

1 **GPU-HADVPPM V1.0: high-efficient parallel GPU design of the**
2 **Piecewise Parabolic Method (PPM) for horizontal advection in**
3 **air quality model (CAMx V6.10)**

4 **Kai Cao¹, Qizhong Wu¹, Lingling Wang², Nan Wang², Huaqiong Cheng¹, Xiao**
5 **Tang³, Dongqing Li¹, and Lanning Wang¹**

6 ¹College of Global Change and Earth System Science, Beijing Normal University,
7 Beijing 100875, China

8 ²Henan Ecological Environmental Monitoring Centre and Safety Center, Henan Key
9 Laboratory of Environmental Monitoring Technology, Zhengzhou 450008, China

10 ³State Key Laboratory of Atmospheric Boundary Layer Physics and Atmospheric
11 Chemistry, Institute of Atmospheric Physics, Chinese Academy of Science, Beijing
12 100029, China

13

14 **Correspondence to:** Qizhong Wu (wqizhong@bnu.edu.cn); Lingling
15 Wang(928216422@qq.com); Lanning Wang (wangln@bnu.edu.cn)

16

17 **Abstract.** With semiconductor technology gradually approaching its physical and
18 thermal limits, Graphics processing unit (GPU) is becoming an attractive solution in
19 many scientific applications due to their high performance. This paper presents an
20 application of GPU accelerators in air quality model. We endeavor to demonstrate an
21 approach that runs a PPM solver of horizontal advection (HADVPPM) for air quality
22 model CAMx on GPU clusters. Specifically, we first convert the HADVPPM to a new
23 Compute Unified Device Architecture C (CUDA C) code to make it computable on the
24 GPU (GPU-HADVPPM). Then, a series of optimization measures are taken, including
25 reducing the CPU-GPU communication frequency, increasing the size of data
26 computation on GPU, optimizing the GPU memory access, and using thread and block
27 indices in order to improve the overall computing performance of CAMx model
28 coupled with GPU-HADVPPM (named as CAMx-CUDA model). Finally, a
29 heterogeneous, hybrid programming paradigm is presented and utilized with the GPU-
30 HADVPPM on GPU clusters with Message Passing Interface (MPI) and CUDA.
31 Offline experiment results show that running GPU-HADVPPM on one NVIDIA Tesla
32 K40m and NVIDIA Tesla V100 GPU can achieve up to 845.4x and 1113.6x

33 acceleration. By implementing a series of optimization schemes, the CAMx-CUDA
34 model resulted in a 29.0x and 128.4x improvement in computational efficiency using a
35 GPU accelerator card on a K40m and V100 cluster, respectively. In terms of the single-
36 module computational efficiency of GPU-HADVPPM, it can achieve 1.3x and 18.8x
37 speedup on NVIDIA Tesla K40m GPU and NVIDA Tesla V100 GPU respectively. The
38 multi-GPU acceleration algorithm enables 4.5x speedup with 8 CPU cores and 8 GPU
39 accelerators on V100 cluster.

40 **1. Introduction**

41 Since the introduction of the personal computer in the late 1980s, the computer
42 and mobile device industry has been one of the most flourishing markets all over the
43 world (Bleichrodt et al., 2012). In recent years, the improvement of the performance of
44 the Central Processing Unit (CPU) is limited by its heat dissipation, the development
45 of Moore's Law has flattened. A common trend in high-performance computing today
46 is the utilization of hardware accelerators that execute codes rich in data parallelism to
47 form high-performance heterogeneous system. GPUs are widely used as accelerators
48 due to high peak performance offered. In the top ten supercomputing list released in
49 December 2022 (<https://www.top500.org/lists/top500/list/2022/11/>, last access: 19
50 December 2022), there are seven heterogeneous supercomputing platforms built with
51 CPU processors and GPU accelerators, of which the top one Frontier at the Oak Ridge
52 National Laboratory uses AMD's third-generation EPYC CPU and AMD Instinct
53 MI250X GPU, and its computing performance reaches Exascale (10^{18} calculations per
54 second) for the first time ([https://www.amd.com/en/press-releases/2022-05-30-world-
55 s-first-exascale-supercomputer-powered-amd-epyc-processors-and-amd](https://www.amd.com/en/press-releases/2022-05-30-world-s-first-exascale-supercomputer-powered-amd-epyc-processors-and-amd), last access:
56 19 December 2022). Such powerful computing performance of the heterogeneous
57 system not only injects new vitality into high-performance computing, but also provides
58 new solutions for improving the performance of numerical models in geoscience.

59 The GPU has proven successful in weather models such as Non-Hydrostatic

60 Icosahedral Model (NIM; Govett et al.,2017), Global/Regional Assimilation and
61 Prediction System (GRAPES; Xiao et al., 2022), and Weather Research and Forecasting
62 model (WRF; Huang et al., 2011; Huang et al., 2012; Mielikainen et al., 2012a;
63 Mielikainen et al., 2012b; Mielikainen et al., 2013a ; Mielikainen et al., 2013b; Price et
64 al., 2014; Huang et al., 2015), ocean models such as LASG/IAP Climate System Ocean
65 Model (LICOM; Jiang et al., 2019; Wang et al., 2021a) and Princeton Ocean Model
66 (POM; Xu et al., 2015), and the Earth System Model of Chinese Academy of Sciences
67 (~~CAS-ESM~~; Wang et al., 2016; Wang et al., 2021b).

删除的内容: CAS-EMS

68 Govett et al., (2017) used Open Accelerator (OpenACC) directives to port the
69 dynamics of NIM to the GPU and achieved 2.5x acceleration. Also using OpenACC
70 directives, Xiao et al., (2022) ported the PRM (Piecewise Rational Method) scalar
71 advection scheme in the GRAPES to the GPU, achieving up to 3.51x faster than 32
72 CPU cores. In terms of the most widely used WRF, several parameterization schemes,
73 such as RRTMG_LW scheme (Price et al., 2014), 5-layer thermal diffusion scheme
74 (Huang et al., 2015), Eta Ferrier Cloud Microphysics scheme (Huang et al., 2012),
75 Goddard Shortwave scheme (Mielikainen et al., 2012a), Kessler cloud microphysics
76 scheme (Mielikainen et al., 2013b), SBU-YLIN scheme (Mielikainen et al., 2012b),
77 WMS5 scheme (Huang et al., 2011), WMS6 scheme (Mielikainen et al., 2013a), etc.,
78 have been ported heterogeneously using CUDA C and achieved 37x~896x acceleration
79 results. The LICOM has carried out heterogeneous porting using OpenACC (Jiang et
80 al., 2019) and Heterogeneous-compute Interface for Portability C (HIP C) technologies,
81 and achieved up to 6.6x and 42x acceleration, respectively (Wang et al., 2021a). For the
82 Princeton Ocean Model, Xu et al., (2015) use CUDA C to carry out heterogeneous
83 porting and optimization, the performance of gpu-POM v1.0 on four GPUs is
84 comparable to that on 408 standard Intel Xeon X5670 CPU cores. In terms of climate
85 system model, Wang et al., (2016) and Wang et al., (2021b) used CUDA Fortran and
86 CUDA C to carry out heterogeneous porting of the RRTMG_SW and RRTMG_LW
87 scheme of the atmospheric component model of the ~~CAS-ESM~~ earth system model,
88 and achieved a 38.88x and 77.78x acceleration respectively.

删除的内容: CAS-EMS

91 Programming a GPU accelerator can be a hard and error-prone process that
92 requires specially designed programming methods, there are three widely used methods
93 for porting program to GPUs as described above. The first method uses the OpenACC
94 directive (<https://www.openacc.org/>, last access: 19 December 2022) which provides a
95 set of high-level directives that enable C/C++ and Fortran programmers to utilize
96 accelerators. The second method uses CUDA Fortran. CUDA Fortran is a software
97 compiler which co-developed by the Portland Group (PGI) and NVIDIA, and tool chain
98 for building performance optimized GPU-accelerated Fortran applications targeting the
99 NVIDIA GPU platform (<https://developer.nvidia.com/cuda-fortran>, last access: 19
100 December 2022). CUDA C involves rewriting the entire program using standard C
101 programming language and low-level CUDA subroutines
102 (<https://developer.nvidia.com/cuda-toolkit>, last access: 19 December 2022) to support
103 the NVIDIA GPU accelerator. Compared to the other two technologies, CUDA C
104 porting scheme is the most complex, but its computational performance is the highest
105 (Mielikainen et al., 2012b; Wahib and Maruyama, 2013; Xu et al., 2015).

106 Air quality models are critical to understanding how the chemistry and
107 composition of atmospheric may change over 21st century, as well as preparing adaptive
108 responses or developing mitigation strategies. Because air quality models need to take
109 into account the complex physicochemical processes that occur in the atmosphere of
110 anthropogenic and naturally emissions, simulations are computationally expensive.
111 Compared to the other geoscientific numerical models, few research have carried out
112 heterogeneous porting of air quality models. In this study, CUDA C scheme was
113 implemented in this paper to carry out the hotspot module porting attempt of CAMx in
114 order to improve the computation efficiency.

