

NEUSTREAM: Bridging Deep Learning Serving and Stream Processing

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 [27]. 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, 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

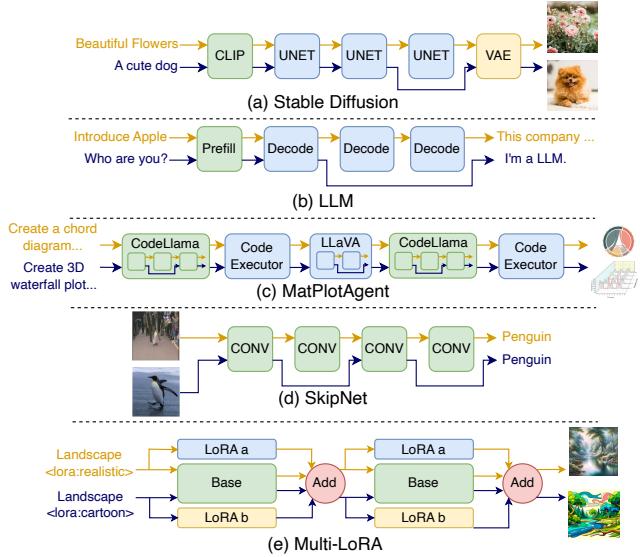
111 PyTorch, making it natural to adapt existing models. These in-
 112 terfaces implicitly require converting control flows in the na-
 113 tive language into data flow connections, i.e., streams, which
 114 allows NEUSTREAM to schedule the requests in each stream
 115 to optimize GPU utilization. Second, to execute all the het-
 116 erogeneous stream-modules on GPUs, NEUSTREAM abstracts
 117 the GPU resource into fine-grained computation units called
 118 *Streaming Processing Units (SPUs)*, where stream-modules
 119 can concurrently execute on different subsets of SPUs. Third,
 120 given the explicit stream-modules and fine-grained SPUs,
 121 NEUSTREAM decomposes the global SLO of each request into
 122 partial SLOs for each individual module. NEUSTREAM then
 123 decouples the whole scheduling space into two levels: an
 124 intra-module level scheduling decision to determine how
 125 many requests to execute in a batch for a specific module,
 126 and an inter-module level scheduling decision to determine
 127 how many resources to allocate across modules. Through
 128 such scheduling decoupling, NEUSTREAM can easily scale
 129 out to multiple GPUs, i.e., through inter-module resource
 130 allocation, and maximize serving throughput, i.e., through
 131 intra-module batch execution.

132 We implemented NEUSTREAM in 5,000 lines of Python code.
 133 NEUSTREAM’s executors for diffusion models and LLM mod-
 134 els are built on top of PyTorch and vLLM [9], respectively. We
 135 evaluated NEUSTREAM for serving mainstream DNN applica-
 136 tions, such as LLM-based chat service, image generation, and
 137 multi-agent tasks, on various NVIDIA GPUs, including the
 138 GeForce RTX 4090, Hopper H100, and RTX A6000. Our eval-
 139 uation demonstrates that NEUSTREAM achieves significantly
 140 higher serving goodput compared to state-of-the-art baseline
 141 systems, e.g., outperforming Clockwork [7] by up to 5.26×
 142 for serving Diffusion models and vLLM by 1.53×–69.31× for
 143 serving OPT models at various request rates, respectively.
 144 Looking forward, the rapid advancements in DNN model
 145 capabilities could enable more complex execution flows in fu-
 146 ture applications. We believe NEUSTREAM provides a scalable,
 147 efficient, and hardware-friendly system design to accommo-
 148 date these ever-increasing applications.

150 2 Motivation

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

153 **Challenge I: Dynamic execution paths.** Most of the ex-
 154 isting serving systems consider *static* DNN models whose
 155 execution paths are fixed, i.e., independent of the input. *Dy-
 156 namic* DNN models, in contrast, can adjust their computa-
 157 tional path with respect to input, and such models are now
 158 becoming prevalent. As illustrated in Figure 1, diffusion [18]
 159 models could work at different numbers of denoising steps,
 160 depending on preferences on generation quality and infer-
 161 ence speed of each input. The LLM [34] generates output
 162 tokens autoregressively on a per-input basis, where each
 163 request may require a different number of decode iterations.
 164



166 **Figure 1.** Examples of dynamic neural networks. They dy-
 167 namically adjust computation paths on per-input basis.
 168

169 The LLM-based multi-agents [28] further exacerbate this
 170 dynamics properties. Skipping and early stopping models
 171 such as Skipnet [26] dynamically decide whether to skip
 172 certain layers based on the input. Multiple Low-Rank Adap-
 173 tation (LoRA [8]) requests share the same base model layer
 174 but enter into different LoRA adapter layers. Dynamic per-
 175 input execution paths in these models pose challenges for
 176 the optimizations that have been key to obtaining high GPU
 177 efficiency. For example, batching cannot be applied if the
 178 execution paths of inputs are not identical. Also, one can-
 179 not *proactively* perform scheduling actions (e.g., when or
 180 whether to execute a request), because dynamic execution
 181 paths have unpredictable latency performance.

182 **Challenge II: Distinct functional phases.** There are dif-
 183 ferent main functional phases during a dynamic model in-
 184 ference. For example, stable diffusion inference involves text
 185 encoding, diffusion and decoding phases, and LLM inference
 186 contains prefilling and decoding phases. We find that differ-
 187 ent phases have distinct batching characteristics. Figure 2
 188 and Figure 3 show the request execution latency on each
 189 phase, which is batch’s latency / batch size (i.e., excluding
 190 queuing time). On a certain phase, as the batch size increases,
 191 the decreasing slope of a curve implies that batching leads
 192 to only a marginal gain in improving the execution latency
 193 (and thus the throughput). Once the gain becomes insignif-
 194 icant, adding more requests to the batch no longer improves
 195 the throughput. Instead, it proportionally extends the total
 196 processing time for the batch, inadvertently delaying all in-
 197 cluded requests. Thus, it is necessary to identify a critical
 198 threshold for batch size in each phase, beyond which the
 199 marginal gain becomes insignificant. Given that the slopes of
 200 different phases are distinct, we should identify the threshold
 201

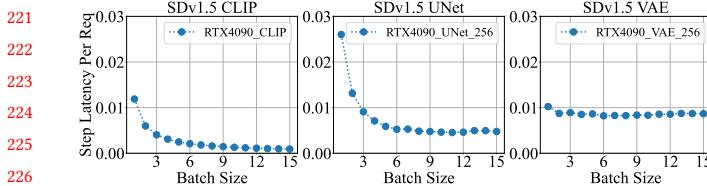


Figure 2. Diffusion phase execution latency on RTX 4090.

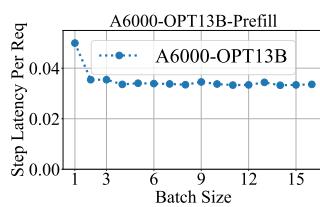


Figure 3. LLM phase execution latency on RTX A6000.

