations on `message.data` tensors in the function with tensor-level programming languages.

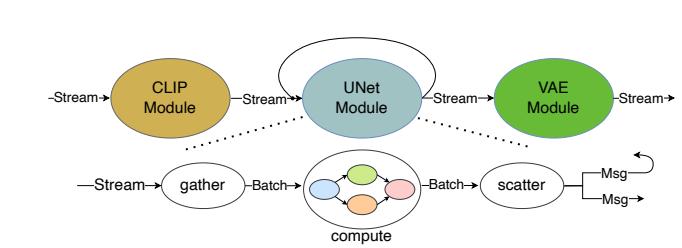
scatter. Given multiple output streams, the dynamism-related message routing operations (i.e., control flow) are decoupled from the compute of a module. scatter provides the programmer a unified function that supports customized rules to route messages to multiple output streams. The function takes the output batches of `compute`, and routes each message in every batch to one output stream based on routing information in the message head. This allows to express the control-flow of the multiple execution paths of dynamic neural networks in Figure 1. To support the iterative processes like diffusion model or LLMs, stream modules could use the same stream as both the output and the input stream. The messages routed to this output stream actually pass back to the input stream for the next iteration.

Stream program. Under our abstraction, a DL program is a collection of stream modules connected by streams. Each module consumes messages from its input streams and produces results on its output streams. The program consists of two blocks: the stream modules definitions and stream modules interconnects. Figure 5 illustrates the program of the Stable Diffusion. It first uses stream modules to define

```

1   class LLM:
2     #stream module definitions
3     class Prefill(StreamModule):
4       def compute(m:List[message])->List[message]
5         ... # processing code
6     class Decode(StreamModule):
7       def compute(m:List[message])->List[message]
8         ... # processing code
9       def scatter(m:List[message]):
10         for msg in m:
11           if msg.header.last_token != "EOS":
12             self.s_out[0].push(msg)
13           else:
14             self.s_out[1].push(msg)
15     #stream initialization
16     stream prompt, prefix, output;
17     #define model pipeline via stream connection
18     prefill = Prefill([prompt], [prefix])
19     decode = Decode([prefix],[prefix, output])

```

Figure 6. Program for LLM inference.**Figure 7.** Module graph for Stable Diffusion.

three phases: text encoder (CLIP), backbone (U-Net), and image decoder (VAE). Next, the program initializes streams and defines the model inference pipeline via module stream connection. The CLIP module transforms the input prompt into a *latent stream* including text embedding and initial noisy latent, which is connected to the U-Net. The UNet module iteratively denoises the latent image messages in the *latent stream*, and outputs the denoised latent messages. We decrease the iteration number of latent messages for each U-Net execution, and route denoised latent messages with the positive number of iterations back to *latent stream* for UNet to execute the next iteration, and leave others to the *output stream* connecting to VAE. The VAE decoder transforms the latent messages back into images. Figure 6 illustrates the program of the LLMs, which contains two stream modules. The prefill module processes a user's prompt to generate the first token (including KV cache) of the response into *prefix stream*. Following it, the decoding module sequentially generates subsequent tokens (and updates KV cache) in multiple steps. In particular, after running one decode step, the termination condition will be checked: if a special stop token is met, we route the token to *output stream*, otherwise, each generated token is appended and fed back into *prefix stream* for the decode module to generate the next token.

386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440

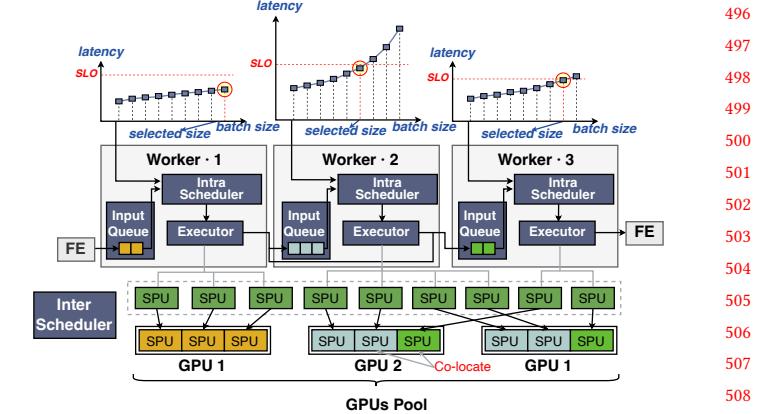
441 3.2 Execution Model

442 The execution of a stream program can be modeled as a
 443 *stream graph*: a hierarchical (two-level) computation graph
 444 of modules communicating via data streams, in which each
 445 module is further represented as a dataflow of lower-level
 446 operations on messages. Figure 7 illustrates the stream graph
 447 corresponding to Figure 5. Each module performs local com-
 448 putation on input streams from its in-edges and produces
 449 output streams as its out-edges. The stream graph abstraction
 450 supports the gather-compute-scatter operations defined
 451 in stream modules. A module can maintain a local state s . Its
 452 execution starts with its initial state s_0 and proceeds in steps.
 453 In each step, the module consumes the next batch from its
 454 input streams by calling `gather`, updates its state and pro-
 455 duces an output batch by executing a series of operations
 456 of `compute`, and routes output messages into downstream
 457 modules by `scatter` according to the stream connection.

458 Different from the dataflow of existing DL frameworks,
 459 the stream graph has two major features. First, instead of
 460 mixing control-flow and dataflow together, the stream graph
 461 decouples the module communication (control flow) and
 462 computation (dataflow) into high-level module graph and
 463 low-level operation graph, respectively. This enables effi-
 464 cient batching at the granularity of modules, given that data
 465 in a single stream have the same static operation graph de-
 466 fined in `compute`. Second, modules are connected via streams
 467 with the continuous flow of data, instead of static tensors
 468 consuming exactly once. This supports a decomposed and
 469 continuous stream computation model, extended with the
 470 ability to customize execution at each module (e.g., batching
 471 and resource allocation).

473 4 NEUSTREAM Scheduling

474 The goal for NEUSTREAM is to maximize the system goodput,
 475 i.e., the number of served requests that meet their SLO re-
 476 quirements per unit time. The hierarchical streaming graph
 477 naturally leads to a two-level scheduling framework: *intra-*
 478 *module* scheduling on how many requests to include in a
 479 batch for a specific module, and *inter-module* scheduling
 480 on how much resource to be allocated across modules. Fig-
 481 ure 8 illustrates the architecture of NEUSTREAM, which is
 482 composed of a frontend and multiple stream workers. Each
 483 worker is responsible for executing a single module. The front-
 484 end continuously sends incoming requests to the first stream
 485 worker. Upon receiving request messages, the executor in
 486 a worker performs the gather-compute-scatter execution
 487 model, sending new messages to downstream workers based
 488 on stream connection. The batching decision is made by the
 489 intra-module scheduler at each worker, utilizing partial SLO
 490 and profiled latency function of each module. The resource
 491 allocation decision is made by the inter-module scheduler,
 492 which aims to balance the maximum goodput permitted by
 493 each module under an appropriate resource allocation.



500 **Figure 8.** NEUSTREAM’s two-level scheduling framework.

501 **Algorithm 1:** Intra-module scheduler for batching
 502 on a module m_i

503 **Input:** Latency function $L_i(b, a_i)$ of batch size b and
 504 resource allocation a_i , partial SLO_i , input stream in

505 **Output:** batch_request

506 **Function** IntraScheduler(L_i, SLO_i, in):

507 batch_request $\leftarrow []$;

508 // get maximum batch size based on partial SLO

509 batch_limit $\leftarrow argmax_b L_i(b, a_i) \leq SLO_i$;

510 **for** request r in in **do**

511 **if** $r.budget \leq 0$ **then**

512 // drop request that fail to meet upstream
 513 partial SLOs

514 $in.remove(r)$;

515 **continue**;

516 **end**

517 **if** length(batch_request) $< batch_limit$ **then**

518 | batch_request.append(r);

519 **end**

520 **end**

521 **return** batch_request;

534 4.1 Intra-module Scheduling

535 **Partial SLOs.** In the presence of control flow, it is unclear
 536 how to define end-to-end latency SLOs for individual re-
 537 quests with dynamic and even unpredictable execution paths.
 538 Since NEUSTREAM decouples routing control-flow from ex-
 539 ecution, we take the approach of assigning *Partial SLOs* to
 540 individual modules, which is similar to the partial SLO con-
 541 cept in microservices [14, 28]. Specifically, given the module
 542 execution path $\varphi = \{m_1, \dots, m_k\}$ of a certain request, we
 543 consider its end-to-end latency $SLO = \sum_{i \in \{1, \dots, |\varphi|\}} SLO_{\varphi(i)}$,
 544 where $SLO_{\varphi(i)}$ is the partial SLO for module $\varphi(i)$.