115 **2. The CAMx model and experiments**

116 **2.1. Model description**

117 CAMx model is a state-of-the air quality model developed by Ramboll Environ

118 (<https://www.camx.com/>, last access: 19 December 2022). CAMx version 6.10 (CAMx
 119 V6.10; ENVIRON, 2014) is chosen in this study, it simulates the emission, dispersion,
 120 chemical reaction, and removal of pollutants by marching the Eulerian continuity
 121 equation forward in time for each chemical species on a system of nested three-
 122 dimensional grids. The Eulerian continuity equation is expressed mathematically in
 123 terrain-following height coordinates as formula (1):

$$124 \quad \frac{\partial c_i}{\partial t} = -\nabla_H \cdot V_H c_i + \left[\frac{\partial(c_i \eta)}{\partial z} - c_i \frac{\partial^2 h}{\partial z \partial t} \right] + \nabla \cdot \rho K \nabla (c_i / \rho)$$

$$125 \quad + \frac{\partial c_i}{\partial t} \Big|_{Emission} + \frac{\partial c_i}{\partial t} \Big|_{Chemistry} + \frac{\partial c_i}{\partial t} \Big|_{Removal} \quad (1)$$

$$126 \quad \nabla_H \cdot \rho V_H = \frac{m^2}{A_{yz}} \frac{\partial}{\partial x} \left(\frac{u A_{yz} \rho}{m} \right) + \frac{m^2}{A_{xz}} \frac{\partial}{\partial y} \left(\frac{v A_{xz} \rho}{m} \right) \quad (2)$$

127 The first term on the right-hand side represents horizontal advection. In the
 128 numerical methods, the equation of horizontal advection (described in formula (2)) is
 129 performed using the area preserving flux-form advection solver of the Piecewise
 130 Parabolic Method (PPM) of Colella and Woodward (1984) as implemented by Odman
 131 and Ingram (1996). The PPM solution of horizontal advection (HADVPPM) was
 132 incorporated into CAMx model because it provides higher order accuracy with minimal
 133 numerical diffusion.

134 In the Fortran code implementation of HADVPPM scheme, the CAMx main
 135 program calls the emistrns program, which mainly performs the physical processes such
 136 as emission, diffusion, advection and dry/wet deposition of pollutants. And then, the
 137 horizontal advection program is invoked by emistrns program to solve the horizontal
 138 advection equation by using the HADVPPM scheme.

139 **2.2. Benchmark performance experiments**

140 The first step of the porting is to test the performance of CAMx benchmark version
 141 and identify the hotspots of the model. On the Intel x86 CPU platform, we launch two
 142 processes concurrently to run the CAMx and take advantage of the Intel Trace Analyzer

143 Collector(ITAC; [https://www.intel.com/content/www/us/en/docs/trace-analyzer-](https://www.intel.com/content/www/us/en/docs/trace-analyzer-collector/get-started-guide/2021-4/overview.html)
144 [collector/get-started-guide/2021-4/overview.html](https://www.intel.com/content/www/us/en/docs/trace-analyzer-collector/get-started-guide/2021-4/overview.html), last access: 19 December 2022) and
145 Intel VTune Profiler(VTune;[https://www.intel.com/content/www/us/en/develop/documentation/vtune-](https://www.intel.com/content/www/us/en/develop/documentation/vtune-help/top.html)
146 [help/top.html](https://www.intel.com/content/www/us/en/develop/documentation/vtune-help/top.html), last access: 19 December 2022) performance analysis tools to collect
147 performance information during CAMx operation.

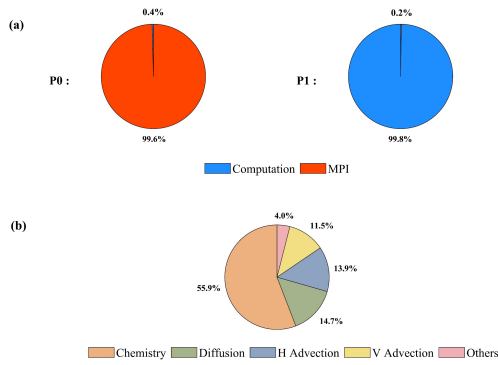
149 The general MPI performance can be reported by the ITAC tool, and MPI load
150 balance information, computation and communication profiling of each process is
151 shown as Fig. 1a. During the running process of CAMx model, Process 0 (P0) spends
152 99.6% of the time on the MPI_Barrier function and only 0.4% of the time on
153 computation, while Process 1(P1) spends 99.8% of its time computation and only 0.2%
154 of its time receiving messages from P0. It is indicated that the parallel design of CAMx
155 model adopts Master-Slave mode, P0 is responsible for inputting and outputting data
156 and calling the MPI_Barrier function to synchronize the process, so there is a lot of
157 MPI waiting time. The other processes are responsible for computation.

158 The VTune tool detects each module's runtime and the most time-consuming
159 functions on P1. As shown in Figure 1b, the top four time-consuming modules are
160 chemistry, diffusion, horizontal advection, and vertical advection in the CAMx model.
161 In the above four modules, the top five most time-consuming programs are ebrate,
162 hadvppm, tridiag, diffus, and ebisolv programs, and the total runtime of P1 is 325.1
163 seconds. Top1 and Top2's most time-consuming programs take 49.4 and 35.6 seconds,
164 respectively.

165 By consideration, the hadvppm program was selected to carry out heterogeneous
166 porting for some reasons. Firstly, the advection module is one of the compulsory
167 modules of the air quality model, which is mainly used to simulate the transport process
168 of air pollutants, and it is also a hotspot module detected by the Intel VTune tool. Then,
169 typical air quality models CAMx, CMAQ, and NAQPMS include advection modules
170 and use the exact PPM advection solver. The heterogeneous version developed in this
171 study can be directly applied to the above models. Furthermore, the weather model (e.g.,

删除的内容: The VTune tool is used to detect the runtime of each module and the most time-consuming functions on P1. As shown in Figure 1b, the top four time-consuming modules are chemistry, diffusion, horizontal advection, and vertical advection in CAMx model. The top five most time-consuming programs and their elapsed time are in Table 1. The total runtime of P1 is 325.1 seconds, and the top five most time-consuming programs are ebrate, hadvppm, tridiag, diffus, and ebisolv program. Top1 and Top2's most time-consuming programs take 49.4 and 35.6 seconds, respectively. By viewing the Fortran code of the above programs, the hadvppm program has few calculation branches, and its calculation process does not involve iterative operations, which satisfies the basic conditions for the program to run on the GPU. Therefore, a GPU acceleration version of the HADVPPM scheme, namely GPU-HADVPPM, is built to improve CAMx performance.^d

189 WRF) also contains an advection module, so this study's heterogeneous porting method
 190 and experience can be used for reference. Therefore, a GPU acceleration version of the
 191 HADVPPM scheme, namely GPU-HADVPPM, is built to improve CAMx
 192 performance.



193
 194 **Figure 1.** The computation performance of the modules in the CAMx model. (a) Computation and
 195 communication profiling of P0 and P1. (b) Overhead proportions of P1. The top four most time-
 196 consuming modules are chemistry, diffusion, horizontal advection, and vertical advection.

删除的内容: **Table 1.** The top five most time-consuming programs on the P1 (Total runtime is 325.1 seconds).[†] ... [1]

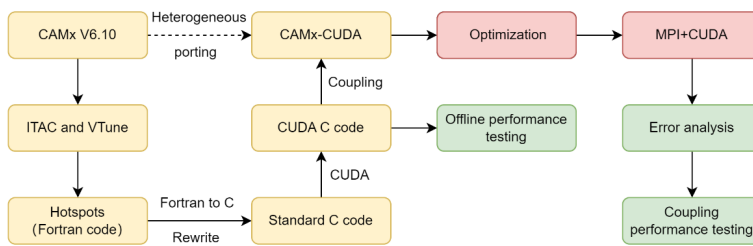
199 **2.3. Porting scheme introduction**

200 The heterogeneous scheme of CAMx-CUDA is shown in Figure 2. The second
 201 time-consuming hadvppm program in the CAMx model was selected to implement the
 202 heterogeneous porting. In order to map the hadvppm program to the GPU, the Fortran
 203 code was converted to standard C code. Then, CUDA programing language, which was
 204 taylor-made for NVIDIA, was added to convert the standard C code into CUDA C for
 205 data-parallel execution on GPU, as GPU-HADVPPM. It prepared the input data for
 206 GPU-HADVPPM by constructing random numbers and tested its offline performance
 207 on the GPU platform.