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[8] adapters enables batching in the shared base model by decomposing the computation between the base model and the LoRA adapters. 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.

```

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.

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.

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

```

331 1
332 2
333 3
334 4
335 5
336 6
337 7
338 8
339 9
340 10
341 11
342 12
343 13
344 14
345 15
346 16
347 17
348 18
349 19
350 20
351 21
352 22
353 23
354 24
355 25
356 26
357 27
358 28
359 29
360 30
361 31
362 32
363 33
364 34
365 35
366 36
367 37
368 38
369 39
370 40
371 41
372 42
373 43
374 44
375 45
376 46
377 47
378 48
379 49
380 50
381 51
382 52
383 53
384 54
385 55
386 56
387 57
388 58
389 59
390 60
391 61
392 62
393 63
394 64
395 65
396 66
397 67
398 68
399 69
400 70
401 71
402 72
403 73
404 74
405 75
406 76
407 77
408 78
409 79
410 80
411 81
412 82
413 83
414 84
415 85
416 86
417 87
418 88
419 89
420 90
421 91
422 92
423 93
424 94
425 95
426 96
427 97
428 98
429 99
430 100
431 101
432 102
433 103
434 104
435 105
436 106
437 107
438 108
439 109
440 110
441 111
442 112
443 113
444 114
445 115
446 116
447 117
448 118
449 119
450 120
451 121
452 122
453 123
454 124
455 125
456 126
457 127
458 128
459 129
460 130
461 131
462 132
463 133
464 134
465 135
466 136
467 137
468 138
469 139
470 140
471 141
472 142
473 143
474 144
475 145
476 146
477 147
478 148
479 149
480 150
481 151
482 152
483 153
484 154
485 155
486 156
487 157
488 158
489 159
490 160
491 161
492 162
493 163
494 164
495 165
496 166
497 167
498 168
499 169
500 170
501 171
502 172
503 173
504 174
505 175
506 176
507 177
508 178
509 179
510 180
511 181
512 182
513 183
514 184
515 185
516 186
517 187
518 188
519 189
520 190
521 191
522 192
523 193
524 194
525 195
526 196
527 197
528 198
529 199
530 200
531 201
532 202
533 203
534 204
535 205
536 206
537 207
538 208
539 209
540 210
541 211
542 212
543 213
544 214
545 215
546 216
547 217
548 218
549 219
550 220
551 221
552 222
553 223
554 224
555 225
556 226
557 227
558 228
559 229
560 230
561 231
562 232
563 233
564 234
565 235
566 236
567 237
568 238
569 239
570 240
571 241
572 242
573 243
574 244
575 245
576 246
577 247
578 248
579 249
580 250
581 251
582 252
583 253
584 254
585 255
586 256
587 257
588 258
589 259
590 260
591 261
592 262
593 263
594 264
595 265
596 266
597 267
598 268
599 269
600 270
601 271
602 272
603 273
604 274
605 275
606 276
607 277
608 278
609 279
610 280
611 281
612 282
613 283
614 284
615 285
616 286
617 287
618 288
619 289
620 290
621 291
622 292
623 293
624 294
625 295
626 296
627 297
628 298
629 299
630 300
631 301
632 302
633 303
634 304
635 305
636 306
637 307
638 308
639 309
640 310
641 311
642 312
643 313
644 314
645 315
646 316
647 317
648 318
649 319
650 320
651 321
652 322
653 323
654 324
655 325
656 326
657 327
658 328
659 329
660 330
661 331
662 332
663 333
664 334
665 335
666 336
667 337
668 338
669 339
670 340
671 341
672 342
673 343
674 344
675 345
676 346
677 347
678 348
679 349
680 350
681 351
682 352
683 353
684 354
685 355
686 356
687 357
688 358
689 359
690 360
691 361
692 362
693 363
694 364
695 365
696 366
697 367
698 368
699 369
700 370
701 371
702 372
703 373
704 374
705 375
706 376
707 377
708 378
709 379
710 380
711 381
712 382
713 383
714 384
715 385
716 386
717 387
718 388
719 389
720 390
721 391
722 392
723 393
724 394
725 395
726 396
727 397
728 398
729 399
730 400
731 401
732 402
733 403
734 404
735 405
736 406
737 407
738 408
739 409
740 410
741 411
742 412
743 413
744 414
745 415
746 416
747 417
748 418
749 419
750 420
751 421
752 422
753 423
754 424
755 425
756 426
757 427
758 428
759 429
760 430
761 431
762 432
763 433
764 434
765 435
766 436
767 437
768 438
769 439
770 440
771 441
772 442
773 443
774 444
775 445
776 446
777 447
778 448
779 449
780 450
781 451
782 452
783 453
784 454
785 455
786 456
787 457
788 458
789 459
790 460
791 461
792 462
793 463
794 464
795 465
796 466
797 467
798 468
799 469
800 470
801 471
802 472
803 473
804 474
805 475
806 476
807 477
808 478
809 479
810 480
811 481
812 482
813 483
814 484
815 485
816 486
817 487
818 488
819 489
820 490
821 491
822 492
823 493
824 494
825 495
826 496
827 497
828 498
829 499
830 500
831 501
832 502
833 503
834 504
835 505
836 506
837 507
838 508
839 509
840 510
841 511
842 512
843 513
844 514
845 515
846 516
847 517
848 518
849 519
850 520
851 521
852 522
853 523
854 524
855 525
856 526
857 527
858 528
859 529
860 530
861 531
862 532
863 533
864 534
865 535
866 536
867 537
868 538
869 539
870 540
871 541
872 542
873 543
874 544
875 545
876 546
877 547
878 548
879 549
880 550
881 551
882 552
883 553
884 554
885 555
886 556
887 557
888 558
889 559
890 560
891 561
892 562
893 563
894 564
895 565
896 566
897 567
898 568
899 569
900 570
901 571
902 572
903 573
904 574
905 575
906 576
907 577
908 578
909 579
910 580
911 581
912 582
913 583
914 584
915 585
916 586
917 587
918 588
919 589
920 590
921 591
922 592
923 593
924 594
925 595
926 596
927 597
928 598
929 599
930 600
931 601
932 602
933 603
934 604
935 605
936 606
937 607
938 608
939 609
940 610
941 611
942 612
943 613
944 614
945 615
946 616
947 617
948 618
949 619
950 620
951 621
952 622
953 623
954 624
955 625
956 626
957 627
958 628
959 629
960 630
961 631
962 632
963 633
964 634
965 635
966 636
967 637
968 638
969 639
970 640
971 641
972 642
973 643
974 644
975 645
976 646
977 647
978 648
979 649
980 650
981 651
982 652
983 653
984 654
985 655
986 656
987 657
988 658
989 659
990 660
991 661
992 662
993 663
994 664
995 665
996 666
997 667
998 668
999 669
1000 670
1001 671
1002 672
1003 673
1004 674
1005 675
1006 676
1007 677
1008 678
1009 679
1010 680
1011 681
1012 682
1013 683
1014 684
1015 685
1016 686
1017 687
1018 688
1019 689
1020 690
1021 691
1022 692
1023 693
1024 694
1025 695
1026 696
1027 697
1028 698
1029 699
1030 700
1031 701
1032 702
1033 703
1034 704
1035 705
1036 706
1037 707
1038 708
1039 709
1040 710
1041 711
1042 712
1043 713
1044 714
1045 715
1046 716
1047 717
1048 718
1049 719
1050 720
1051 721
1052 722
1053 723
1054 724
1055 725
1056 726
1057 727
1058 728
1059 729
1060 730
1061 731
1062 732
1063 733
1064 734
1065 735
1066 736
1067 737
1068 738
1069 739
1070 740
1071 741
1072 742
1073 743
1074 744
1075 745
1076 746
1077 747
1078 748
1079 749
1080 750
1081 751
1082 752
1083 753
1084 754
1085 755
1086 756
1087 757
1088 758
1089 759
1090 760
1091 761
1092 762
1093 763
1094 764
1095 765
1096 766
1097 767
1098 768
1099 769
1100 770
1101 771
1102 772
1103 773
1104 774
1105 775
1106 776
1107 777
1108 778
1109 779
1110 780
1111 781
1112 782
1113 783
1114 784
1115 785
1116 786
1117 787
1118 788
1119 789
1120 790
1121 791
1122 792
1123 793
1124 794
1125 795
1126 796
1127 797
1128 798
1129 799
1130 800
1131 801
1132 802
1133 803
1134 804
1135 805
1136 806
1137 807
1138 808
1139 809
1140 810
1141 811
1142 812