545 **Batch scheduling algorithm.** As long as all modules meet
 546 their respective partial SLOs, the request SLO is guaran-
 547 teed. This approach ensures independent batch scheduling

551 across individual modules. Algorithm 1 describes the
 552 procedure of intra-module batch scheduling for one execution
 553 step on a specific module m_i . The algorithm receives the fol-
 554 lowing inputs for m_i : (1) $L_i(b, a)$: profiled execution latency
 555 of batch size b and resource allocation a (detailed in Sec-
 556 tion 4.2), which includes the duration of executing compute
 557 and scatter, (2) SLO_i : partial SLO (in latency), and (3) in : the
 558 input stream. For every step, if the stream in is non-empty,
 559 the scheduler determines a batch: we first obtain the maxi-
 560 mum permitted batch size $batch_limit_i = argmax_b L_i(b, a) \leq$
 561 SLO_i without violating the partial-SLO of this module (line 3).
 562 Then for each request in the input stream, we add a request
 563 into the batch if it satisfies:

C1. The request meets the partial SLOs on the upstream
 565 modules. In particular, let $\varphi_r = \{m_1, \dots, m_k\}$ be the execu-
 566 tion path of previous steps of a request r before executing
 567 the module m_i at the current step, and t and t_0 be the current
 568 time and request issued time, respectively. We require the
 569 request budget $B_r = \sum_{j \in \{1, \dots, |\varphi_r|\}} SLO_{\varphi_r(j)} - (t - t_0) \geq 0$. In
 570 line with other SLO-aware scheduling approaches [8, 11, 35]
 571 that drop requests unlikely to meet their SLOs, we removes
 572 requests with negative budget from the stream to avoid con-
 573 tinuously blocking the subsequent incoming requests (line
 574 5-8). However, in more sophisticated systems, these removed
 575 requests may not be dropped entirely; instead, they can be
 576 placed in a low-priority queue or escalated to an upper-level
 577 scheduler for re-dispatch.

C2. The current batch size is smaller than the maximum
 579 permitted batch size $batch_limit$ (line 9-11).

The above two conditions together achieve a large good-
 581 put at a specific module m_i by ensuring that the requests in
 582 the batch can still meet partial SLOs after the execution of
 583 this module, and the GPU operates at high efficiency. In the
 584 stream of each module, we sort the request messages based
 585 on their arrival time at the frontend. One could also priori-
 586 tize requests with tight SLOs based on the request budget.
 587 NeuStream supports an adaptive batch size without wait-
 588 ing for the number of requests in the input stream to reach
 589 max permitted batch size, in order to avoid SLO violation of
 590 waiting requests.

4.2 Inter-module Scheduling

597 **Streaming Processing Units (SPUs).** To facilitate resource
 598 allocation, we introduce an abstraction of SPUs, which can
 599 be considered as a uniform partition of physical devices, as il-
 600 lustrated in Figure 8. Given a total of N SPUs and M modules,
 601 the inter-scheduler determines the number of SPUs a_i for
 602 each module m_i . Let k_i be the average number of execution
 603 times of a specific module m_i in the execution path. Given
 604 the batch scheduling of Algorithm 2 and resource allocation
 605

Algorithm 2: Inter-moduler scheduler for resource allocation across modules

606 **Input:** stream graph G, function partial SLO for each
 607 module, total number of SPUs N , devices D

608 **Output:** [SPU_allot, SPU_map]

609 **1 Function** InterScheduler(G, SLO, N):

610 // assign initial SPUs satisfying memory limit

611 SPU_allot \leftarrow Initialize(G, SLO);

612 **while** $N \geq \Delta$ **do**

613 // map SPU allotment to physical devices

614 sort(SPU_allot);

615 **for** $i \in \{1, \dots, M\}$ **do**

616 SPU_map[i] = best_fit(SPU_allot[i], D);

617 **if** $SPU_map[i] = \emptyset$ **then**

618 | $SPU_map[i] = partition(SPU_allot[i], D)$;

619 **end**

620 $L_i(b, a_i) = get_profile_latency(SPUs_map[i], D)$;

621 **end**

622 // assign SPUs to bottleneck module

623 $m = argmin_{i \in \{1, \dots, M\}} g_i(SPU_allot[i], SLO_i)$;

624 SPU_allot(m).increase(Δ);

625 $N = N - \Delta$;

626 **end**

627 **return** [SPU_allot, SPU_map];

630 a_i , we define the normalized goodput $g_i(a_i, SLO_i)$ of m_i as:

$$632 g_i(a_i, SLO_i) = \frac{batch_limit_i}{SLO_i \times k_i} = \frac{argmax_b L_i(b, a_i) \leq SLO_i}{SLO_i \times k_i} \quad (1)$$

635 The goal of inter-scheduler is to optimize the overall pipeline
 636 by maximizing the minimum module normalized goodput:

$$638 \begin{aligned} & \max \min_{i \in \{1, \dots, M\}} g_i(a_i, SLO_i) \\ & \text{s.t. } \sum_i a_i \leq N \forall i \in 1, \dots, M \\ & \quad a_i \cdot mem \geq Mem_i \forall i \in 1, \dots, M \text{ (Memory limit)} \end{aligned} \quad (2)$$

643 where mem is the memory capacity of a SPU and Mem_i is
 644 the memory footprint of module m_i .

645 **Allocation algorithm.** Algorithm 2 presents pseudocode
 646 for our iterative SPU allocation process. The algorithm first
 647 initializes all modules by allocating the minimum number of
 648 SPUs a_i that satisfies the memory limit condition for each
 649 module m_i (line 2). This step ensures that each module can
 650 be put on the SPUs under the memory constraint. Then, the
 651 algorithm finds the modules that yield the minimum normal-
 652 ized goodput, and allocates Δ ($\Delta \geq 1$) SPU to that module
 653 (line 12-14). To this end, for each module, the algorithm maps
 654 its allocated SPUs onto physical devices (e.g., GPUs) through
 655 module placement detailed below (line 4-9), and gets the
 656 profiled latency function $L_i(b, a_i)$ of the module (line 10).
 657 Based on the latency functions and equation (1), the algo-
 658 rithm computes normalized goodput g_i of each module m_i .

661 The algorithm continues allocating Δ SPUs at a time until
 662 all the SPUs are allocated.

663 **Module placement.** When mapping SPUs to physical devices
 664 (line 4-9 in Algorithm 2), the scheduler considers modules in descending order of SPU allocation, maps SPUs of
 665 a module onto a device that has the least available SPUs to
 666 accommodate the module (line 5). If no single device has
 667 sufficient available resources to accommodate a module, we
 668 uniformly map SPUs of a module onto the fewest possible
 669 devices to reduce communication costs (line 7-9). In this
 670 case, we rely on model-parallel partition tools to partition a
 671 module (operation graph) across multiple devices. We can
 672 enumerate possible combinations of data parallelism and
 673 model parallelism, and use the profiled latency function to
 674 identify the configuration that yields the best module good-
 675 put, i.e., $\max argmax_b L_i(b, a_i) \leq SLO_i$.
 676

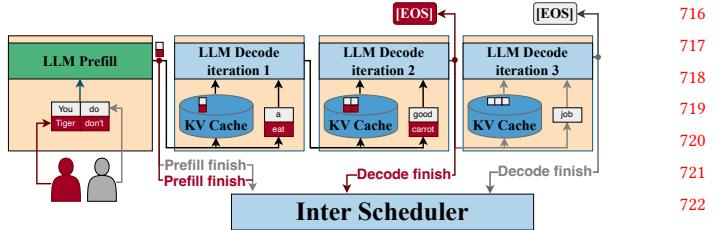
677 **Colocating modules.** The scheduler allows the two mod-
 678 ules to be annotated as *co-located*, and maps the SPUs of two
 679 co-located modules on the same set of devices to reduce the
 680 communication and increase the data locality. Given a work-
 681 load, NEUSTREAM could automatically search for the mod-
 682 ules to be co-located: For a pair of originally non-collocated
 683 modules (m_i, m_j), the scheduler re-maps SPUs by annotating
 684 them as *co-located*, and treats the annotation valid only if
 685 co-locating m_i and m_j could improve their overall normal-
 686 ized goodput, i.e., $\min\{g_i, g_j\}$, based on the profiled latency
 687 functions and equation (1).