删除的内容: The heterogeneous scheme of CAMx-CUDA is shown in Figure 2. The second time-consuming program hadvppm in CAMx model, was selected to implement the heterogeneous porting. In order to map the hadvppm program to the GPU, the Fortran code of hadvppm program is converted to standard C code. Then, CUDA programing language which is tailor-made for NVIDIA was added to convert the standard C code into CUDA C for data-parallel execution on GPU, as GPU-HADVPPM. It prepares the input data for GPU-HADVPPM by constructing random numbers, and tests its offline performance on GPU platform.[†]

208 After coupling GPU-HADVPPM to CAMx model, the advection module code was

222 optimized according to the characteristics of GPU architecture to improve the overall
 223 computational efficiency on CPU-GPU heterogeneous platform. And then, the multi-
 224 CPU core and multi-GPU card acceleration algorithm was adopted to improve the
 225 parallel extensibility of heterogeneous computing. Finally, the coupling performance
 226 test is implemented after verifying the different CAMx model simulation results.



227
 228 **Figure 2.** Heterogeneous porting scheme of CAMx-CUDA model.

229 2.4. Hardware components and software environment of the testing system

230 The experiments are conducted on two GPU clusters: K40m and V100.
 231 hardware components and software environment of the two clusters are listed in Table
 232 1. The K40m cluster is equipped with two 2.5GHz 16-core Intel Xeon E5-2682 v4 CPU
 233 processors and one NVIDIA Tesla K40m GPU card on each node. The NVIDIA Tesla
 234 K40m GPU has 2880 CUDA cores with 12GB of memory. The V100 cluster contains
 235 two 2.7GHz 24-core Intel Xeon Platinum 8168 processors and eight NVIDIA Tesla
 236 V100 GPU cards with 5120 CUDA cores and 16GB memory on each card.

237 **Table 1.** Configurations of GPU cluster.

	Hardware components	
	CPU	GPU
K40m cluster	Intel Xeon E5-2682 v4 CPU @2.5GHz, 16 cores	NVIDIA Tesla K40m, 2880 CUDA cores, 12GB memory
V100 cluster	Intel Xeon Platinum 8168 CPU @2.7 GHz, 24 cores	NVIDIA Tesla V100, 5120 CUDA cores, 16GB memory
	Software environment	
	Compiler and MPI	Programming Model
K40m cluster	Intel-2021.4.0	CUDA-10.2

删除的内容: 2

删除的内容: 2

V100 cluster	Intel-2019.1.144	CUDA-10.0
--------------	------------------	-----------

240 For Fortran and standard C programming, Intel Toolkit (including compiler and
 241 MPI library) version 2021.4.0 and version 2019.1.144 are employed for compiling on
 242 Intel Xeon E4-2682 v4 CPU and Intel Xeon Platinum 8168 CPU, respectively. And
 243 then, CUDA version 10.2 and version 10.0 are employed on NVIDIA Tesla K40m GPU
 244 and NVIDIA Tesla V100 GPU. CUDA (NVIDIA, 2020) is an extension of the C
 245 programming language that offers direct programming of the GPUs. In CUDA
 246 programming, what is called a kernel is actually a subroutine that can be executed on
 247 the GPU. The underlying code in the kernel is divided into a series of threads, each with
 248 a unique "ID" number that can simultaneously process different data through a single-
 249 instruction multiple-thread (SIMT) parallel mode. These threads are grouped into
 250 equal-sized thread blocks, which are organized into a grid.

251 **3. Porting and optimization of CAMx advection module on heterogeneous**
 252 **platform**

253 **3.1. Mapping HADVPPM scheme to GPU**

254 **3.1.1. Manual code translation from Fortran to standard C**

255 As the CAMx V6.10 code was written in Fortran 90, we rewrote the hadvppm
 256 program from Fortran to CUDA C. As an intermediate conversion step, we refactor the
 257 original Fortran code using standard C. During the refactoring, some considerations are
 258 listed in Table 2;

259 (1) The subroutine name refactored with standard C must be followed by an
 260 underscore identifier, which can only be recognized when Fortran calls.

261 (2) In Fortran language, the parameters are transferred by memory address by
 262 default. In the case of mixed programming in Fortran and standard C, parameters
 263 transferred by Fortran are processed by the pointer in standard C.

264 (3) Variable precision types defined in standard C must be strictly consistent with

删除的内容: 3

266 those in Fortran.

267 (4) Some built-in functions in Fortran are not available in standard C and need to
268 be defined in standard C macro definitions.

269 (5) For multidimensional arrays, Fortran and standard C follow column-major and
270 row-major order in-memory read and write, respectively;

271 (6) Array subscripts in Fortran and standard C are indexed from any integer and 0,
272 respectively.

273 **Table 2.** Some considerations during Fortran to C refactoring.

删除的内容: 3

	Fortran code	C code
Function name	<i>subroutine hadvppm()</i>	<i>void hadvppm()</i>
Parameter passing	<i>hadvppm(nn,dt,dx,con,vel,area,areav, flxarr,mynn)</i>	<i>hadvppm(int *nn,float *dt, float *dx,float *con,float *vel,float *area,float *areav, float *flxarr,int *mynn)</i>
Variable precision	<i>real(kind=8) x</i>	<i>double x</i>
Built-in functions	<i>max</i>	<i>#define Max(a, b) ((a)>(b)?(a):(b))</i>
Memory read and write for multidimensional array	Column-major	Row-major
Array subscript index	Starting from any integer	Starting from 0

274

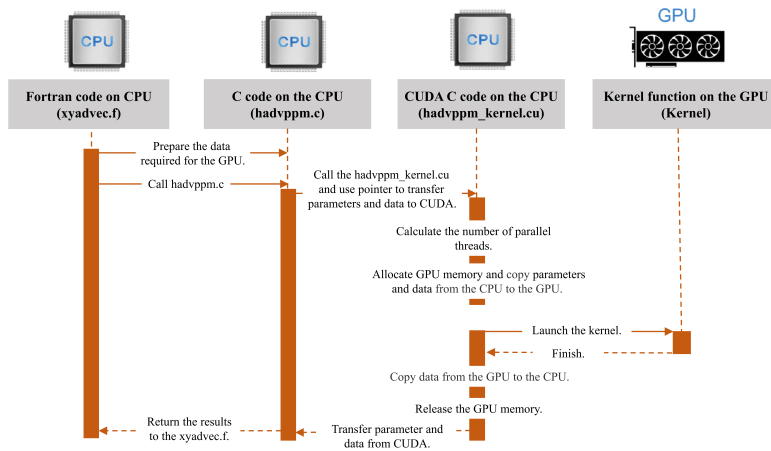
275 3.1.2. Converting standard C code into CUDA C

276 After refactoring the Fortran code of the hadvppm program with standard C,
277 CUDA was used to convert the C code into CUDA C to make it computable on the
278 GPU. A standard C program using CUDA extensions distributes a large number of
279 copies of the kernel functions into available multiprocessors and executes them
280 simultaneously on the GPU.

281 Figure 3 shows the implementation process of the GPU-HADVPPM. As

283 mentioned in Sect.2.1, xyadvec program calls the hadvppm program to solve the
 284 horizontal advection function. Since the rewritten CUDA program cannot be called
 285 directly by Fortran program (xyadvec.f), we add an intermediate subroutine
 286 (hadvppm.c) as an interface to transfer the parameters and data required for GPU
 287 computing from xyadvec Fortran program to hadvppm_kernel CUDA C program.

288 A CUDA program automatically uses numerous threads on GPU to execute kernel
 289 functions. Therefore, the hadvppm_kernel CUDA C program first calculates the
 290 number of parallel threads according to the array dimension. And then allocate GPU
 291 memory, and copy parameters and data from the CPU to the GPU. As the CUDA
 292 program launches a large number of parallel threads to execute kernel functions
 293 simultaneously, the computation results will be copied from the GPU back to the CPU.
 294 Finally, the GPU memory is released, and data computed on the GPU is returned to the
 295 xyadvec program via hadvppm C program.



296
 297 **Figure 3.** The calling and computation process of the GPU-HADVPPM on the CPU-GPU
 298 heterogeneous platform.

299 3.2. Coupling and optimization of GPU-HADVPPM scheme on a single GPU

300 After the hadvppm program was rewritten with standard C and CUDA, the
 301 implementation process of HADVPPM scheme is loaded from the CPU to the GPU.

302 And then, we coupled the GPU-HADVPPM to CAMx model. For ease of description,
303 we will refer to this original heterogeneous version of CAMx as CAMx-CUDA V1.0.
304 In the CAMx-CUDA V1.0, four external loops are nested when hadvppm C program is
305 called by the xyadvec program. It will result in the widespread data transfers from the
306 CPU to the GPU over the PCIe bus within a time step, making the computation of the
307 CAMx-CUDA V1.0 inefficient.

308 Therefore, we optimize the xyadvec Fortran program to significantly reduce the
309 frequency of data transmission between CPU and GPU, increase the amount of data
310 computation on GPU, and improve the total computing efficiency of the CAMx on
311 CPU-GPU heterogeneous platforms. In the original CAMx-CUDA V1.0, four external
312 loops outside of hadvppm C program and several one-dimensional arrays are computed
313 before calling hadvppm C program. Then the CPU will frequently launch the GPU and
314 transfer data to it within a time step. When the code optimization is completed, three or
315 four-dimensional arrays required for GPU computation within a time step will be sorted
316 before calling the hadvppm C program, and then the CPU will package and transfer the
317 arrays to the GPU in batches. The example of xyadvec Fortran program optimization
318 was shown in Figure S1.