1143 813
1144 814
1145 815
1146 816
1147 817
1148 818
1149 819
1150 820
1151 821
1152 822
1153 823
1154 824
1155 825
1156 826
1157 827
1158 828
1159 829
1160 830
1161 831
1162 832
1163 833
1164 834
1165 835
1166 836
1167 837
1168 838
1169 839
1170 840
1171 841
1172 842
1173 843
1174 844
1175 845
1176 846
1177 847
1178 848
1179 849
1180 850
1181 851
1182 852
1183 853
1184 854
1185 855
1186 856
1187 857
1188 858
1189 859
1190 860
1191 861
1192 862
1193 863
1194 864
1195 865
1196 866
1197 867
1198 868
1199 869
1200 870
1201 871
1202 872
1203 873
1204 874
1205 875
1206 876
1207 877
1208 878
1209 879
1210 880
1211 881
1212 882
1213 883
1214 884
1215 885
1216 886
1217 887
1218 888
1219 889
1220 890
1221 891
1222 892
1223 893
1224 894
1225 895
1226 896
1227 897
1228 898
1229 899
1230 900
1231 901
1232 902
1233 903
1234 904
1235 905
1236 906
1237 907
1238 908
1239 909
1240 910
1241 911
1242 912
1243 913
1244 914
1245 915
1246 916
1247 917
1248 918
1249 919
1250 920
1251 921
1252 922
1253 923
1254 924
1255 925
1256 926
1257 927
1258 928
1259 929
1260 930
1261 931
1262 932
1263 933
1264 934
1265 935
1266 936
1267 937
1268 938
1269 939
1270 940
1271 941
1272 942
1273 943
1274 944
1275 945
1276 946
1277 947
1278 948
1279 949
1280 950
1281 951
1282 952
1283 953
1284 954
1285 955
1286 956
1287 957
1288 958
1289 959
1290 960
1291 961
1292 962
1293 963
1294 964
1295 965
1296 966
1297 967
1298 968
1299 969
1300 970
1301 971
1302 972
1303 973
1304 974
1305 975
1306 976
1307 977
1308 978
1309 979
1310 980
1311 981
1312 982
1313 983
1314 984
1315 985
1316 986
1317 987
1318 988
1319 989
1320 990
1321 991
1322 992
1323 993
1324 994
1325 995
1326 996
1327 997
1328 998
1329 999
1330 1000
1331 1001
1332 1002
1333 1003
1334 1004
1335 1005
1336 1006
1337 1007
1338 1008
1339 1009
1340 1010
1341 1011
1342 1012
1343 1013
1344 1014
1345 1015
1346 1016
1347 1017
1348 1018
1349 1019
1350 1020
1351 1021
1352 1022
1353 1023
1354 1024
1355 1025
1356 1026
1357 1027
1358 1028
1359 1029
1360 1030
1361 1031
1362 1032
1363 1033
1364 1034
1365 1035
1366 1036
1367 1037
1368 1038
1369 1039
1370 1040
1371 1041
1372 1042
1373 1043
1374 1044
1375 1045
1376 1046
1377 1047
1378 1048
1379 1049
1380 1050
1381 1051
1382 1052
1383 1053
1384 1054
1385 1055
1386 1056
1387 1057
1388 1058
1389 1059
1390 1060
1391 1061
1392 1062
1393 1063
1394 1064
1395 1065
1396 1066
1397 1067
1398 1068
1399 1069
1400 1070
1401 1071
1402 1072
1403 1073
1404 1074
1405 1075
1406 1076
1407 1077
1408 1078
1409 1079
1410 1080
1411 1081
1412 1082
1413 1083
1414 1084
1415 1085
1416 1086
1417 1087
1418 1088
1419 1089
1420 1090
1421 1091
1422 1092
1423 1093
1424 1094
1425 1095
1426 1096
1427 1097
1428 1098
1429 1099
1430 1100
1431 1101
1432 1102
1433 1103
1434 1104
1435 1105
1436 1106
1437 1107
1438 1108
1439 1109
1440 1110
1441 1111
1442 1112
1443 1113
1444 1114
1445 1115
1446 1116
1447 1117
1448 1118
1449 1119
1450 1120
1451 1121
1452 1122
1453 1123
1454 1124
1455 1125
1456 1126
1457 1127
1458 1128
1459 1129
1460 1130
1461 1131
1462 1132
1463 1133
1464 1134
1465 1135
1466 1136
1467 1137
1468 1138
1469 1139
1470 1140
1471 1141
1472 1142
1473 1143
1474 1144
1475 1145
1476 1146
1477 1147
1478 1148
1479 1149
1480 1150
1481 1151
1482 1152
1483 1153
1484 1154
1485 1155
1486 1156
1487 1157
1488 1158
1489 1159
1490 1160
1491 1161
1492 1162
1493 1163
1494 1164
1495 1165
1496 1166
1497 1167
1498 1168
1499 1169
1500 1170
1501 1171
1502 1172
1503 1173
1504 1174
1505 1175
1506 1176
1507 1177
1508 1178
1509 1179
1510 1180
1511 1181
1512 1182
1513 1183
1514 1184
1515 1185
1516 1186
1517 1187
1518 1188
1519 1189
1520 1190
1521 1191
1522 1192
1523 1193
1524 1194
1525 1195
1526 1196
1527 1197
1528 1198
1529 1199
1530 1200
1531 1201
1532 1202
1533 1203
1534 1204
1535 1205
1536 1206
1537 1207
1538 1208
1539 1209
1540 1210
1541 1211
1542 1212
1543 1213
1544 1214
1545 1215
1546 1216
1547 1217
1548 1218
1549 1219
1550 1220
1551 1221
1552 1222
1553 1223
1554 1224
1555 1225
1556 1226
1557 1227
1558 1228
1559 1229
1560 1230
1561 1231
1562 1232
1563 1233
1564 1234
1565 1235
1566 1236
1567 1237
1568 1238
1569 1239
1570 1240
1571 1241
1572 1242
1573 1243
1574 1244
1575 1245
1576 1246
1577 1247
1578 1248
1579 1249
1580 1250
1581 1251
1582 1252
1583 1253
1584 1254
1585 1255
1586 1256
1587 1257
1588 1258
1589 1259
1590 1260
1591 1261
1592 1262
1593 1263
1594 1264
1595 1265
1596 1266
1597 1267
1598 1268
1599 1269
1600 1270
1601 1271
1602 1272
1603 1273
1604 1274
1605 1275
1606 1276
1607 1277
1608 1278
1609 1279
1610 1280
1611 1281
1612 1282
1613 1283
1614 1284
1615 1285
1616 1286
1617 1287
1618 1288
1619 1289
1620 1290
1621 1291
1622 1292
1623 1293
1624 1294
1625 1295
1626 1296
1627 1297
1628 1298
1629 1299
1630 1300
1631 1301
1632 1302
1633 1303
1634 1304
1635 1305
1636 1306
1637 1307
1638 1308
1639 1309
```

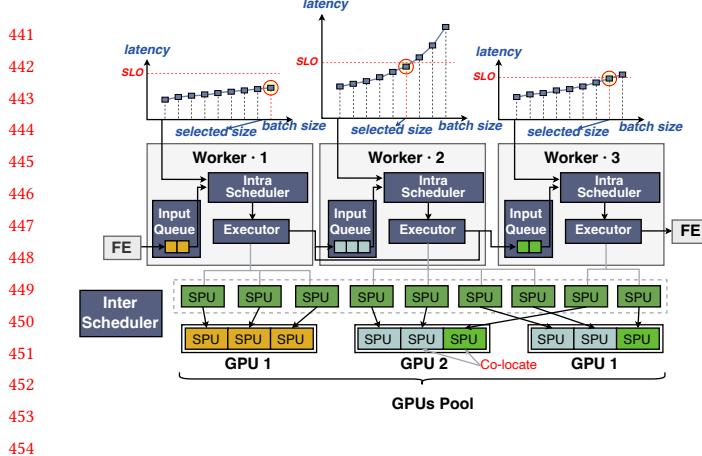


Figure 8. NEUSTREAM’s two-level scheduling framework.

in a single stream have the same static operation graph defined in compute. Second, modules are connected via streams with the continuous flow of data, instead of static tensors consuming exactly once. This supports a decomposed and continuous stream computation model, extended with the ability to customize execution at each module (e.g., batching and resource allocation).

4 NEUSTREAM Scheduling

The goal for NEUSTREAM is to maximize the system goodput, i.e., the number of served requests that meet their SLO requirements per unit time. The hierarchical streaming graph naturally leads to a two-level scheduling framework: *intra*-module scheduling on how many requests to include in a batch for a specific module, and *inter*-module scheduling on how much resource to be allocated across modules. Figure 8 illustrates the architecture of NEUSTREAM, which is composed of a frontend and multiple stream workers. Each worker is responsible for executing a single module. The frontend continuously sends incoming requests to the first stream worker. Upon receiving request messages, the executor in a worker performs the gather-compute-scatter execution model, sending new messages to downstream workers based on stream connection. The batching decision is made by the intra-module scheduler at each worker, utilizing partial SLO and profiled latency function of each module. The resource allocation decision is made by the inter-module scheduler, which aims to balance the maximum goodput permitted by each module under an appropriate resource allocation.

4.1 Intra-module Scheduling

Partial SLOs. In the presence of control flow, it is unclear how to define end-to-end latency SLOs for individual requests with dynamic and even unpredictable execution paths. Since NEUSTREAM decouples routing control-flow from execution, we take the approach of assigning *Partial SLOs* to individual modules, which is similar to the partial SLO concept in microservices [13, 25]. Specifically, given the module

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