690 5 Implementation

691 NEUSTREAM was implemented with 5K lines of Python code.

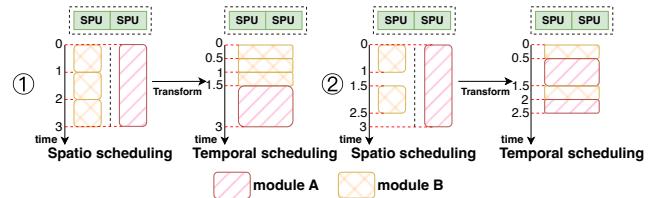
692 **Module execution.** We built the executor of diffusion mod-
 693 els and LLM models on top of PyTorch [2] and vLLM[10],
 694 respectively. For diffusion models, we directly use PyTorch
 695 for sub-module parallelism across multiple cores on a single
 696 GPU. For larger LLMs, we use vLLM to provide sub-module
 697 parallelism across multiple GPUs. Each module is executed
 698 in its own process by default, however, co-located modules
 699 can be executed in the same process to optimize module
 700 switching and communication overhead (e.g., only passing
 701 pointers to the output data that's already on the device). The
 702 module invocation number depends on the workload, which
 703 can be determined statically before execution in diffusion
 704 or dynamically during execution in LLMs with the EOS to-
 705 ken being detected. To compute the normalized goodput, we
 706 could use the average invocation number over the requests
 707 that have already been served.

708 **State management.** Each module can optionally manage
 709 the state for each request that can be used across differ-
 710 ent messages. The state of a module is a set of key/value
 711 pairs, which can be stored in a key-value store and operated
 712 by a module executor via get, put, and delete operations.
 713 NeuStream's inter-module scheduler tracks the *progress* of a
 714



716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770

Figure 9. adjust font size State management for LLM.



734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770

Figure 10. adjust figure font Spatial-to-temporal scheduling transformation with EFTF algorithm in a single device.

specific request r and notifies a module to evict its state once the module completes the processing of r . In cases where a module is repeatedly invoked within a loop in the stream graph, the module is only considered finished when r exits the loop. Figure 9 illustrates how NEUSTREAM handles state in LLMs. The decode module is stateful by treating the KVCache as the state of a request and updating the state (KVCache) at each iteration by inserting the key and the value of the new token. To reduce memory fragmentation, following PagedAttention [10], NeuStream partitions the KVCache into blocks along the sequence length dimension, with each block containing KVCache for multiple tokens. These blocks are represented as keys in the key-value store. At each decoding step, the decode module retrieves the addresses of these blocks, performs attention directly on the distributed blocks, and appends the KVCache of the newly generated token to the last block. Once the inter-module scheduler tracks that the decode module (loop) has completed processing for a certain request (e.g., encountering an EOS token), it notifies the decode module's executor to evict the state of that request.

Streaming Processing Units. We could leverage spatial partitioning such as Multi-process-service (MPS[16]) to split a GPU resource into multiple SPUs. However, it cannot guarantee performance isolation and often results in unpredictable throughput, and also cannot adapt to workload variations. Thus, we transform the spatial scheduling into an equivalent or even better temporal scheduling by leveraging the predictability in the latency of a single module. Specifically, we adopt an Earliest Finish Time First (EFTF) algorithm, if modules are executed by the SPUs on the same device, we order the batches of different modules in increasing order of finish time in *spatial* scheduling, and schedule

them in this order with transformed *temporal* scheduling. Figure 10 gives an illustrative example. In spatial scheduling, two SPUs are located in the same device, and each is responsible for a module. By contrast, as a module takes up the two SPUs in temporal scheduling, the execution time of a single module is halved. In the first case, the batch of module A arrives continuously, and we insert three module A’s batches before module B’s batch in temporal scheduling based on their finish time in original spatial scheduling. In the second case, the second module A’s batch arrives after the scheduling of module B’s batch, but with an earlier finish time in original spatial scheduling. So we allow module A’s batch to be executed concurrently with module B’s using a new CUDA stream with high priority. We see that each batch’s finish time is no later than that in spatial scheduling.

Stream. We borrow a similar concept of data streaming from traditional streaming systems [4, 27] and additionally perform specific optimizations for DNNs: the data flit in a DNN typically consists of large tensors, so we pass only data references in the stream to avoid heavy data copying. To this end, only the metadata is put into the stream to facilitate scheduling, while the actual tensors are fetched as needed.

Profiling. Given a total of K SPUs on a device, if a module is executed on n devices with a_i SPUs, we first profile the latency function of $L_i(b)$ by taking the entire n devices, and estimates the latency function $L_i(b, a_i) = \frac{L_i(b)nK}{a_i}$ for the get_profile_latency interface.

6 Evaluation

6.1 Experiment Setup

Models. We evaluate NEUSTREAM on three types of modern model families listed in Table 1. For diffusion, we select widely used models for different tasks, which include text-to-image (Stable Diffusion [19]), image restoration (Palette [22]), and class-conditional image generation (DiT [18]). For LLM model, we select the OPT [36]. We choose several model sizes to test performance under different numbers of devices. For multi-agent, we use MatPlotAgent [31] designed to automate scientific data visualization tasks. The model is comprised of three LLMs: a Code Llama [21] to generate preliminary code, LLaVA [13] to generate image-text question/answer, and the second Code Llama to refine the initial code.

Cluster testbed. We evaluate NEUSTREAM using three types of servers equipped with different accelerators: The first one is a local workstation with Intel Xeon Gold 6230R CPUs and NVIDIA GeForce RTX 4090 (24GB) GPUs, running on Ubuntu 20.04 with CUDA 11.8. The second one is a cloud server with Intel Xeon Platinum 8480C CPUs, 8*NVIDIA Hopper H100 (80GB) GPUs with NVLink enabled, running on Ubuntu 22.04 with CUDA 12.0. The third one is a cloud

Name	Type	Task	Year
StableDiffusion [19]	Diffusion	Text to Image	2022
Palette [22]	Diffusion	Image Restoration	2022
DiT_XL/2 [18]	Diffusion	Image Generation	2023
DiT_S/2	Diffusion	Image Generation	2023
OPT-6.7B [36]	LLM	Text Generation	2022
OPT-13B	LLM	Text Generation	2022
OPT-30B	LLM	Text Generation	2022
OPT-66B	LLM	Text Generation	2022
MatPlotAgent [31]	Multi-agent	Multi-agent Assistant	2024

Table 1. Experiments setting.

server with AMD EPYC 7513 and 4 RTX A6000 (48GB) GPUs, running on Ubuntu 18.04 with cuda 12.2.

Offered load. Although Microsoft Azure has issued two real traces: Microsoft Azure function trace 2019 (MAF1) [23] and 2021 (MAF2) [37], those existing datasets do not include timestamps for the newly released diffusion models and LLMs. We focus our attention on various traffic patterns and SLOs that have a significant effect on batching behavior. Similar to prior work [11], we generate the request arrival times using Gamma Process, which is parameterized by rate and coefficient of variance (CV). By scaling the rate and CV, we can control the workload’s rate and burstiness, respectively. For diffusion models, we use two types of image sizes (256×256 , 512×512) and use a uniform iteration distribution over the recommended range, e.g., [30 – 50] for Stable Diffusion [1]. For LLM models, we use LMSYS-Chat-1M [38], a large-scale dataset containing real-world conversations with 25 state-of-the-art LLMs. For multi-agent, we use MatPlotBench [31] consisting of complex data visualization problems.

Metrics. We use *goodput* to measure the performance of the systems, which is the number of served requests that meet their SLO requirements per unit of time. We set partial SLO of each module based on its inference time on a specific hardware. Diffusion model serving focuses on request-level SLO which is the sum of partial SLOs, whereas LLM serving focuses on token-level partial SLOs on the time to first token (TTFT) and average time per output token (TPOT) for prefill and decode modules, respectively. Each experiment was repeated 5 times.