319 The details of four different versions are shown in Table 3. In the CAMx-CUDA
320 V1.0, the Fortran code of the HADVPPM scheme was rewritten using standard C and
321 CUDA, and the xyadvec program is not optimized. The dimensions of the c1d variable
322 array transmitted to GPU in the X and Y directions are 157 and 145 in this case,
323 respectively. In CAMx-CUDA V1.1 and CAMx-CUDA V1.2, the c1d variable
324 transmitted from CPU to GPU are expanded to two (about 23,000 numbers) and four
325 dimensions (about 27.4 million numbers) by optimizing the xyadvec Fortran program
326 and hadvppm_kernel CUDA C program, respectively.

327 The order in which data is accessed in GPU memory affects the computational
328 efficiency of the code. In the CAMx-CUDA V1.3 of the Table 4, we further optimized
329 the order in which data is accessed in GPU memory based on the order in which it is
330 stored in memory, and eliminated unnecessary assignment loops that were added due

删除的内容: 4

332 to the difference in memory read order between Fortran and C.

333 As described in Sect.2.4, a thread is the basic unit of parallelism in CUDA
 334 programming. The structure of threads is organized into a three-level hierarchy. The
 335 highest level is a grid, which consists of three-dimensional thread blocks. The second
 336 level is a block, which also consists of three-dimensional threads. Built-in CUDA
 337 variable *threadIdx.x* determines a unique thread "ID" number inside a thread block.
 338 Similarly, built-in variable *blockIdx.x* and *blockIdx.y* determine which block to execute
 339 on, and the size of the block is determined by using the built-in variable *blockDim.x*.
 340 For the two-dimensional horizontal grid points, many threads and blocks can be
 341 organized so that each CUDA thread computes the results for different spatial positions
 342 simultaneously.

343 Before the CAMx-CUDA V1.4, the loops for three-dimension spatial grid points
 344 (i,j,k) are replaced by index computations only using thread index ($i = threadIdx.x +$
 345 $blockIdx.x * blockDim.x$), to use thread indexes only computes the grid point in the x or
 346 y direction simultaneous. In order to take full advantage of thousands of threads in the
 347 GPU, we implement thread and block indices ($i = threadIdx.x + blockIdx.x * blockDim.x;$
 348 $j = blockIdx.y$) to compute all horizontal grid points (i,j) simultaneous in the CAMx-
 349 CUDA V1.4. This is permitted because there are no interactions among horizontal grid
 350 points.

351 **Table 3.** The details of different CAMx-CUDA versions during optimization.

Version	Major revisions	Amount of data computation on GPU
CAMx-CUDA V1.0	The Fortran code of the HADVPPM subroutine was rewritten using standard C and CUDA, and <i>xyadvec.f</i> was not optimized.	157 and 145 in the x direction and y direction for the c1d variable, respectively.
CAMx-CUDA V1.1	Optimize <i>xyadec.f</i> and <i>hadvppm_kernel.cu</i> to expand the dimension of the array transmitted to the GPU from 1-dimensional to 2-dimensional.	157×145, about 23,000 numbers for the c2d variable.
CAMx-CUDA V1.2	Based on the CAMx-CUDA V1.1, the dimension of the array transmitted to the GPU is extended from 2 to 4 dimensions.	157×145×14×86, about 27.4 million numbers for the c4d variable.

删除的内容: 4

CAMx-CUDA V1.3	Based on the CAMx-CUDA V1.2, the order of GPU memory access is optimized and unnecessary assignment loops are eliminated.	157×145×14×86, about 27.4 million numbers for the c4d variable.
CAMx-CUDA V1.4	Based on the CAMx-CUDA V1.3, using thread and block indices ($i = threadIdx.x + blockIdx.x * blockDim.x; j = blockIdx.y$).	157×145×14×86, about 27.4 million numbers for the c4d variable.

353

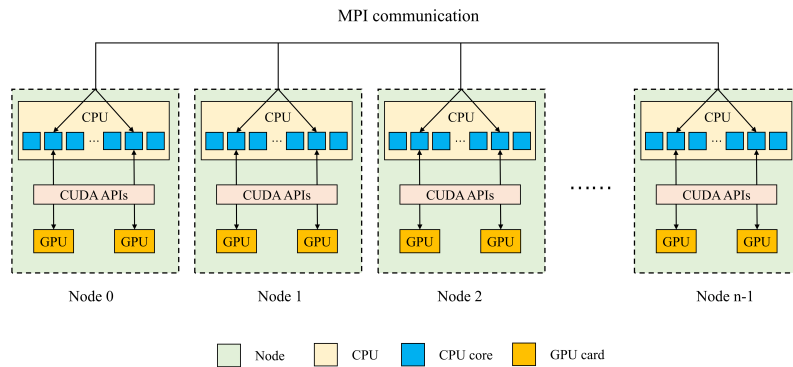
354 3.3. MPI+CUDA acceleration algorithm of CAMx-CUDA on multiple GPUs

355 Generally, super-large clusters have thousands of compute nodes. The current
 356 CAMx V6.10, implemented by adopting MPI communication technology, typically
 357 runs on dozens of compute nodes. Once the GPU-HADVPPM is coupled into the
 358 CAMx, it also has to run on multiple compute nodes which equipped one or more GPUs
 359 on each node. To make full use of multi-core and multi-GPU supercomputers and
 360 further improve the overall computational performance of the CAMx-CUDA, we adopt
 361 a parallel architecture with an MPI+CUDA hybrid paradigm, that is, the collaborative
 362 computing strategy of multiple CPU cores and multiple GPU cards is adopted during
 363 the operation of CAMx-CUDA model. Adopt this strategy, the GPU-HADVPPM can
 364 run on multiple GPUs, the Fortran code of other modules in CAMx-CUDA model can
 365 run on multiple CPU cores.

366 As is shown in Figure 4., after the simulated region is subdivided by MPI, a CPU
 367 core is responsible for the computation of a subregion. In order to improve the total
 368 computational performance of the CAMx-CUDA model, we further used the NVIDIA
 369 CUDA library to obtain the number of GPUs per node, and then used MPI process ID
 370 and remainder function to determine the GPU ID to be launched by each node. Finally,
 371 we used NVIDIA CUDA library `cudaSetDevice` to configure a GPU card for each CPU
 372 core.

373 According to the benchmark performance experiments, the parallel design of
 374 CAMx adopts Master-Slave mode, P0 is responsible for inputting and outputting data.
 375 If two processes (P0 and P1) were launched, only the P1 and its configured GPU

376 participate in integration.



377

378 **Figure 4.** An example of parallel architecture with an MPI+CUDA hybrid paradigm on multiple
379 GPUs.

380 4. Experimental results

381 The validation and evaluation of porting the HADVPPM scheme from the CPU to
382 the GPU platform were conducted using offline and coupling performance experiments.
383 First, we validated the result between different CAMx versions, and then the offline
384 performance of the GPU-HADVPPM on a single GPU was tested by offline experiment.
385 Finally, the coupling performance experiments illustrate its potential in three
386 dimensions with varying chemical regimes. Sect.4.2 and Sect.4.4, the CAMx version
387 of the HADVPPM scheme written by Fortran language, standard C, and CUDA C, is
388 named F, C, and CUDA C, respectively.

删除的内容: The validation and evaluation of porting the HADVPPM scheme from CPU to GPU platform were conducted using offline and coupling performance experiments. First, we validate the result between different CAMx versions, and then the offline performance of the GPU-HADVPPM on a single GPU was tested by offline experiment. Finally, the coupling performance experiments illustrate its potential in three dimensions with varying chemical regimes. In Sect.4.2 and Sect.4.4, the CAMx version of the HADVPPM scheme written by Fortran language, standard C, and CUDA C are named as F, C, and CUDA C, respectively.

389 4.1. Experimental setup

390 The test case is a 48h simulation covering the Beijing, Tianjin and part region of
391 Hebei province. The horizontal resolution is 3km with 145×157 grid boxes. The
392 model adopted 14 vertical layers. The simulation started at 12:00 UTC, 01 November

405 2020, and ended at 12:00 UTC, 03 November 2020. The meteorological fields driving
406 the CAMx model were provided by the Weather Research and Forecasting (WRF;
407 Skamarock et al., 2008) model. The Sparse Matrix Operator Kernel Emission (SMOKE;
408 Houyoux and Vukovich, 1999) version 2.4 model is used to provide gridded emission
409 data for the CAMx model. The emission inventories (Sun et al., 2022) include the
410 regional emissions in East Asia that were obtained from the Transport and Chemical
411 Evolution over the Pacific (TRACE-P; Streets et al., 2003; Streets et al., 2006) project,
412 30-min(about 55.6km at mid-latitude) spatial resolution Intercontinental Chemical
413 Transport Experiment-Phase B (INTEX-B; Zhang et al., 2009) and the updated regional
414 emission inventories in North China. The physical and chemical numerical methods
415 selected during CAMx model integration are listed in Table S2.