```

Input: Latency function  $L_i(b, a_i)$  of batch size  $b$  and
        resource allocation  $a_i$ , partial  $SLO_i$ , input stream  $in$ 
Output: batch_request
1 Function IntraScheduler( $L_i, SLO_i, in$ ):
  2   batch_request  $\leftarrow [ ]$ ;
    // get maximum batch size based on partial SLO
  3   batch_limit  $\leftarrow \text{argmax}_b L_i(b, a_i) \leq SLO_i$ ;
  4   for request  $r$  in  $in$  do
  5     if  $r.budget \leq 0$  then
      // drop request that fail to meet upstream
      // partial SLOs
      in.remove( $r$ );
      continue;
    end
    if length(batch_request) < batch_limit then
      | batch_request.append( $r$ );
    end
  end
  return batch_request;

```

execution path $\varphi = \{m_1, \dots, m_k\}$ of a certain request, we consider its end-to-end latency $SLO = \sum_{i \in \{1, \dots, |\varphi|\}} SLO_{\varphi(i)}$, where $SLO_{\varphi(i)}$ is the partial SLO for module $\varphi(i)$.

Batch scheduling algorithm. As long as all modules meet their respective partial SLOs, the request SLO is guaranteed. This approach ensures independent batch scheduling across individual modules. Algorithm 1 describes the procedure of intra-module batch scheduling for one execution step on a specific module m_i . The algorithm receives the following inputs for m_i : (1) $L_i(b, a)$: profiled execution latency of batch size b and resource allocation a (detailed in Section 4.2), which includes the duration of executing compute and scatter, (2) SLO_i : partial SLO (in latency), and (3) in : the input stream. For every step, if the stream in is non-empty, the scheduler determines a batch: we first obtain the maximum permitted batch size $batch_limit_i = \text{argmax}_b L_i(b, a) \leq SLO_i$ without violating the partial-SLO of this module (line 3). Then for each request in the input stream, we add a request into the batch if it satisfies:

C1. The request meets the partial SLOs on the upstream modules. In particular, let $\varphi_r = \{m_1, \dots, m_k\}$ be the execution path of previous steps of a request r before executing the module m_i at the current step, and t and t_0 be the current time and request issued time, respectively. We require the request budget $B_r = \sum_{j \in \{1, \dots, |\varphi_r|\}} SLO_{\varphi_r(j)} - (t - t_0) \geq 0$. In line with other SLO-aware scheduling approaches [7, 10, 32] that drop requests unlikely to meet their SLOs, we remove requests with negative budget from the stream to avoid continuously blocking the subsequent incoming requests (line 5-8). However, in more sophisticated systems, these removed

551 requests may not be dropped entirely; instead, they can be
 552 placed in a low-priority queue or escalated to an upper-level
 553 scheduler for re-dispatch.

554 **C2.** The current batch size is smaller than the maximum
 555 permitted batch size *batch_limit* (line 9-11).

556 The above two conditions together achieve a large good-
 557 put at a specific module m_i by ensuring that the requests in
 558 the batch can still meet partial SLOs after the execution of
 559 this module, and the GPU operates at high efficiency. In the
 560 stream of each module, we sort the request messages based
 561 on their arrival time at the frontend. One could also prioritize
 562 requests with tight SLOs based on the request budget.
 563

564 4.2 Inter-module Scheduling

565 **Streaming Processing Units (SPUs).** To facilitate resource
 566 allocation, we introduce an abstraction of SPUs, which can
 567 be considered as a uniform partition of physical devices, as il-
 568 lustrated in Figure 8. Given a total of N SPUs and M modules,
 569 the inter-scheduler determines the number of SPUs a_i for
 570 each module m_i . Let k_i be the average number of execution
 571 times of a specific module m_i in the execution path. Given
 572 the batch scheduling of Algorithm 2 and resource allocation
 573 a_i , we define the normalized goodput $g_i(a_i, SLO_i)$ of m_i as:

$$575 \quad g_i(a_i, SLO_i) = \frac{batch_limit_i}{SLO_i \times k_i} = \frac{argmax_b L_i(b, a_i) \leq SLO_i}{SLO_i \times k_i} \quad (1)$$

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

$$578 \quad \begin{aligned} & \max \min_{i \in \{1, \dots, M\}} g_i(a_i, SLO_i) \\ & \text{s.t. } \sum_i a_i \leq N \quad \forall i \in 1, \dots, M \\ & \quad a_i \cdot mem \geq Mem_i \quad \forall i \in 1, \dots, M \quad (\text{Memory limit}) \end{aligned} \quad (2)$$

585 where mem is the memory capacity of a SPU and Mem_i is
 586 the memory footprint of module m_i .

588 **Allocation algorithm.** Algorithm 2 presents pseudocode
 589 for our iterative SPU allocation process. The algorithm first
 590 initializes all modules by allocating the minimum number of
 591 SPUs a_i that satisfies the memory limit condition for each
 592 module m_i (line 2). This step ensures that each module can
 593 be put on the SPUs under the memory constraint. Then, the
 594 algorithm finds the modules that yield the minimum normal-
 595 ized goodput, and allocates Δ ($\Delta \geq 1$) SPU to that module
 596 (line 12-14). To this end, for each module, the algorithm maps
 597 its allocated SPUs onto physical devices (e.g., GPUs) through
 598 module placement detailed below (line 4-9), and gets the
 599 profiled latency function $L_i(b, a_i)$ of the module (line 10).
 600 Based on the latency functions and equation (1), the algo-
 601 rithm computes normalized goodput g_i of each module m_i .
 602 The algorithm continues allocating Δ SPUs at a time until
 603 all the SPUs are allocated.

606 Algorithm 2: Inter-moduler scheduler for resource 607 allocation across modules