Baselines. We choose the most representative baselines with a focus on batching: for Diffusion with predictable execution times, Clockwork serves as a strong baseline using a proactive scheduling approach. The original Clockwork implementation is not fully functional for supporting modern diffusion models. It uses TVM as the model execution backend. Our testing shows that the Clockwork-compatible TVM version lacks some operator support for diffusion. Thus, we re-implemented Clockwork by replacing TVM with PyTorch. Similar to prior work (e.g., AlpaServe[11], Shepherd[35]), the same PyTorch backend makes systems fairly evaluated, ensuring that performance differences are solely due to the

881 scheduling design. For LLMs and Multi-Agents, we select
 882 vLLM [10], which is a state-of-the-art(SOTA) baseline sup-
 883 porting continuous batching [32], PagedAttention [10], chun-
 884 ked prefill [3], and other advanced optimizations. We adopt
 885 the latest vLLM version (v0.5.4) released in August 2024.
 886

888 6.2 Results on Diffusion Models

889 For Diffusion models that can fit into a GPU, we focus on a
 890 single GPU placement.

892 **Goodput under varying offered load.** Figure 11 shows
 893 the results of serving diffusion models with 256x256 image
 894 size on RTX 4090 GPU. In each row, we change one parame-
 895 ter to see how workload variation affects the performance of
 896 each serving system. Figure 11’s first row shows the goodput
 897 with varying rates. NEUSTREAM provides significantly higher
 898 goodput due to batching more requests on the module level
 899 (e.g., up to 5.26× at 4 requests/s in the DiT-S/2 model). Fig-
 900 ure 11’s 2nd row varies the CV of the workloads. The traffic
 901 becomes more bursty with a higher CV, which increases the
 902 possibility of SLO violation, resulting in quick goodput re-
 903 duction in Clockwork. By contrast, the more bursty increases
 904 the batching opportunities in NEUSTREAM, which alleviates
 905 the negative effects of high bursty requests (outperforming
 906 Clockwork by 1.37 ~ 4.04× across models at a high CV of 4).
 907 Figure 11’s 3rd row shows the effect of different SLO. SLO
 908 scale equals 1 means the latency requirement is the single
 909 request inference latency. As SLO scale increases, the latency
 910 requirement becomes looser. We see that even SLO is only
 911 slightly larger than the model inference time, NEUSTREAM
 912 gets a much higher throughput due to module-level batching
 913 (up to 3.13× at a small SLO scale of 1.2).

914 The computational cost of StableDiffusion depends on the
 915 input image size. For large 512×512 images, as shown in
 916 Figure 12, it almost saturates the RTX 4090 GPU even with
 917 a single request, limiting the batching gain of NEUSTREAM.
 918 For smaller 256x256 images, the GPU is not fully utilized, al-
 919 lowing us to benefit more from NeuStream’s batching. How-
 920 ever, we observe that NEUSTREAM can effectively scale up
 921 to a more powerful H100 GPU (up to 1.90× as compared to
 922 Clockwork across offered loads). This is because the H100’s
 923 computation capability is significantly better than RTX 4090,
 924 and it can benefit from batching even in the 512x512 image
 925 generation task.

926 **Batching efficiency.** To fully understand NEUSTREAM’s
 927 benefits on batching, we plot the running batch size on DiT
 928 module for DiT-S/2 along the time in Figure 13(a). On average,
 929 the batch size of NEUSTREAM is 5.35× larger than Clockwork
 930 due to efficiently batching requests that (1) arrive at different
 931 times; (2) require different number of iterations to finish;
 932 or (3) start with different number of iterations. Figure 13(b)
 933 increases the batching opportunity of Clockwork by setting
 934 the same number of iterations for each request, i.e., disabling
 935

the above case (3). We see that NEUSTREAM still has a larger
 936 batch size (1.67× on average larger).

938 6.3 Results on LLM Models

939 For LLMs, our evaluation covers multi-resource scenarios
 940 (e.g., OPT-30B/-66B across 2/4 GPUs). We focus on the good-
 941 put of decode module that has a much larger impact on the
 942 end-to-end latency.

943 **Goodput under varying offered load.** Figure 14 shows
 944 the results of serving OPT models on RTX A6000 GPUs. The
 945 1st row shows that, as the request load increases, vLLM good-
 946 put starts to decrease significantly. This is because vLLM
 947 greedily prioritizes prefill requests and the use of preemp-
 948 tion leads to resource starvation of decode module. In con-
 949 trast, NEUSTREAM is SLO-aware and balances the resource
 950 provision across the modules. For OPT-13B, NEUSTREAM
 951 outperforms vLLM by 1.53× at 2 requests/s. Under an even
 952 higher rate of 4 requests/s, NEUSTREAM achieves 37.21× more
 953 goodput than vLLM, due to decode resource starvation and
 954 prefill congestion (caused by insufficient memory for KV
 955 cache) in vLLM. As shown in Figure 14’s 2nd and 3rd rows,
 956 NEUSTREAM performs better than vLLM to handle more
 957 burstiness (a higher CV) and stringent SLOs (a smaller SLO
 958 scale) that aggravates the resource contention, e.g., 2.25× at
 959 the CV of 4 and 7.98× at SLO scale of 1.2 on OPT-13B. Fig-
 960 ure 14’s 2nd and 3rd columns show NEUSTREAM consistently
 961 outperforms vLLM when scaling the serving to accommodate
 962 larger models on multiple devices (e.g., NEUSTREAM achieves
 963 up to 69.31× at 3 requests/s when serving the largest OPT-
 964 66B model on 4 devices).

965 Figure 15 further shows the result of serving OPT models
 966 on more powerful NVIDIA H100 GPUs, where each system
 967 is able to handle a relatively higher offered load as compared
 968 to that in A6000 GPUs. We see that NEUSTREAM can still
 969 achieve significantly more goodput as compared to vLLM
 970 on various traffic patterns and SLOs (e.g., up to 11.44× at
 971 8 requests/s when serving the largest OPT-66B model on 4
 972 H100 GPUs).

973 **Normalized goodput of modules.** For different LLMs,
 974 Figure 16 shows the normalized goodput of each module
 975 on A6000. We see that NEUSTREAM’s allocation ensures a
 976 balanced processing speed between stages, since our inter-
 977 module scheduling explicitly allocates SPUs to maximize the
 978 minimum module normalized goodput. By contrast, vLLM
 979 exhibits a significant imbalance between the two stages, as it
 980 greedily assigns resources to prefill with the use of preemp-
 981 tion. Thus, vLLM overall goodput is limited by the slower
 982 (or bottleneck) decode module.

983 6.4 Results on Multi-agent Model

984 For the multi-agent model, different LLM agents are placed
 985 on different GPUs.

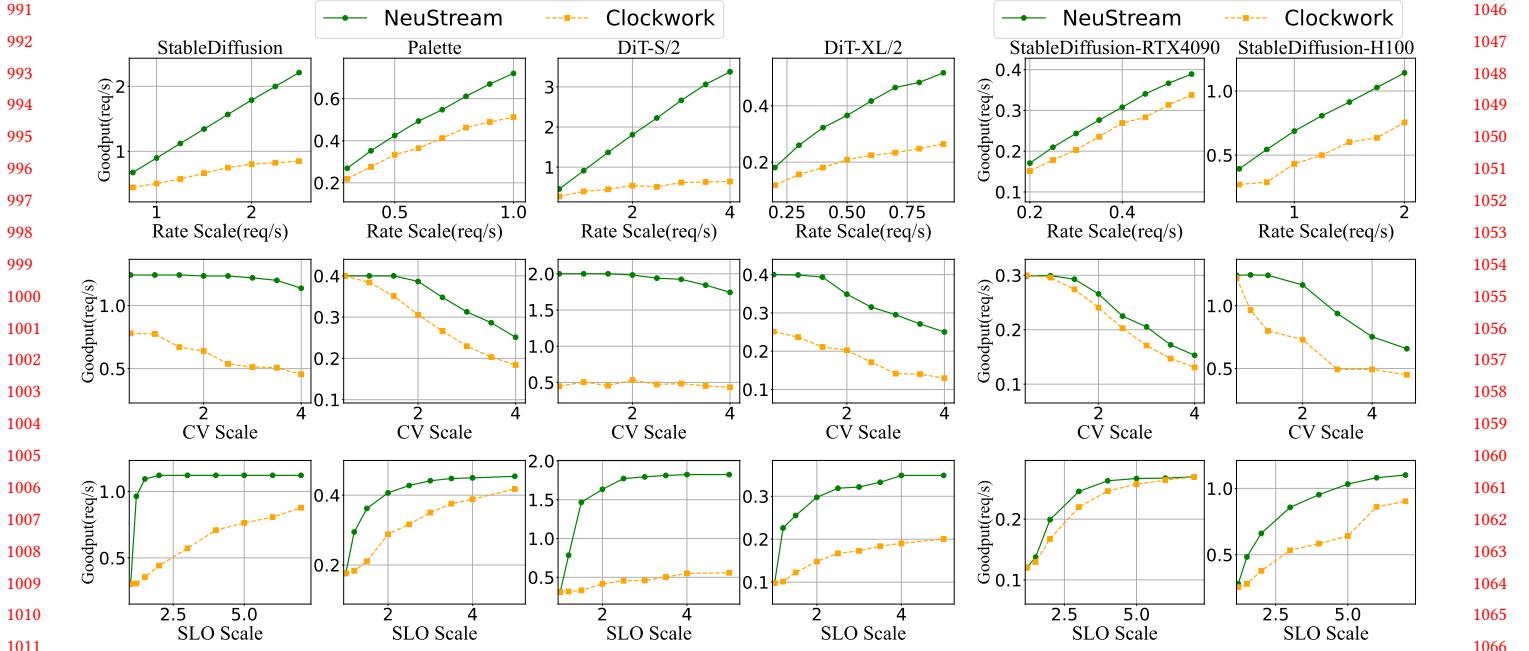


Figure 11. Experiments of generating 256*256 images on NVIDIA RTX4090.