删除的内容:

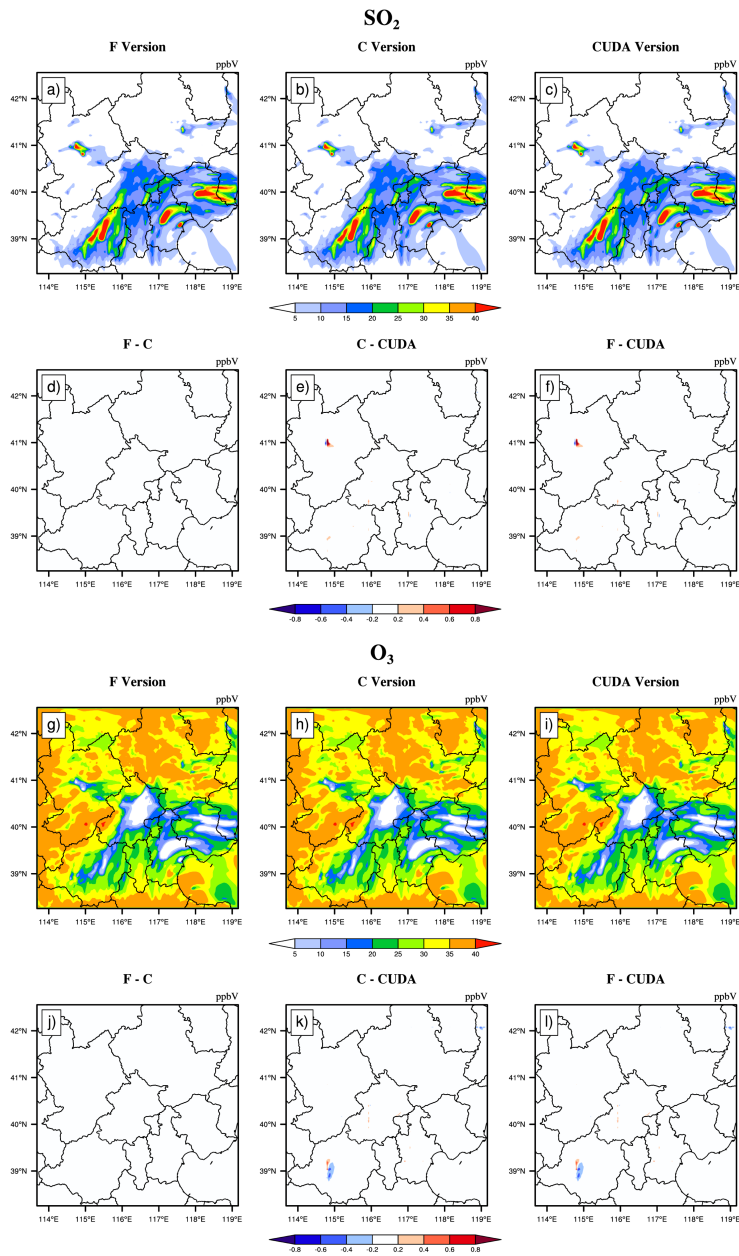
416 4.2. Error analysis

417 The hourly concentration of different CAMx simulations (Fortran, C, and CUDA
418 C versions) are compared to verify the usefulness of the CUDA C version of CAMx for
419 the numerical precision of scientific usage. Here, we chose six major species, i.e., SO₂,
420 O₃, NO₂, CO, H₂O₂ and PSO₄ after 48h integration to verify the results. Due to the
421 differences in programming languages and hardware, the simulation results are affected
422 during the porting process. Figure 5~7 present the spatial distribution of SO₂, O₃, NO₂,
423 CO, H₂O₂ and PSO₄, as well as the absolute errors (AEs) of their concentrations from
424 different CAMx versions. The species' spatial patterns of the three CAMx versions are
425 visually very similar. Especially between the Fortran and C versions, the AEs in all grid
426 boxes are in the range of ± 0.01 ppbV (the unit of PSO₄ is $\mu\text{g} \cdot \text{m}^{-3}$). During the porting
427 process, the primary error comes from converting standard C to CUDA C, and the main
428 reason was related to the hardware difference between the CPU and GPU. Due to the
429 slight difference in data operation and accuracy between CPU and GPU
430 (NVIDIA,2023), the concentration variable of hadvppm program appears to have
431 minimal negative values (about $-10^{-9} \sim -10^{-4}$) when integrating on GPU. In order to
432 allow the program to continue running, we forcibly replace these negative values with

删除的内容: ↵

435 10^{-9} . The absolute errors between the simulation results are caused by the negative
436 values are replaced by positive values. In general, for SO₂, O₃, NO₂, H₂O₂ and PSO₄,
437 the AEs in the majority of grid boxes are in the range of ± 0.8 ppbV or $\mu\text{g} \cdot \text{m}^{-3}$
438 between the standard C and CUDA C versions; for CO, because its background
439 concentration is higher, the AEs of standard C and CUDA C versions are outside that
440 range which falls into the range of -8 and 8 ppbV in some grid boxes and shows more
441 obvious AEs than other species.

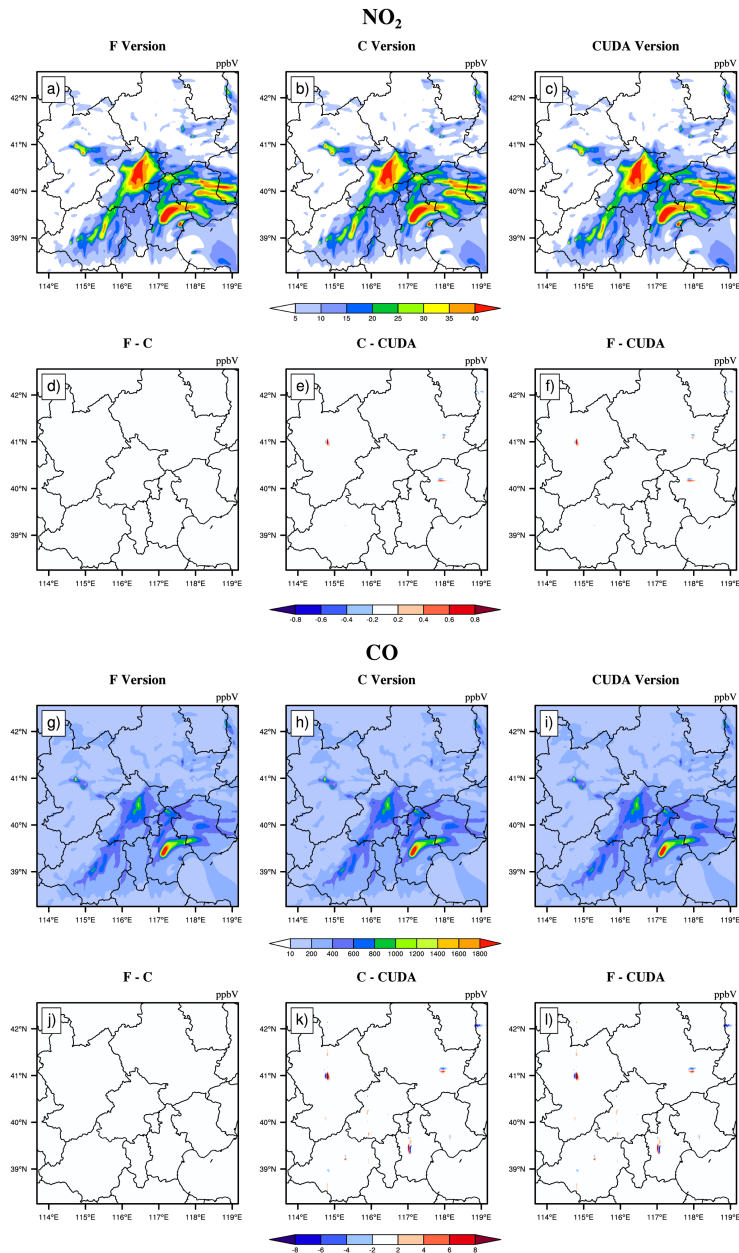
删除的内容: .