```
608 Input: stream graph G, function partial SLO for each
  609 module, total number of SPUs  $N$ , devices  $D$ 
  610 Output: [SPU_allot, SPU_map]
  611 Function InterScheduler( $G, SLO, N$ ):
  612   // assign initial SPUs satisfying memory limit
  613   SPU_allot  $\leftarrow$  Initialize( $G, SLO$ );
  614   while  $N \geq \Delta$  do
  615     // map SPU allotment to physical devices
  616     sort(SPU_allot);
  617     for  $i \in \{1, \dots, M\}$  do
  618       SPU_map[i] = best_fit(SPU_allot[i], D);
  619       if SPU_map[i] =  $\emptyset$  then
  620         | SPU_map[i] = partition(SPU_allot[i], D);
  621       end
  622        $L_i(b, a_i) = get\_profile\_latency(SPU\_map[i], D)$ ;
  623     end
  624     // assign SPUs to bottleneck module
  625      $m = argmin_{i \in \{1, \dots, M\}} g_i(SPU\_allot[i], SLO_i)$ ;
  626     SPU_allot(m).increase( $\Delta$ );
  627      $N = N - \Delta$ ;
  628   end
  629   return [SPU_allot, SPU_map];
  630
```

580 **Module placement.** When mapping SPUs to physical de-
 581 vices (line 4-9 in Algorithm 2), the scheduler considers mod-
 582 ules in descending order of SPU allocation, maps SPUs of
 583 a module onto a device that has the least available SPUs to
 584 accommodate the module (line 5). If no single device has
 585 sufficient available resources to accommodate a module, we
 586 uniformly map SPUs of a module onto the fewest possible
 587 devices to reduce communication costs (line 7-9). In this
 588 case, we rely on model-parallel partition tools to partition a
 589 module (operation graph) across multiple devices. We can
 590 enumerate possible combinations of data parallelism and
 591 model parallelism, and use the profiled latency function to
 592 identify the configuration that yields the best module good-
 593 put, i.e., $\max argmax_b L_i(b, a_i) \leq SLO_i$.

594 **Colocating modules.** The scheduler allows the two mod-
 595 ules to be annotated as *co-located*, and maps the SPUs of two
 596 co-located modules on the same set of devices to reduce the
 597 communication and increase the data locality. Given a work-
 598 load, NEUSTREAM could automatically search for the mod-
 599 ules to be co-located: For a pair of originally non-collocated
 600 modules (m_i, m_j) , the scheduler re-maps SPUs by annotating
 601 them as *co-located*, and treats the annotation valid only if
 602 co-locating m_i and m_j could improve their overall normal-
 603 ized goodput, i.e., $\min\{g_i, g_j\}$, based on the profiled latency
 604 functions and equation (1).

605 5 Implementation

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

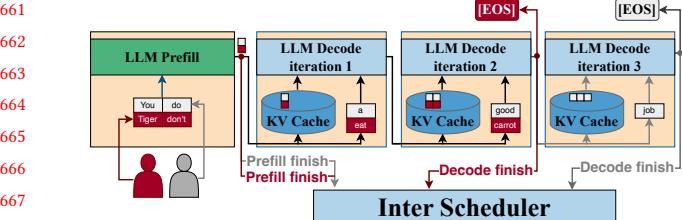


Figure 9. State management for LLM.

Module execution. We built the executor of diffusion models and LLM models on top of PyTorch [2] and vLLM[9], respectively. For diffusion models, we directly use PyTorch for sub-module parallelism across multiple cores on a single GPU. For larger LLMs, we use vLLM to provide sub-module parallelism across multiple GPUs. Each module is executed in its own process by default, however, co-located modules can be executed in the same process to optimize module switching and communication overhead (e.g., only passing pointers to the output data that's already on the device). The module invocation number depends on the workload, which can be determined statically before execution in diffusion or dynamically during execution in LLMs with the EOS token being detected. To compute the normalized goodput, we could use the average invocation number over the requests that have already been served.

State management. Each module can optionally manage the state for each request that can be used across different messages. The state of a module is a set of key/value pairs, which can be stored in a key-value store and operated by a module executor via get, put, and delete operations. NeuStream's inter-module scheduler tracks the *progress* of a 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 [9], 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.

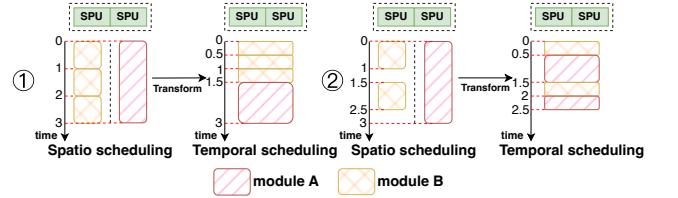


Figure 10. Spatial-to-temporal scheduling transformation with EFTF algorithm in a single device.

Streaming Processing Units. We could leverage spatial partitioning such as Multi-process-service (MPS[15]) 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 [3, 24] 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.

771	Name	Type	Task	Year
772	StableDiffusion [18]	Diffusion	Text to Image	2022
773	Palette [21]	Diffusion	Image Restoration	2022
774	DiT_XL/2 [17]	Diffusion	Image Generation	2023
775	DiT_S/2	Diffusion	Image Generation	2023
776	OPT-6.7B [34]	LLM	Text Generation	2022
777	OPT-13B	LLM	Text Generation	2022
778	OPT-30B	LLM	Text Generation	2022
779	OPT-66B	LLM	Text Generation	2022
780	MatPlotAgent [28]	Multi-agent	Multi-agent Assistant	2024

Table 1. Experiments setting.

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 [18]), image restoration (Palette [21]), and class-conditional image generation (DiT [17]). For LLM model, we select the OPT [34]. We choose several model sizes to test performance under different numbers of devices. For multi-agent, we use MatPlotAgent [28] designed to automate scientific data visualization tasks. The model is comprised of three LLMs: a Code Llama [20] to generate preliminary code, LLaVA [12] 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 server with AMD EPYC 7513 and 4 RTX A6000 (48GB) GPUs, running on Ubuntu 18.04 with cuda 12.2.