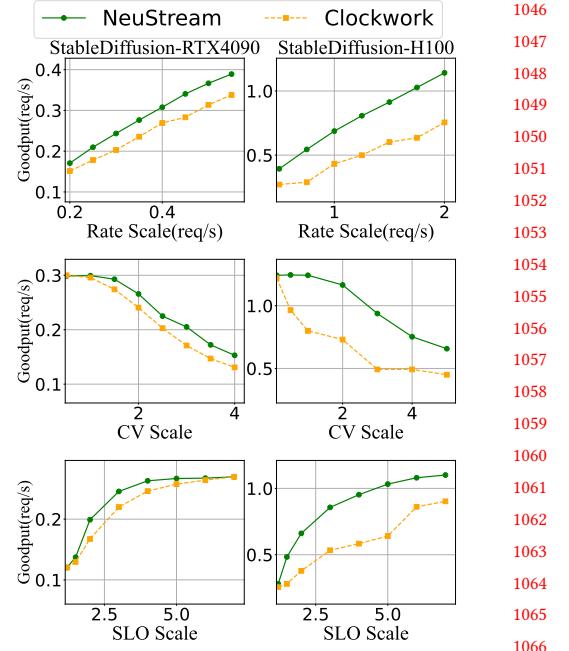


Figure 12. Generating 512*512 images on RTX4090 and H100, respectively.

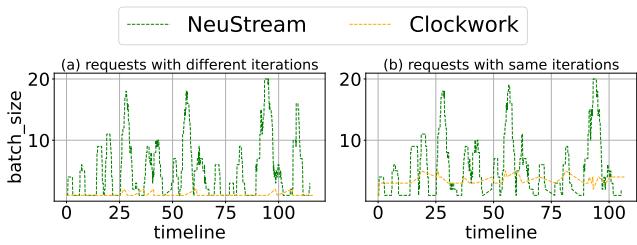


Figure 13. Batch size of DiT module along the time, where requests start with different/same number of iterations.

Goodput under varying offered load. Figure 17 and Figure 18 show the goodput of decode module for each LLM agent on RTX 4090 and H100, respectively. NEUSTREAM significantly outperforms vLLM across LLM agents (e.g., achieving 13.46× at 1 requests/s, up to 1.54× at CV of 4, and 2.93× at SLO scale of 1.4 on RTX 4090 GPUs, meanwhile, achieving 1.38× at 2.2 requests/s, up to 1.2× at CV of 4, and 2.11× at SLO scale of 1.6 on NVIDIA H100 GPUs).

Normalized goodput of modules. Figure 19 shows the goodput of each LLM agent’s module on RTX 4090. We see that NEUSTREAM achieves balanced speeds across stages of LLM agents, alleviating the bottleneck. Both the normalized goodput of prefill/ decode stages of vLLM is less than NEUSTREAM. This is because imbalanced prefill and decode phases in vLLM lead to high prefill congestion due to insufficient memory to hold new KV cache.

	Cost/Model	SD	OPT	M-agents
Profile	4.32 min	5.2 min	11.37 min	
Schedule	1.02%	1.95%	5.31%	
Comm.	colocate	<0.1%	<0.1%	/
	non-colocate	2.6%	12.10%	0.1%
	PCIe	0.66%	0.98%	/
	NVLink	0.81%	<0.1%	0.22%
Switch	3.62%	3.39%	4.79%	
Memory				

Table 2. Execution Cost for Stable Diffusion (SD), OPT-30B (OPT), and multi-agents (M-agents).

6.5 Execution Cost

Table 2 reports the execution cost of NEUSTREAM on different types of models.

Profiling. Profiling the latency function is a one-time *pre-processing* stage (e.g., only a few minutes), which is negligible compared to the subsequent long serving time (e.g., hours and days). For diffusion, profiling only depends on batch size. For LLMs, the profiling further depends on input token length. To reduce overhead, we use profiling and interpolation to figure out the latency function as DistServe [39].

Scheduling. Running the scheduling algorithm might incur an overhead of batching decisions, resource allocation and state tracking. We see that NEUSTREAM only incurs a slight scheduling cost (< 2%) in LLMs, and still very little overhead (< 5.3%) even more LLMs are involved in multi-agents.

Communication. We test different inter-communication scenarios: For co-located case, the communication overhead

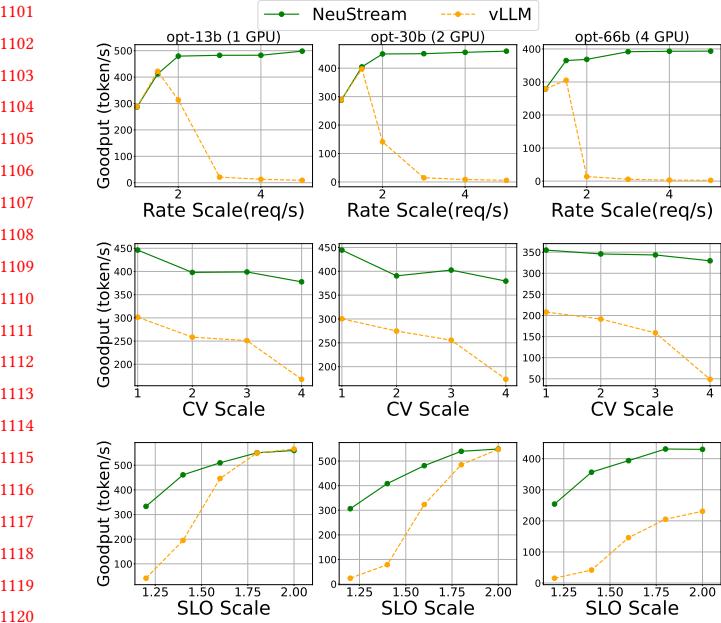


Figure 14. Goodput of LLM decode module on A6000.

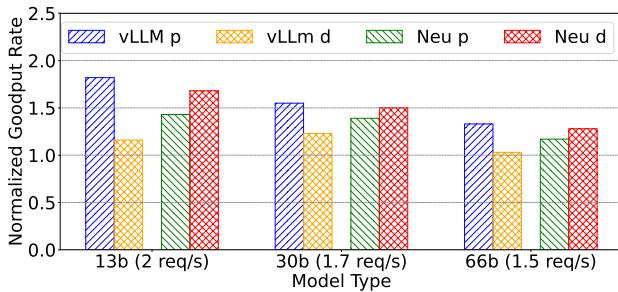


Figure 16. Module's normalized goodputs of each LLM.

is not visible in end-to-end request latency as modules only pass data pointers. For non-colocated case, we place each module (or agent) on a different GPU for SD and LLM (or multi-agents). We still observe low overheads for diffusion and multi-agents due to only transferring small activation messages. For LLM, transferring KV cache via PCIe in A6000 could incur a relatively large overhead, which can be alleviated via either co-locating or transferring via NVLink in H100.

Switching. We see that the overhead of switching between modules is little for all the models (e.g., <1%) as the most context of module resides in the GPU memory when switching between modules (e.g., the model weight used by both prefill and decode modules).