443

444 **Figure 5.** SO₂ and O₃ concentrations outputted by CAMx model for Fortran, standard C, and CUDA

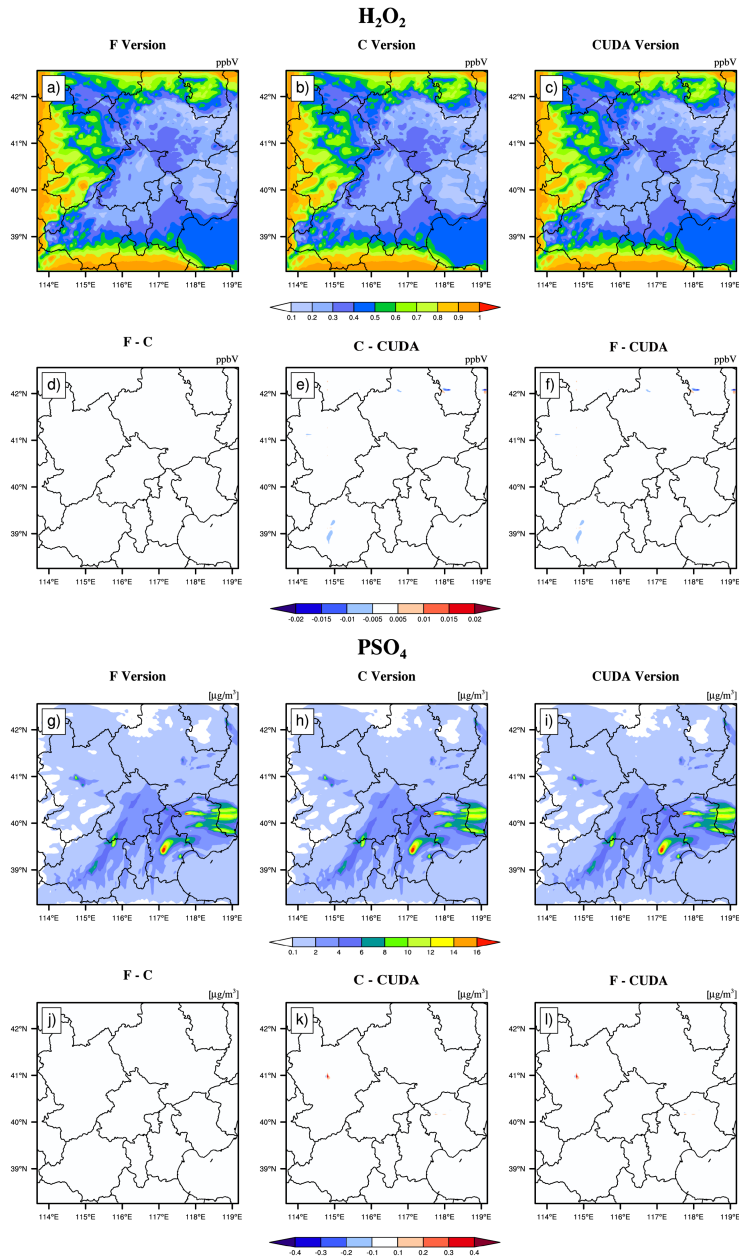
445 C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard C
446 versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output
447 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output
448 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output
449 concentration differences of Fortran and CUDA C versions.



450

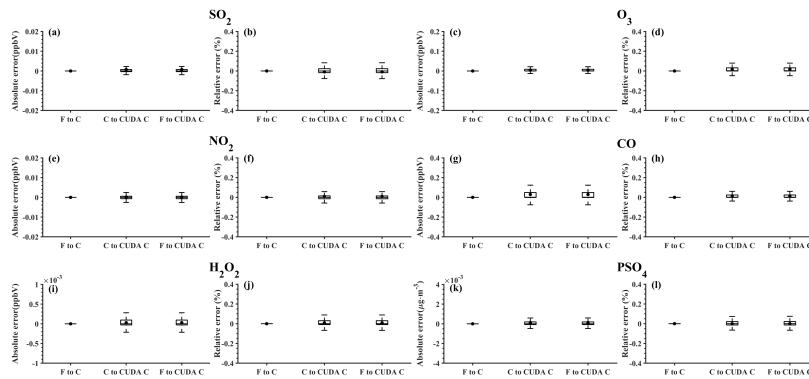
451 **Figure 6.** NO₂ and CO concentrations outputted by CAMx model for Fortran, standard C, and

452 CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard
453 C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output
454 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output
455 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output
456 concentration differences of Fortran and CUDA C versions.



458 **Figure 7.** H₂O₂ and PSO₄ concentrations outputted by CAMx model for Fortran, standard C, and
 459 CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard
 460 C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output
 461 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output
 462 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output
 463 concentration differences of Fortran and CUDA C versions.

464 Figure 8. shows the boxplot of AEs and relative error (REs) in all grid boxes for
 465 the six species during the porting process. As described above, the AEs and REs
 466 introduced by the Fortran to standard C code refactoring process are significantly small,
 467 and the primary error comes from converting standard C to CUDA C. Statistically, the
 468 average of AEs (REs) of SO₂, O₃, NO₂, CO, H₂O₂ and PSO₄ were -0.0009 ppbV (-
 469 0.01%), 0.0004 ppbV (-0.004%), 0.0005 ppbV (0.008%), 0.03 ppbV (0.01%),
 470 2.1×10^{-5} ppbV (-0.01%) and 0.0002 $\mu\text{g} \cdot \text{m}^{-3}$ (0.0023%), respectively between
 471 the Fortran and CUDA C versions. In terms of time series, the regionally averaged time
 472 series of the three versions are almost consistent (as is shown in Figure S2), and the
 473 maximum AEs for the above six species are 0.001ppbV, 0.005 ppbV, 0.002 ppbV,
 474 0.03ppbV, 0.0001 ppbV and 0.0002 $\mu\text{g} \cdot \text{m}^{-3}$, respectively, between the Fortran and
 475 CUDA C versions.



476 **Figure 8.** The distributions of absolute errors and relative errors for SO₂, O₃, NO₂, CO, H₂O₂ and
 477 PSO₄ in all of the grid boxes after 48 hours of integration.
 478

479 Figure 9. presents the regionally averaged time series and AEs of SO₂, O₃, NO₂,

480 CO, H₂O₂ and PSO₄. The time series between different versions is almost consistent,
 481 and the maximum AEs for above six species are 0.001ppbV, 0.005 ppbV, 0.002 ppbV,
 482 0.03ppbV, 0.0001 ppbV and 0.0002 $\mu\text{g} \cdot \text{m}^{-3}$, respectively between the Fortran and
 483 CUDA C versions.

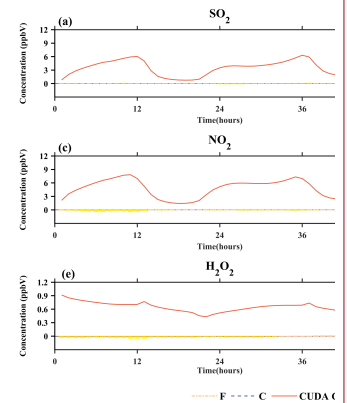
484 It is difficult to verify the scientific applicability of the results from CUDA C
 485 version because the programming language and hardware are different between the
 486 Fortran and CUDA C version. Here, we used the evaluation method of Wang et al.
 487 (2021a) to compute the root mean square errors (RMSEs) of SO₂, O₃, NO₂, CO, H₂O₂
 488 and PSO₄ between the Fortran and CUDA C versions, which are 0.0007 ppbV, 0.001
 489 ppbV, 0.0002 ppbV, 0.0005 ppbV, 0.00003 ppbV, and 0.0004 $\mu\text{g} \cdot \text{m}^{-3}$ respectively,
 490 much smaller than the spatial variation of the whole region, which is 7.0 ppbV
 491 (approximately 0.004%), 9.7 ppbV (approximately 0.003%), 7.4 ppbV (approximately
 492 0.003%), 142.2 ppbV (approximately 0.006%), 0.2ppbV (approximately 0.015%) and
 493 1.7 $\mu\text{g} \cdot \text{m}^{-3}$ (approximately 0.004%). It is indicated that the bias between CUDA C
 494 and Fortran version of the above six species is negligible compared with their own
 495 spatial changes, and the results of the CUDA C version are generally acceptable for
 496 research.

498 4.3. Offline performance comparison of GPU-HADVPPM

499 As described in the Sect. 4.2, we validate that the CAMx model result of the
 500 CUDA C version can be generally acceptable for scientific research. We tested the
 501 offline performance of the HADVPPM and GPU-HADVPPM scheme on 1 CPU core
 502 and 1 GPU card, respectively. There are 7 variables input into the HADVPPM program,
 503 which are nn, dt, dx, con, vel, area, and areav, and their specific meanings are shown in
 504 Table S1.

505 Firstly, we use random_number function in Fortran to create random single-
 506 precision floating-point numbers of different sizes for the above 7 variables, and then
 507 transmit these random numbers to the hadvppm Fortran program and hadvppm_kernel
 508 CUDA C program for computation, respectively. Finally, test the offline performance

带格式的: 下标
 带格式的: 下标
 带格式的: 下标
 带格式的: 下标
 带格式的: 下标
 带格式的: 下标



删除的内容:

删除的内容: **Figure 9.** Time series and AEs of SO₂, O₃, NO₂, CO, H₂O₂ and PSO₄ outputted by CAMx model for Fortran, standard C, and CUDA C versions.^d

513 of the HADVPPM and GPU-HADVPPM on the CPU and GPU platforms. During the
514 offline performance experiments, we used two different CPUs and GPUs described in
515 the Sect. 2.4., and the experimental results are shown in Figure 9.