Offered load. Similar to prior work [10], 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 [35], a large-scale dataset containing real-world conversations with 25 state-of-the-art LLMs. For multi-agent, we use MatPlotBench [28] 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. Clockwork was originally implemented with a legacy TVM backend that has limited operator coverage. For a fair comparison of the scheduling design, similar to prior work [10, 33], we re-implemented Clockwork using the latest PyTorch backend. For LLMs, we select vLLM [9] as a strong baseline, which supports iteration-level scheduling [29] and PagedAttention [9]. We adopt the latest vLLM version (v0.5.4) released in August 2024.

6.2 Results on Diffusion Models

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

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

The computational cost of StableDiffusion depends on the input image size. For large 512×512 images, as shown in Figure 12, it almost saturates the RTX GPU even with a single request, limiting the batching gain of NEUSTREAM. However, we observe that NEUSTREAM can effectively scale up to a more powerful H100 GPU (up to 1.90× as compared to Clockwork across offered loads).

Batching efficiency. To fully understand NEUSTREAM’s benefits on batching, we plot the running batch size on DiT module for DiT-S/2 along the time in Figure 13(a). On average,

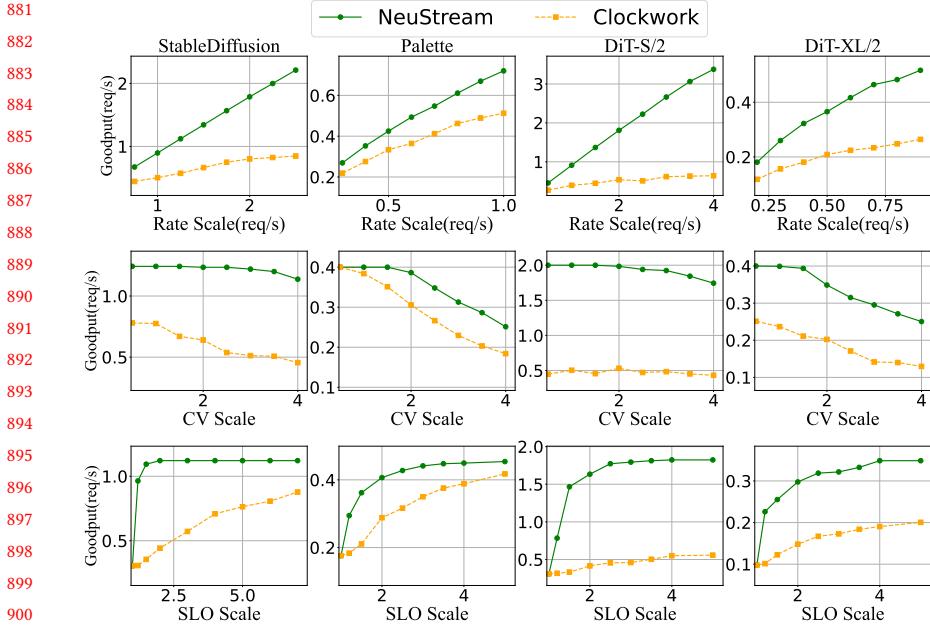


Figure 11. Experiments of generating 256×256 images on NVIDIA RTX4090.

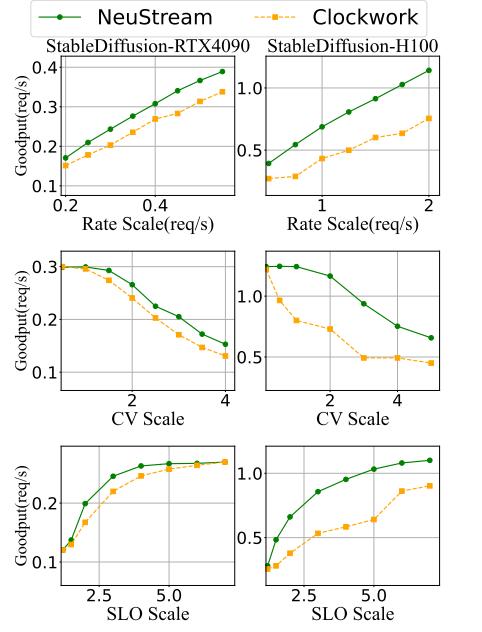


Figure 12. Generating 512×512 images on RTX4090 and H100, respectively.

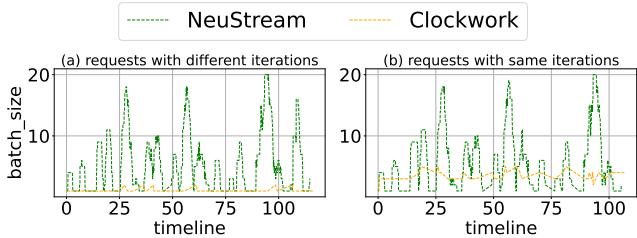


Figure 13. Batch size of DiT module along the time, where requests start with different/same number of iterations.

the batch size of NEUSTREAM is $5.35 \times$ larger than Clockwork due to efficiently batching requests that (1) arrive at different times; (2) require different number of iterations to finish; or (3) start with different number of iterations. Figure 13(b) increases the batching opportunity of Clockwork by setting the same number of iterations for each request, i.e., disabling the above case (3). We see that NEUSTREAM still has a larger batch size ($1.67 \times$ on average larger).

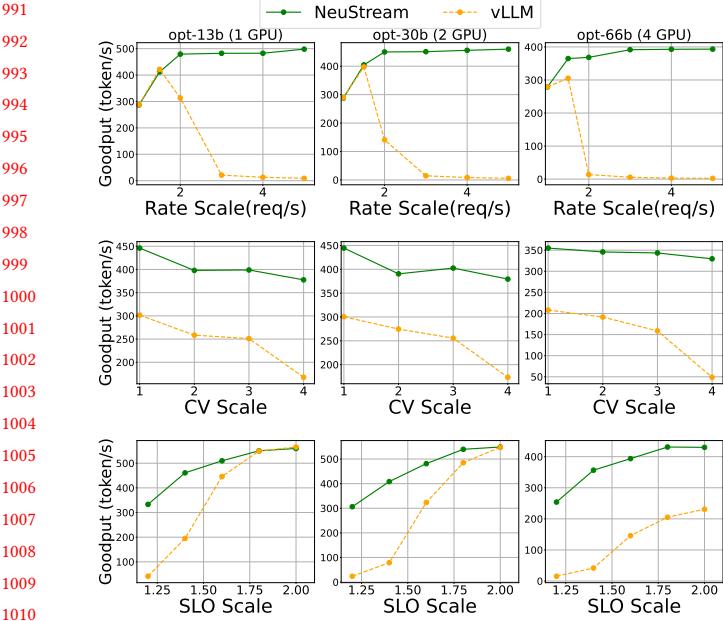
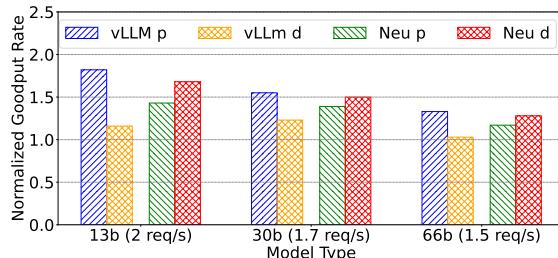
6.3 Results on LLM Models

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

Goodput under varying offered load. Figure 14 shows the results of serving OPT models on RTX A6000 GPUs. The

1st row shows that, as the request load increases, vLLM goodput starts to decrease significantly. This is because vLLM greedily prioritizes prefill requests and the use of preemption leads to resource starvation of decode module. In contrast, NEUSTREAM is SLO-aware and balances the resource provision across the modules. For OPT-13B, NEUSTREAM outperforms vLLM by $1.53 \times$ at 2 requests/s. Under an even higher rate of 4 requests/s, NEUSTREAM achieves $37.21 \times$ more goodput than vLLM, due to decode resource starvation and prefill congestion (caused by insufficient memory for KV cache) in vLLM. As shown in Figure 14's 2nd and 3rd rows, NEUSTREAM performs better than vLLM to handle more burstiness (a higher CV) and stringent SLOs (a smaller SLO scale) that aggravates the resource contention, e.g., $2.25 \times$ at the CV of 4 and $7.98 \times$ at SLO scale of 1.2 on OPT-13B. Figure 14's 2nd and 3rd columns show NEUSTREAM consistently outperforms vLLM when scaling the serving to accommodate larger models on multiple devices (e.g., NEUSTREAM achieves up to $69.31 \times$ at 3 requests/s when serving the largest OPT-66B model on 4 devices).

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

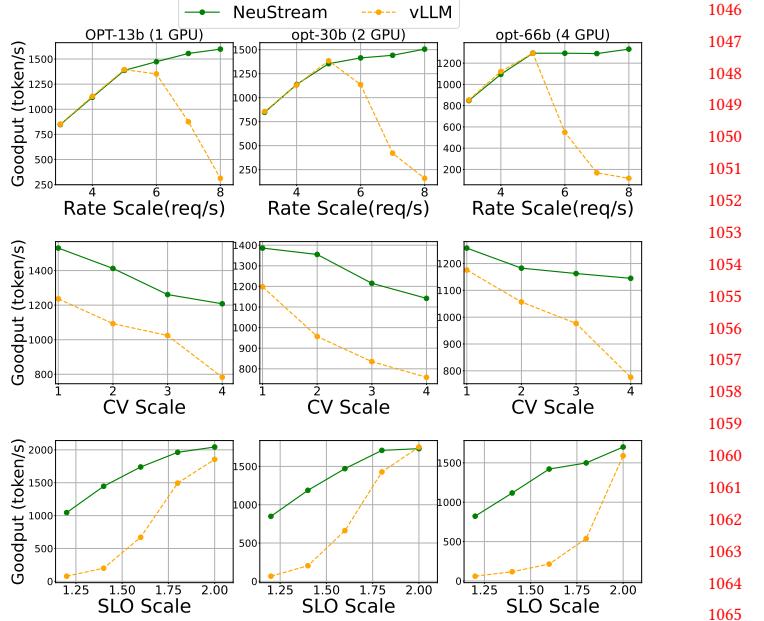
**Figure 14.** Goodput of LLM decode module on A6000.**Figure 16.** Module's normalized goodputs of each LLM.

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

6.4 Results on Multi-agent Model

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

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×

**Figure 15.** Goodput of LLM decode module on H100.

	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%
Memory	3.62%	3.39%	4.79%	

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

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.

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

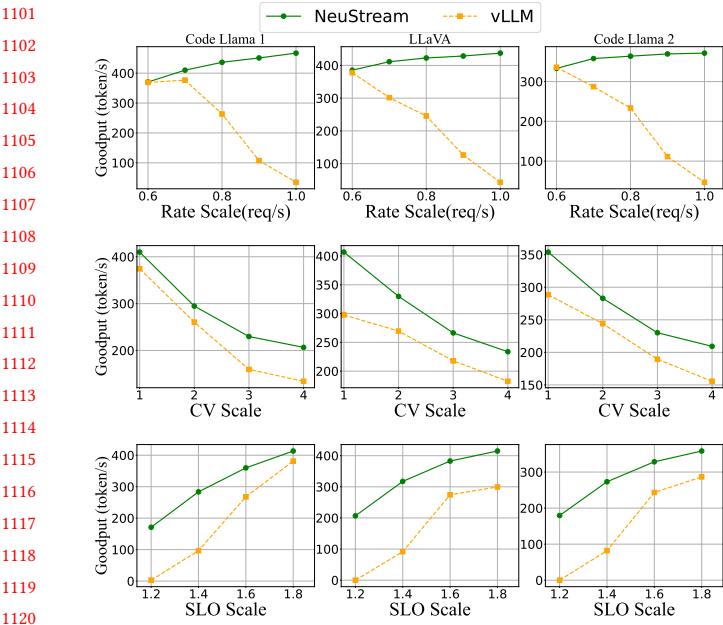


Figure 17. Goodput of agent’s decode module on RTX4090.

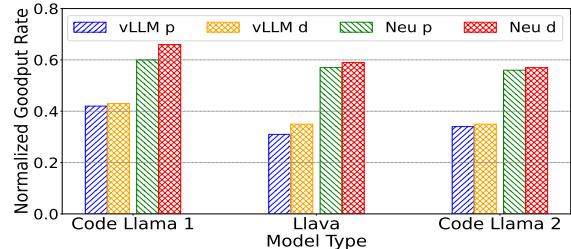


Figure 19. module normalized goodputs in MatPlotAgent

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 [36].

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 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.

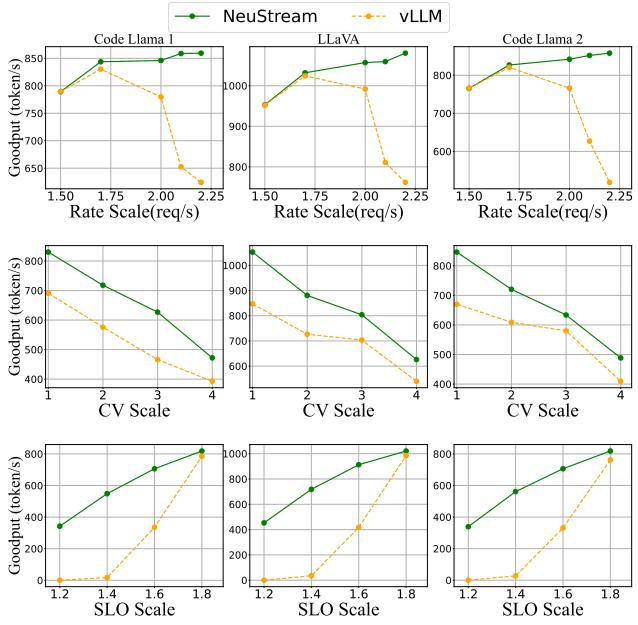


Figure 18. Goodput of agent’s decode module on H100.