Memory. With the proposed model of time-sharing to implement the SPU execution, more tasks of heterogeneous GPU kernels are concurrently executed. The memory overhead refers to the extra memory consumption (e.g., activation

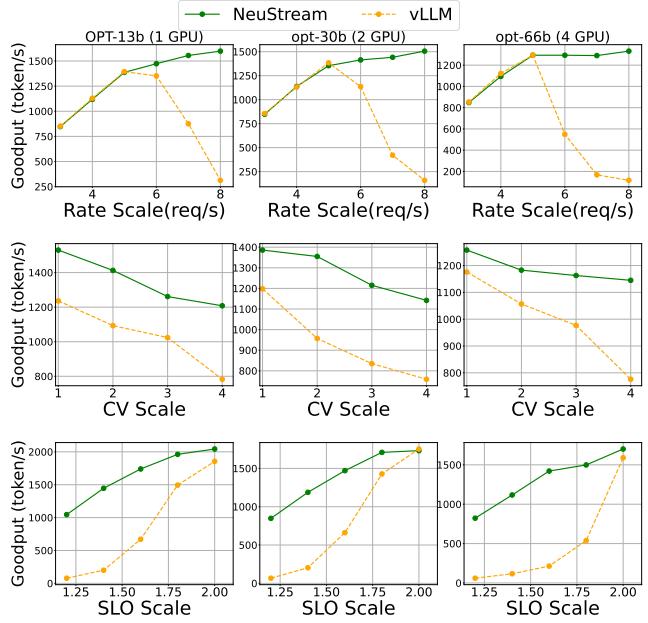
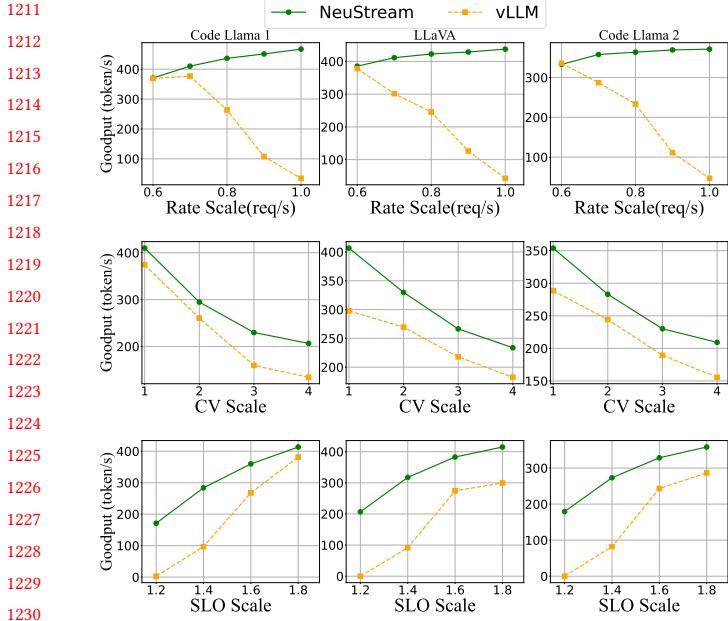
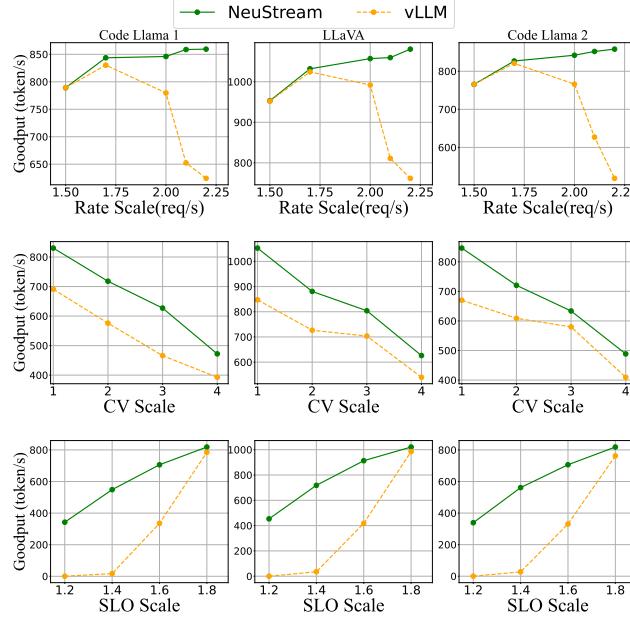
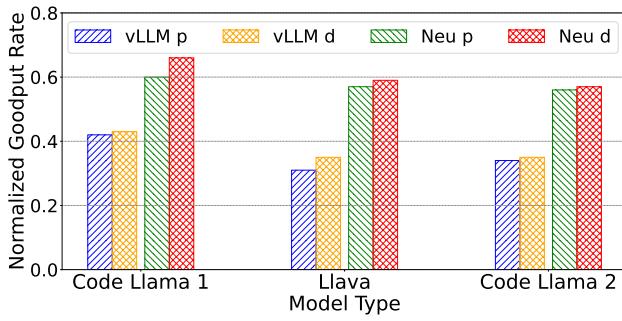


Figure 15. Goodput of LLM decode module on H100.

values) for concurrent kernels, as compared to the baseline memory usage without time sharing (e.g., kernel executed sequentially). We see that time sharing does not introduce significant memory overhead in our workload (e.g., < 5%), since the workspace mainly consists of small activations. Also, our module granularity (e.g., sub-models) does not split individual or fused kernels, without yielding additional memory traffic.

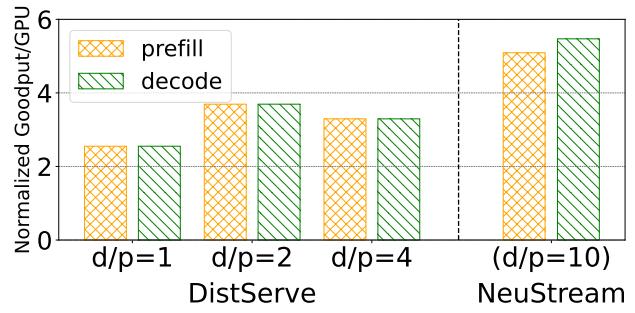
7 Related work

DNN serving systems. Clockwork [8] leverages the predictable execution times of individual DNN inferences. Shepherd [35] aggregates request streams into moderately sized groups to improve predictability. Systems like Nexus[24], INFaaS [20], gputlet [5], Gslice [7], focus on the serving of multiple static DNN models with spatio-temporal sharing of GPUs. AlpaServe [11] further introduces model parallelism for statistical multiplexing. However, these systems fail to take into account the dynamic property of modern DNNs, resulting in missed opportunities for optimization. Systems for dynamic DNNs have also been studied in recent years. BrainStorm [6] proposes a new data abstraction called Cell to optimize the dynamism of different parts within a single input, e.g., a token inside a request for Mixture-of-Experts (MoE). Nimble [25] and Cocktailer [34] express and execute dynamic neural networks from a compiler perspective. NEUSTREAM is orthogonal to these works by exploring the dynamism and batching opportunity across multiple inputs and supports SLO-aware DNN serving.

**Figure 17.** Goodput of agent’s decode module on RTX4090.**Figure 18.** Goodput of agent’s decode module on H100.**Figure 19.** module normalized goodputs in MatPlotAgent

Specialized systems for LLM serving. LLM inference serving is a fast-developing area. Orca [32] introduces iteration-level scheduling to increase the GPU utilization. vLLM [10] proposes a novel memory management strategy PagedAttention for KVCache to increase the throughput of LLM serving. These systems are optimized and specialized for LLMs. our NEUSTREAM abstraction not only naturally enables the aforementioned optimizations for LLM models but also generalizes to other DNN models (e.g., as shown in Figure 1).

Streaming systems. Traditional streaming processing systems such as Naiad[15], Spark [33], Flink [4], StreamScope [12], Apache Storm [27], etc., process continuous data flows in real-time. However, these streaming systems primarily handle I/O-intensive data such as application logs, sensor events, etc. In contrast, DNN applications are often compute-intensive, requiring the execution of tensor-based inputs in large batches

**Figure 20.** Normalized per-GPU goodput of instances with OPT-13B model on LMSYS-Chat-1M dataset and H100.

on GPUs to achieve high throughput. This leads us to different trade-offs and design options for NEUSTREAM.