删除的内容: 10

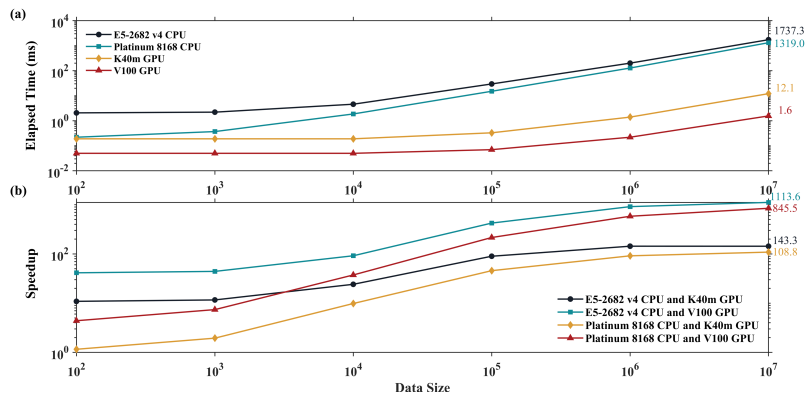
516 On the CPU platform, the wall time of hadvppm Fortran program does not change
517 significantly when the data size is less than 1000. With the increase in the data size, its
518 wall time increases linearly. When the data size reaches 10^7 , the wall time of the
519 hadvppm Fortran program on Intel Xeon E5-2682v4 and Intel Platinum 8168 CPU
520 platforms is 1737.3ms and 1319.0ms, respectively. On the GPU platform, the
521 reconstructed and extended CUDA C program implements parallel computation of
522 multiple grid points by executing a large number of kernel function copies, so the
523 computational efficiency of hadvppm_kernel CUDA C code on it is significantly
524 improved. In the size of 10^7 random numbers, the hadvppm_kernel CUDA C program
525 takes only 12.1ms and 1.6ms to complete the computation on the NVIDIA Tesla K40m
526 and NVIDIA Tesla V100 GPU.

527 Figure 9 (b) shows the speedup of HADVPPM and GPU-HADVPPM on CPU
528 platform and GPU platform under different data sizes. When mapping the HADVPPM
529 scheme to GPU, the computational efficiency under different data size is not only
530 significantly improved, but also the larger the data size, the more obvious the
531 acceleration effect of the GPU-HADVPPM. For example, in the size of 10^7 random
532 numbers, the GPU-HADVPPM achieved 1113.6x and 845.4x acceleration on the
533 NVIDIA Tesla V100 GPU, respectively, compared to the two CPU platforms. Although
534 the K40m GPU's single-card computing performance is slightly lower than that of the
535 V100 GPU, GPU-HADVPPM can also achieve up to 143.3x and 108.8x acceleration.

删除的内容: 10

536 As described in Sect. 3.2, the thread is the most basic unit of GPU for parallel
537 computing. Each dimension of the three-dimensional block can contain a maximum
538 number of threads of 1024,1024, and 64, respectively. Each dimension of the three-
539 dimensional grid can contain a maximum number of blocks of $2^{31} - 1$, 65535, and
540 65535. It is theoretically possible to distribute a large number of copies of kernel
541 functions into tens of billions of threads for parallel computing without exceeding the

544 GPU memory. In the offline performance experiments, the GPU achieved up to 10
 545 million threads of parallel computing, while the CPU can only use serial cyclic
 546 computation. Therefore, GPU-HADVPPM achieves a maximum acceleration of about
 547 1100x without I/O. In addition to this study, the GPU-based SBU-YLIN scheme in the
 548 WRF model can achieve 896x acceleration compared to the Fortran implementation
 549 running on the CPU (Mielikainen et al., 2012b).



550
 551 **Figure 2.** The offline performance of the HADVPPM and GPU-HADVPPM scheme on CPU and
 552 GPU. The unit of the wall times for the offline performance experiments is millisecond(ms).

删除的内容: 10

553 **4.4. Coupling performance comparison of GPU-HADVPPM with different GPU**
 554 **configurations**

555 **4.4.1. CAMx-CUDA on a single GPU**

556 Offline performance results show that the larger the data size, the more obvious
 557 the acceleration effect of GPU-HADVPPM scheme. After coupling the GPU-
 558 HADVPPM to CAMx without changing the advection module algorithm, the overall
 559 computational efficiency of CAMx-CUDA model is extremely low, and it takes about
 560 621 minutes to complete one-hour integration on the V100 cluster. Therefore, according
 561 to the optimization scheme in Sect. 3.2, by optimizing the algorithm of xyadvec Fortran
 562 program, we gradually increase the size of data transmitted and reduce the frequency
 563 of data transmission between CPU and GPU. When the data transmission frequency

565 between CPU and GPU is reduced to 1 within one time-step, we further optimize the
566 GPU memory access order on GPU card, eliminate unnecessary assignment loops
567 before kernel functions launched and use thread and block indices.

568 Table 4 lists the total elapsed time for different versions of CAMx-CUDA model
569 during the optimization, as described in Section 3.2. Since the xyadvec program in the
570 CAMx-CUDA V1.0 is not optimized, it is extremely computationally inefficient when
571 starting two CPU processes and configuring a GPU card for P1. On the K40m and V100
572 cluster, it takes 10829 seconds and 37237 seconds respectively to complete 1-hour
573 simulation.

574 By optimizing the algorithm of xyadvec Fortran program and hadvppm_kernel
575 CUDA C program, the frequency of data transmission between CPU and GPU was
576 decreased, and the overall computing efficiency was improved after GPU-HADVPPM
577 coupling to CAMx-CUDA model. In CAMx-CUDA V1.2, the frequency of data
578 transmission between CPU-GPU within one time step is reduced to 1, and the elapsed
579 time on the two heterogeneous clusters is 1207 seconds and 548 seconds, respectively,
580 and the speedup is 9.0x and 68.0x compared to the CAMx-CUDA V1.0.

581 GPU memory access order can directly affect the overall computational
582 efficiency of GPU-HADVPPM on the GPU. In CAMx-CUDA V1.3, we have optimized
583 the memory access order of hadvppm_kernel CUDA C program on the GPU and
584 eliminated unnecessary assignment loops before kernel functions launched, which
585 further improved the CAMx-CUDA model computational efficiency, resulting in 12.7x
586 and 94.8x speedups.

587 Using thread and block indices to compute horizontal grid points simultaneous can
588 greatly improve the computational efficiency of GPU-HADVPPM and thus reduce the
589 overall elapsed time of CAMx-CUDA model. CAMx-CUDA V1.4 further reduces the
590 elapsed time by 378 seconds and 103 seconds respectively on K40m cluster and V100
591 cluster compared with CAMx-CUDA V1.3, and achieving up to 29.0x and 128.4x
592 speedup compared with CAMx-CUDA V1.0.

593 Table 4. Total elapsed time for different versions of CAMx-CUDA during the optimization. The

删除的内容: 5

删除的内容: CAMx-V1.0

删除的内容: 5

597 unit of elapsed time for experiments is seconds (s).

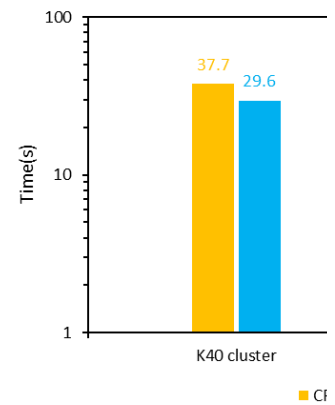
Versions	K40m cluster		V100 cluster	
	Elapsed Time	Speedup	Elapsed Time	Speedup
CAMx-CUDA V1.0	10829	1.0	37237	1.0
CAMx-CUDA V1.1	1403	7.7	1082	34.4
CAMx-CUDA V1.2	1207	9.0	548	68.0
CAMx-CUDA V1.3	751	12.7	393	94.8
CAMx-CUDA V1.4	373	29.0	290	128.4

598

599

600 In terms of the single module computational efficiency of HADVPPM and GPU-
 601 HADVPPM, we further coupling test the computational performance of the Fortran
 602 version HADVPPM on the CPU, C version HADVPPM on the CPU, and CUDA C
 603 version GPU-HADVPPM in CAMx-CUDA V1.4 (GPU-HADVPPM V1.4) on the
 604 GPU, using system_clock functions in the Fortran language and cudaEvent_t in
 605 CUDA programming. The specific results are shown in Figure 10. On the K40m
 606 cluster, it takes 37.7 seconds and 51.4 seconds to launch the Intel Xeon E5-2682 v4
 607 CPU to run Fortran and C version HADVPPM, the C version is 26.7% slower than
 608 the Fortran version. After the CUDA technology was used to convert the C code into
 609 CUDA C, the CUDA C version took 29.6 seconds to launch an NVIDIA Telsa K40m
 610 GPU to run GPU-HADVPPM V1.4, with 1.3x and 1.7x acceleration. On the V100
 611 cluster, the Fortran, the C, and the CUDA C version are computationally more
 612 efficient than those on the K40m cluster. It takes 30.1 seconds and 45.2 seconds to
 613 launch Intel Xeon Platinum 8168 CPU to run Fortran and C version HADVPPM and
 614 1.6 seconds to run the GPU-HADVPPM V1.4 using an NVIDIA V100 GPU. The
 615 computational efficiency of the CUDA C version is 18.8x and 28.3x higher than