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. Although time-sharing of SPUs increases heterogeneous GPU kernels to be concurrently executed, it 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 [7] leverages the predictable execution times of individual DNN inferences. Shepherd [32] aggregates request streams into moderately sized groups to improve predictability. Systems like Nexus[22], INFaaS [19], gplet [4], Gslice [6], focus on the serving of multiple static DNN models with spatio-temporal sharing of GPUs. AlpaServe [10] 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 [5] 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 [23] and Cocktailer [31] express and execute dynamic neural networks from a compiler perspective.

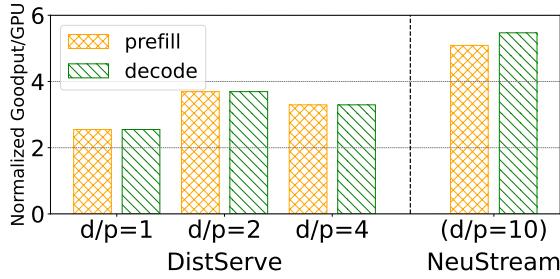


Figure 20. Normalized per-GPU goodput of instances with OPT-13B model on LMSYS-Chat-1M dataset and H100.

NEUSTREAM is orthogonal to these works by exploring the dynamism and batching opportunity across multiple inputs and supports SLO-aware DNN serving.

Specialized systems for LLM serving. LLM inference serving is a fast-developing area. Orca [29] introduces iteration-level scheduling to increase the GPU utilization. vLLM [9] 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[14], Spark [30], Flink [3], StreamScope [11], Apache Storm [24], 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 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., DistServe [36] and Splitwise [16]), 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 goodput of each instance in DistServe, when serving OPT-13B on LMSYS-Chat-1M dataset used in Sec. 6.3. Here, D/P gives the resource usage ratio of decoding instance to prefill instance. Note that the minimum granularity of resource allocation in DistServe is the entire GPU. Given the relatively

low average input-output ratio of requests, we observe that orchestrating prefill and decoding in DistServe is challenging at the coarse-grained GPU allocation model: under tight SLO constraint, insufficient GPU allocation for decoding instance (e.g., D/P=1) at the consumer side would limit the full processing capability and resource utilization of prefill instance at the producer side, however, increasing GPU allocation for decoding instance faces the growing scaling overhead of model parallelism (e.g., the per-GPU goodput even degrades when increasing D/P ratio from 2 to 4). By contrast, NEUSTREAM can minimize the mismatch between the processing capability of P/D instances at fine-grained SPUs abstraction and allocation.

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

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

Suitable scenarios. When employing NEUSTREAM, an interesting question is what are the scenarios one benefits from using it? As illustrated in Sec. 2, NEUSTREAM is useful for serving models with dynamic execution paths or distinct functional phases. However, serving static models with homogeneous phases might not benefit from NEUSTREAM.

9 Conclusion

NEUSTREAM addresses the inefficiencies of serving modern DNN applications with dynamic execution patterns on GPUs. By introducing the concepts of streams and stream-modules, NEUSTREAM enables a continuous data flow execution model that effectively maximizes GPU utilization while meeting stringent latency SLOs. The abstraction of GPU resources into Streaming Processing Units and the decomposition of scheduling into intra-module and inter-module levels further

enhance its scalability and scheduling efficiency. Our evaluation of NEUSTREAM demonstrates its superior performance in serving complex DNN applications, such as LLM-based chat, image generation, and multi-agents. As DNN models continue to evolve and become increasingly complex, NEUSTREAM offers a scalable and efficient solution for future advancements in deep learning applications.

References

- [1] HuggingFace. Diffusers doc. https://huggingface.co/docs/diffusers/v0.27.0/en/stable_diffusion.
- [2] PyTorch. <https://pytorch.org>.
- [3] 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).
- [4] Seungbeom Choi, Sunho Lee, Yeonjae Kim, Jongse Park, Youngjin Kwon, and Jaehyun 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.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] 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).
- [9] 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 *Proceedings of the 29th Symposium on Operating Systems Principles*. 611–626.
- [10] Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E Gonzalez, et al. 2023. AlpaServe: Statistical multiplexing with model parallelism for deep learning serving. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. 663–679.
- [11] Wei Lin, Zhengping Qian, Junwei Xu, Sen Yang, Jingren Zhou, and Lidong Zhou. 2016. {StreamScope}: Continuous Reliable Distributed Processing of Big Data Streams. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*. 439–453.
- [12] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744* (2023).
- [13] Amirhossein Mirhosseini, Sameh Elnikety, and Thomas F Wenisch. 2021. Parslo: A gradient descent-based approach for near-optimal partial slo allotment in microservices. In *Proceedings of the ACM Symposium on Cloud Computing*. 442–457.