8 Discussion

Disaggregated LLMs. Although the idea of decomposing a LLM workload has emerged in more recent work (e.g., Dist-Serve [39] and Splitwise [17]), NEUSTREAM presents a generic decomposition abstraction and system that can be used for different workloads, not just specialized for LLMs. Moreover, based on this abstraction, NEUSTREAM offers additional optimization opportunities, such as optimizing resource allocation and SLO-aware scheduling at fine-grained streaming modules and SPUs. For example, given a pair of prefilling and decoding instances, DistServe allows to assign different degrees of model parallelism (e.g., tensor parallelism) to each instance. Figure 20 shows the normalized per-GPU

1321 goodput of each instance in DistServe, when serving OPT-
 1322 13B on LMSYS-Chat-1M dataset used in Sec. 6.3. Here, D/P
 1323 gives the resource usage ratio of decoding instance to prefill
 1324 instance. Note that the minimum granularity of resource al-
 1325 location in DistServe is the entire GPU. Given the relatively
 1326 low average input-output ratio of requests, we observe that
 1327 orchestrating prefill and decoding in DistServe is challeng-
 1328 ing at the coarse-grained GPU allocation model: under tight
 1329 SLO constraint, insufficient GPU allocation for decoding in-
 1330 stance (e.g., D/P=1) at the consumer side would limit the
 1331 full processing capability and resource utilization of prefill
 1332 instance at the producer side, however, increasing GPU al-
 1333 location for decoding instance faces the growing scaling
 1334 overhead of model parallelism (e.g., the per-GPU goodput
 1335 even degrades when increasing D/P ratio from 2 to 4). By
 1336 contrast, NEUSTREAM can minimize the mismatch between
 1337 the processing capability of P/D instances at fine-grained
 1338 SPUs abstraction and allocation.

1339 **Handle large heterogeneous cluster.** NEUSTREAM could
 1340 scale to a large-scale serving cluster by treating it as multiple
 1341 independent copies of a small cluster and scheduling each in-
 1342 dividual small cluster. NeuStream can allocate resources for
 1343 individual small clusters guided by the profiled normalized
 1344 goodput, e.g., it finds the modules in a certain small cluster
 1345 that yields the minimum normalized goodput, and dynam-
 1346 ically allocates SPU to that module (and thus the cluster).
 1347 Also, in a heterogeneous cluster, different types of GPUs can
 1348 be abstracted into different numbers of SPUs, which can still
 1349 be handled by NEUSTREAM.

1350 **Program re-writing and correctness.** Our module inher-
 1351 its the PyTorch module definition, where developers only
 1352 need to decouple the control-flow and data-flow into scatter
 1353 and compute functions, respectively. Since the control-flow
 1354 code constitutes only a small portion of the overall code, the
 1355 rewriting effort is not significant, e.g., <7% LOC re-writing
 1356 in Stable Diffusion. We currently adopt a common numerical
 1357 testing approach to assess the correctness of the rewritten
 1358 code by constructing a test set, e.g., in Stable Diffusion, the
 1359 same prompt/seeds should produce identical outputs within
 1360 a numerical error bound before/after the code changes.

1361 **Suitable scenarios.** When employing NEUSTREAM, an in-
 1362 teresting question is what are the scenarios one benefits
 1363 from using it? As illustrated in Sec. 2, NEUSTREAM is useful
 1364 for serving models with dynamic execution paths or dis-
 1365 tinct functional phases. However, serving static models with
 1366 homogeneous phases might not benefit from NEUSTREAM,
 1367 since in this scenario, the batching mechanism design is sim-
 1368 ple and straightforward. **Currently, NEUSTREAM focus on a static environment. One potential future direction is to**
 1369 **extend state management to a more dynamic context when**
 1370 **moving from one node to another, by leveraging the recent**
 1371 **work in LLM serving systems [26].**

1372
 1373
 1374
 1375

9 Conclusion

1376 NEUSTREAM addresses the inefficiencies of serving modern
 1377 DNN applications with dynamic execution patterns on GPUs.
 1378 By introducing the concepts of streams and stream-modules,
 1379 NEUSTREAM enables a continuous data flow execution model
 1380 that effectively maximizes GPU utilization while meeting
 1381 stringent latency SLOs. The abstraction of GPU resources
 1382 into Streaming Processing Units and the decomposition of
 1383 scheduling into intra-module and inter-module levels further
 1384 enhance its scalability and scheduling efficiency. Our evalua-
 1385 tion of NEUSTREAM demonstrates its superior performance
 1386 in serving complex DNN applications, such as LLM-based
 1387 chat, image generation, and multi-agents. As DNN mod-
 1388 els continue to evolve and become increasingly complex,
 1389 NEUSTREAM offers a scalable and efficient solution for future
 1390 advancements in deep learning applications.

10 Acknowledgement

1391 We thank anonymous reviewers and our shepherd, Mangpo
 1392 Phothilimthana, for their extensive suggestions. This work
 1393 was supported by National Natural Science Foundation of
 1394 China under Grant No. 92464301.

References

- [1] HuggingFace. Diffusers doc. https://huggingface.co/docs/diffusers/v0.27.0/en/stable_diffusion.
- [2] PyTorch. <https://pytorch.org>.
- [3] Amey Agrawal, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S Gulavani, and Ramachandran Ramjee. 2023. Sarathi: Efficient llm inference by piggybacking decodes with chunked prefills. *arXiv preprint arXiv:2308.16369* (2023).
- [4] Paris Carbone, Asterios Katsifodimos, Stephan Ewen, Volker Markl, Seif Haridi, and Kostas Tzoumas. 2015. Apache flink: Stream and batch processing in a single engine. *The Bulletin of the Technical Committee on Data Engineering* 38, 4 (2015).
- [5] 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 22)*. 199–216.
- [6] Weihao Cui, Zhenhua Han, Lingji Ouyang, Yichuan Wang, Ningxin Zheng, Lingxiao Ma, Yuqing Yang, Fan Yang, Jilong Xue, Lili Qiu, et al. 2023. Optimizing dynamic neural networks with brainstorm. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. 797–815.
- [7] Aditya Dhakal, Sameer G Kulkarni, and KK Ramakrishnan. 2020. Gslice: controlled spatial sharing of gpus for a scalable inference platform. In *Proceedings of the 11th ACM Symposium on Cloud Computing*. 492–506.
- [8] Arpan Gujarati, Reza Karimi, Safya Alzayat, Wei Hao, Antoine Kaufmann, Ymir Vigfusson, and Jonathan Mace. 2020. Serving DNNs like clockwork: Performance predictability from the bottom up. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 443–462.
- [9] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).