- [14] Derek G. Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martín Abadi. 2013. Naiad: a timely dataflow system.
- In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles* (Farmington, Pennsylvania) (*SOSP '13*). Association for Computing Machinery, New York, NY, USA, 439–455. <https://doi.org/10.1145/2517349.2522738>
- [15] NVIDIA Corporation. 2024. NVIDIA Multi-Process Service. <https://docs.nvidia.com/deploy/mps/index.html> Accessed: 2024-05-21.
- [16] Pratyush Patel, Esha Choukse, Chaojie Zhang, Íñigo Goiri, Aashaka Shah, Saeed Maleki, and Ricardo Bianchini. 2023. Splitwise: Efficient generative llm inference using phase splitting. *arXiv preprint arXiv:2311.18677* (2023).
- [17] William Peebles and Saining Xie. 2023. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4195–4205.
- [18] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [19] Francisco Romero, Qian Li, Neeraja J Yadwadkar, and Christos Kozyrakis. 2021. {INFaaS}: Automated model-less inference serving. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*. 397–411.
- [20] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).
- [21] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings*. 1–10.
- [22] Haichen Shen, Lequn Chen, Yuchen Jin, Liangyu Zhao, Bingyu Kong, Matthai Philipose, Arvind Krishnamurthy, and Ravi Sundaram. 2019. Nexus: A GPU cluster engine for accelerating DNN-based video analysis. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. 322–337.
- [23] Haichen Shen, Jared Roesch, Zhi Chen, Wei Chen, Yong Wu, Mu Li, Vin Sharma, Zachary Tatlock, and Yida Wang. 2021. Nimble: Efficiently compiling dynamic neural networks for model inference. *Proceedings of Machine Learning and Systems* 3 (2021), 208–222.
- [24] Ankit Toshniwal, Siddarth Taneja, Amit Shukla, Karthik Ramasamy, Jignesh M Patel, Sanjeev Kulkarni, Jason Jackson, Krishna Gade, Maosong Fu, Jake Donham, et al. 2014. Storm@twitter. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. 147–156.
- [25] Bhuvan Urgaonkar, Prashant Shenoy, Abhishek Chandra, Pawan Goyal, and Timothy Wood. 2008. Agile dynamic provisioning of multi-tier internet applications. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)* 3, 1 (2008), 1–39.
- [26] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E Gonzalez. 2018. Skipnet: Learning dynamic routing in convolutional networks. In *Proceedings of the European conference on computer vision (ECCV)*. 409–424.
- [27] Canwen Xu and Julian McAuley. 2022. A survey on dynamic neural networks for natural language processing. *arXiv preprint arXiv:2202.07101* (2022).
- [28] Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong, Xu Han, Yukun Yan, Zhenghao Liu, Zhixing Tan, Pengyuan Liu, Dong Yu, et al. 2024. MatPlotAgent: Method and Evaluation for LLM-Based Agentic Scientific Data Visualization. *arXiv preprint arXiv:2402.11453* (2024).
- [29] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. 2022. Orca: A distributed serving system for Transformer-Based generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. 521–538.

1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430

- 1431 [30] Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott
1432 Shenker, and Ion Stoica. 2010. Spark: Cluster computing with working
1433 sets. In *2nd USENIX Workshop on Hot Topics in Cloud Computing*
1434 (*HotCloud 10*).
1435 [31] Chen Zhang, Lingxiao Ma, Jilong Xue, Yining Shi, Ziming Miao,
1436 Fan Yang, Jidong Zhai, Zhi Yang, and Mao Yang. 2023. Cocktailler:
1437 Analyzing and Optimizing Dynamic Control Flow in Deep Learn-
1438 ing. In *17th USENIX Symposium on Operating Systems Design and*
1439 *Implementation (OSDI 23)*. USENIX Association, Boston, MA, 681–
1440 699. <https://www.usenix.org/conference/osdi23/presentation/zhang-chchen>
1441 [32] Hong Zhang, Yupeng Tang, Anurag Khandelwal, and Ion Stoica.
1442 2023. {SHEPHERD}: Serving {DNNs} in the wild. In *20th USENIX*
1443 *Symposium on Networked Systems Design and Implementation*
1444 (*NSDI 23*). 787–808.
1445 [33] Hong Zhang, Yupeng Tang, Anurag Khandelwal, and Ion Stoica.
2023. SHEPHERD: Serving DNNs in the Wild. In *20th USENIX*
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485 *Symposium on Networked Systems Design and Implementation*
1486 (*NSDI 23*). USENIX Association, Boston, MA, 787–808. <https://www.usenix.org/conference/nsdi23/presentation/zhang-hong>
1487 [34] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya
1488 Chen, Shuhui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Vic-
1489 toria Lin, et al. 2022. Opt: Open pre-trained transformer language
1490 models. *arXiv preprint arXiv:2205.01068* (2022).
1491 [35] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan
1492 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric P.
1493 Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024. LMSYS-
1494 Chat-1M: A Large-Scale Real-World LLM Conversation Dataset. (2024).
1495 *arXiv:2309.11998* [cs.CL]
1496 [36] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xu-
1497 anzhe Liu, Xin Jin, and Hao Zhang. 2024. DistServe: Disaggregating
1498 Prefill and Decoding for Goodput-optimized Large Language Model
1499 Serving. *arXiv preprint arXiv:2401.09670* (2024).
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540