- 1431 [10] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin
1432 Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica.
1433 2023. Efficient memory management for large language model serv-
1434 ing with pagedattention. In *Proceedings of the 29th Symposium on
1435 Operating Systems Principles*. 611–626.
- 1436 [11] Zhuohan Li, Liammin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng,
1437 Xin Jin, Yanping Huang, Zifeng Chen, Hao Zhang, Joseph E Gonzalez,
1438 et al. 2023. AlpaServe: Statistical multiplexing with model parallelism
1439 for deep learning serving. In *17th USENIX Symposium on Operating
1440 Systems Design and Implementation (OSDI 23)*. 663–679.
- 1441 [12] Wei Lin, Zhengping Qian, Junwei Xu, Sen Yang, Jingren Zhou, and
1442 Lidong Zhou. 2016. {StreamScope}: Continuous Reliable Distributed
1443 Processing of Big Data Streams. In *13th USENIX Symposium on
1444 Networked Systems Design and Implementation (NSDI 16)*. 439–453.
- 1445 [13] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jea Lee. 2023. Im-
1446 proved baselines with visual instruction tuning. *arXiv preprint
1447 arXiv:2310.03744* (2023).
- 1448 [14] Amirhossein Mirhosseini, Sameh Elnikety, and Thomas F Wenisch.
1449 2021. Parslo: A gradient descent-based approach for near-optimal
1450 partial slo allotment in microservices. In *Proceedings of the ACM
1451 Symposium on Cloud Computing*. 442–457.
- 1452 [15] Derek G. Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul
1453 Barham, and Martin Abadi. 2013. Naiad: a timely dataflow system.
1454 In *Proceedings of the Twenty-Fourth ACM Symposium on Operating
1455 Systems Principles* (Farmington, Pennsylvania) (SOSP '13). Association
1456 for Computing Machinery, New York, NY, USA, 439–455. <https://doi.org/10.1145/2517349.2522738>
- 1457 [16] NVIDIA Corporation. 2024. NVIDIA Multi-Process Service. <https://docs.nvidia.com/deploy/mps/index.html> Accessed: 2024-05-21.
- 1458 [17] Pratyush Patel, Esha Choukse, Chaojie Zhang, Íñigo Goiri, Aashaka
1459 Shah, Saeed Maleki, and Ricardo Bianchini. 2023. Splitwise: Effi-
1460 cient generative llm inference using phase splitting. *arXiv preprint
1461 arXiv:2311.18677* (2023).
- 1462 [18] William Peebles and Saining Xie. 2023. Scalable diffusion mod-
1463 els with transformers. In *Proceedings of the IEEE/CVF International
1464 Conference on Computer Vision*. 4195–4205.
- 1465 [19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser,
1466 and Björn Ommer. 2022. High-resolution image synthesis with la-
1467 tent diffusion models. In *Proceedings of the IEEE/CVF conference on
1468 computer vision and pattern recognition*. 10684–10695.
- 1469 [20] Francisco Romero, Qian Li, Neeraja J Yadwadkar, and Christos
1470 Kozyrakis. 2021. {INFaaS}: Automated model-less inference serv-
1471 ing. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*.
1472 397–411.
- 1473 [21] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat,
1474 Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémie Rapin,
1475 et al. 2023. Code llama: Open foundation models for code. *arXiv
1476 preprint arXiv:2308.12950* (2023).
- 1477 [22] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan
1478 Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022.
1479 Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022
1480 conference proceedings*. 1–10.
- 1481 [23] Mohammad Shahrad, Rodrigo Fonseca, Íñigo Goiri, Gohar Chaudhry,
1482 Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark
1483 Russinovich, and Ricardo Bianchini. 2020. Serverless in the wild:
1484 Characterizing and optimizing the serverless workload at a large cloud
1485 provider. In *2020 USENIX annual technical conference (USENIX ATC
1486 20)*. 205–218.
- 1487 [24] Haichen Shen, Lequn Chen, Yuchen Jin, Liangyu Zhao, Bingyu Kong,
1488 Matthai Philipose, Arvind Krishnamurthy, and Ravi Sundaram. 2019.
1489 Nexus: A GPU cluster engine for accelerating DNN-based video
1490 analysis. In *Proceedings of the 27th ACM Symposium on Operating
1491 Systems Principles*. 322–337.
- 1492 [25] Haichen Shen, Jared Roesch, Zhi Chen, Wei Chen, Yong Wu, Mu Li, Vin
1493 Sharma, Zachary Tatlock, and Yida Wang. 2021. Nimble: Efficiently
1494 compiling dynamic neural networks for model inference. *Proceedings
1495 of Machine Learning and Systems* 3 (2021), 208–222.
- 1496 [26] Biao Sun, Ziming Huang, Hanyu Zhao, Wencong Xiao, Xinyi Zhang,
1497 Yong Li, and Wei Lin. 2024. Llumnix: Dynamic Scheduling for Large
1498 Language Model Serving. *arXiv preprint arXiv:2406.03243* (2024).
- 1499 [27] Ankit Toshniwal, Siddarth Taneja, Amit Shukla, Karthik Ramasamy,
1500 Jignesh M Patel, Sanjeev Kulkarni, Jason Jackson, Krishna Gade,
1501 Maosong Fu, Jake Donham, et al. 2014. Storm@twitter. In *Proceedings
1502 of the 2014 ACM SIGMOD international conference on Management
1503 of data*. 147–156.
- 1504 [28] Bhuvan Urgaonkar, Prashant Shenoy, Abhishek Chandra, Pawan
1505 Goyal, and Timothy Wood. 2008. Agile dynamic provisioning of
1506 multi-tier internet applications. *ACM Transactions on Autonomous
1507 and Adaptive Systems (TAAS)* 3, 1 (2008), 1–39.
- 1508 [29] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E Gonzalez.
1509 2018. Skipnet: Learning dynamic routing in convolutional networks. In
1510 *Proceedings of the European conference on computer vision (ECCV)*.
1511 409–424.
- 1512 [30] Canwen Xu and Julian McAuley. 2022. A survey on dynamic
1513 neural networks for natural language processing. *arXiv preprint
1514 arXiv:2202.07101* (2022).
- 1515 [31] Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong, Xu Han, Yukun Yan,
1516 Zhenghao Liu, Zhixing Tan, Pengyuan Liu, Dong Yu, et al. 2024. Mat-
1517 PlotAgent: Method and Evaluation for LLM-Based Agentic Scientific
1518 Data Visualization. *arXiv preprint arXiv:2402.11453* (2024).
- 1519 [32] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and
1520 Byung-Gon Chun. 2022. Orca: A distributed serving system for
1521 Transformer-Based generative models. In *16th USENIX Symposium
1522 on Operating Systems Design and Implementation (OSDI 22)*. 521–
1523 538.
- 1524 [33] Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott
1525 Shenker, and Ion Stoica. 2010. Spark: Cluster computing with working
1526 sets. In *2nd USENIX Workshop on Hot Topics in Cloud Computing
1527 (HotCloud 10)*.
- 1528 [34] Chen Zhang, Lingxiao Ma, Jilong Xue, Yining Shi, Ziming Miao,
1529 Fan Yang, Jidong Zhai, Zhi Yang, and Mao Yang. 2023. Cocktailler:
1530 Analyzing and Optimizing Dynamic Control Flow in Deep Learn-
1531 ing. In *17th USENIX Symposium on Operating Systems Design and
1532 Implementation (OSDI 23)*. USENIX Association, Boston, MA, 681–
1533 699. <https://www.usenix.org/conference/osdi23/presentation/zhangchen>
- 1534 [35] Hong Zhang, Yupeng Tang, Anurag Khandelwal, and Ion Stoica.
1535 2023. {SHEPHERD}: Serving {DNNs} in the wild. In *20th USENIX
1536 Symposium on Networked Systems Design and Implementation
1537 (NSDI 23)*. 787–808.
- 1538 [36] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya
1539 Chen, Shuhui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language
1540 models. *arXiv preprint arXiv:2205.01068* (2022).
- 1541 [37] Yanqi Zhang, Íñigo Goiri, Gohar Irfan Chaudhry, Rodrigo Fonseca,
1542 Sameh Elnikety, Christina Delimitrou, and Ricardo Bianchini. 2021.
1543 Faster and cheaper serverless computing on harvested resources.
1544 In *Proceedings of the ACM SIGOPS 28th Symposium on Operating
1545 Systems Principles*. 724–739.
- 1546 [38] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan
1547 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric P.
1548 Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024. LMSYS-
1549 Chat-1M: A Large-Scale Real-World LLM Conversation Dataset. (2024).
1550 *arXiv:2309.11998 [cs.CL]*
- 1551 [39] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xu-
1552 anzhe Liu, Xin Jin, and Hao Zhang. 2024. DistServe: Disaggregating
1553 Prefill and Decoding for Goodput-optimized Large Language Model
1554

1541 Serving. [arXiv preprint arXiv:2401.09670](#) (2024).

A Artifact Appendix

Abstract

1542 NEUSTREAM is a runtime system for serving deep learning
 1543 workloads. This artifact reproduces the main results of the
 1544 evaluation on NVIDIA H100, A6000, and RTX 4090 GPUs.

Scope

1545 This artifact will validate the following claims:

- 1546 • Performance of the Diffusion models. By reproducing
 the experiments of Figure 11, Figure 12 and Figure 13.
- 1547 • Performance of the LLM models. By reproducing the
 experiments of Figure 14, Figure 15 and Figure 16.
- 1548 • Performance of the Multi-Agent model. By reproducing
 the experiments of Figure 17, Figure 18 and Figure 19.

Contents

1549 This artifacts includes the source code to reproduce the main
 1550 experiment results shown in the paper. We provide yaml
 1551

1552 files to setup conda environments for Diffusions, LLMs and
 1553 Multi-Agent models. For each figure mentioned above, we
 1554 provide a script to reproduce its result. Since there are many
 1555 different model test cases to run to fully reproduce the results,
 1556 which will cost a long time, we also provide a portion of pre-
 1557 run logs in Diffusion models for NVIDIA RTX 4090 GPU.
 1558 Please refer to the README.pdf file in the repository for
 1559 more details.

Hosting

1560 The artifact is hosted at [github repository](#)¹. The repository
 1561 holds some model parameters so please use git lfs system,
 1562 Please use git to clone the repository.

Requirements

1563 This artifacts requires NVIDIA H100, A6000, RTX 4090 GPUs
 1564 with CUDA driver supporting CUDA runtime larger than
 1565 12.0.

1¹<https://github.com/Fjallraven-hc/NeuStream-AE>

1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619
 1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650