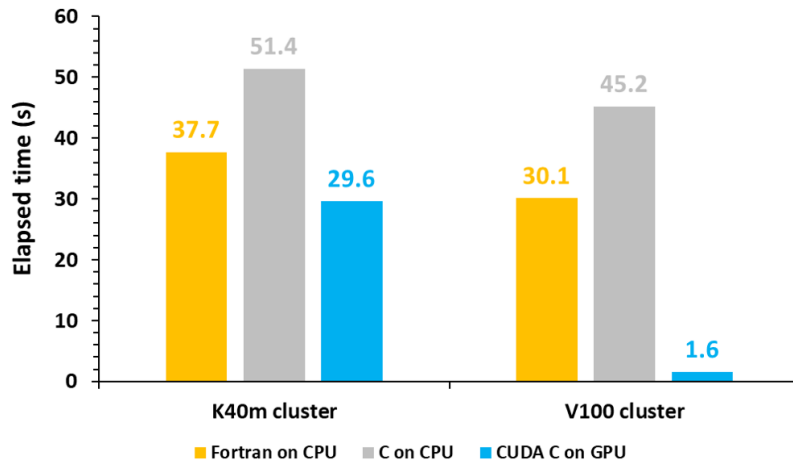
删除的内容: In terms of the single-module computational efficiency of HADVPPM and GPU-HADVPPM, we further test their performance in CPU and GPU using system_clock functions in the Fortran language and cudaEvent_t in CUDA programming. Figure 11. show the elapsed time of HADVPPM and GPU-HADVPPM in CAMX-CUDA V1.4 (GPU-HADVPPM V1.4) on K40m cluster and V100 cluster. On K40m cluster, it takes 37.7 seconds and 29.6 seconds to launch the Intel Xeon E5-2682 v4 CPU and a NVIDIA Tesla K40m GPU to run HADVPPM and GPU-HADVPPM, respectively, with 1.3x acceleration. On the V100 cluster, it takes 30.6 seconds to launch the Intel Xeon Platinum 8168 CPU to complete the HADVPPM operation, and only 1.6 seconds to run the GPU-HADVPPM using a NVIDIA V100 GPU after porting, with a speedup of 19.4x.



删除的内容: **Figure 11.** The elapsed time of HADVPPM and GPU-HADVPPM V1.4 on CPU and GPU. The unit is seconds (s).

删除的内容: 11

635 Fortran and C versions.



636

637 **Figure 10.** The elapsed time of the Fortran version HADVPPM on the CPU, the C version
638 HADVPPM on the CPU, and CUDA C version GPU-HADVPPM V1.4 on the GPU. The unit is
639 seconds (s).

640 4.4.2. CAMx-CUDA on multiple GPUs

641 To make full use of multi-core and multi-GPU in the heterogeneous cluster,
642 MPI+CUDA acceleration algorithm was implemented to improve the total
643 computational performance of the CAMx-CUDA model. Two different compile flags
644 were implemented in this study before comparing the computational efficiency of
645 CAMx-CUDA V1.3 and V1.4 on multiple GPUs, namely *-mieee-fp* and *-fp-model*
646 *precise*. The *-mieee-fp* compile flag comes from the *Makefile* of the official CAMx
647 version, which uses the IEEE standard to compare floating-point numbers. Its
648 computation accuracy is higher, but the efficiency is slower. The *-fp-model precise*
649 compile flag control the balance between precision and efficiency of floating-point
650 calculations, and it can force the compiler to use the vectorization of some calculations
651 under the value-safe. The experiment results show that *-fp model precise* compile flag
652 is 41.4% faster than *-mieee-fp*, and the AEs of the simulation results are less than
653 $\pm 0.05\text{ppbV}$ (Figure S3). Therefore, the *-fp model precise* compile flag is implemented

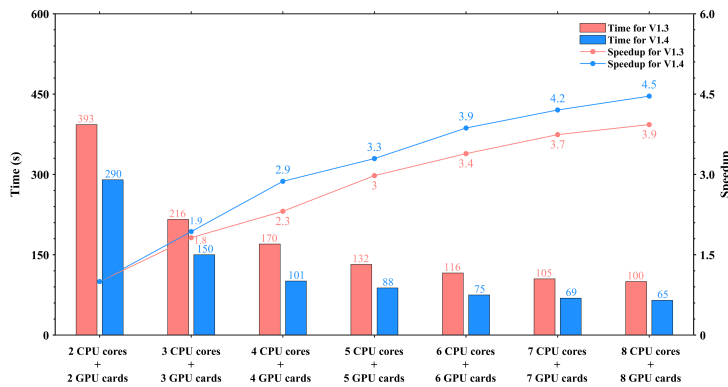
删除的内容: ↵

带格式的: 两端对齐

删除的内容: 2

656 when comparing the computational efficiency of CAMx-CUDA V1.3 and V1.4 on
 657 multiple GPU cards. Figure 11 shows the total elapsed time and speedup of CAMx-
 658 CUDA V1.3 and V1.4 on the V100 cluster. The total elapsed time decreases as the
 659 number of CPU cores and GPU cards increases. When starting 8 CPU cores and 8 GPU
 660 cards, the speedup of CAMx-CUDA V1.4 is increased from 3.9x to 4.5x compared with
 661 V1.3, and the computational efficiency is increased by 35.0%.

删除的内容: 12



662
 663 **Figure 11.** The total elapsed time and speedup of CAMx-CUDA V1.3 and V1.4 on multiple
 664 GPUs. The unit of elapsed time for experiments is seconds (s).

删除的内容: 12

665 5. Conclusions and discussion

666 GPU accelerators are playing an increasingly important role in high-performance
 667 computing. In this study, a GPU acceleration version of the PPM solver (GPU-
 668 HADVPPM) of horizontal advection for air quality model is developed, that can be run
 669 on GPU accelerators using the standard C programming language and CUDA
 670 technology. Offline performance experiments results show that K40m and V100 GPU
 671 can achieve up to 845.4x and 1113.6x speedup, respectively, and the larger the data
 672 input to the GPU, the more obvious the acceleration effect. After coupling GPU-
 673 HADVPPM to CAMx model, a series of optimization measures are taken, including

676 reducing the CPU-GPU communication frequency, increasing the size of data
677 computation on GPU, optimizing the GPU memory access order, and using thread and
678 block indices to improve the overall computing performance of CAMx-CUDA model.
679 Using a single GPU card, the optimized CAMx-CUDA V1.4 model improves the
680 computing efficiency by 29.0x and 128.4x on the K40m cluster and the V100 cluster,
681 respectively. In terms of the single-module computational efficiency of GPU-
682 HADVPPM, it can achieve 1.3x and 18.8x speedup on NVIDIA Tesla K40m GPU and
683 NVIDIA Tesla V100 GPU respectively. To make full use of multi-core and multi-GPU
684 supercomputers and further improve the total computational performance of CAMx-
685 CUDA model, a parallel architecture with an MPI+CUDA hybrid paradigm is presented.
686 After implementing the acceleration algorithm, the total elapsed time decreases as the
687 number of CPU cores and GPU cards increases, and it can achieve up to 4.5x speedup
688 when launch 8 CPU cores and 8 GPU cards compared with 2 CPU cores and 2 GPU
689 cards.

690 However, there are some limitations of the current approach which are as follows:

691 1) We currently implemented thread and block co-indexing to compute horizontal
692 grid points in parallel. Given the CAMx model 3-dimensional grid computing
693 characteristics, 3-dimensional thread and block co-indexing will be considered to
694 compute 3-dimensional grid points in parallel.

695 2) The communication bandwidth of data transfer is one of the main issues for
696 restricting the computing performance of CUDA C codes on GPUs. This restriction not
697 only holds true for GPU-HADVPPM, but also WRF module as well (Mielikainen et al.,
698 2012b; Mielikainen et al., 2013b; Huang et al., 2013). In this study, data transmission
699 efficiency between CPU and GPU is improved only by reducing communication
700 frequency. In the future, more technologies, such as pinned memory (Wang et al.,2016),
701 will be considered to solve the communication bottleneck between CPU and GPU.

702 3) In order to further improve the overall computational efficiency of the CAMx
703 model, the heterogeneous porting scheme proposed in this study will be considered to
704 carry out the heterogeneous porting of other CAMx modules in the future.

删除的内容： The communication bandwidth of data transfer is one of the main issues for restricting the computing performance of CUDA C codes on GPUs. This restriction not only holds true for GPU-HADVPPM, but also WRF module as well (Mielikainen et al., 2012b; Mielikainen et al., 2013b; Huang et al., 2013). Data transfer efficiency between CPU and GPU can be optimized. ↵

The results of this offline performance experiment shows that the larger the amount of data transferred to the GPU, the more obvious the acceleration effect. However, the number of 3D grids points in the coupling test case in this paper is only 145×157×14, a larger simulation case can be used. ↵

The computation of HADVPPM is just a small part of the whole CAMx model. When CAMx model will be completely implemented on GPU, the inputs for GPU-HADVPPM do not have to be transferred from CPU. Similarly, outputs of GPU-HADVPPM will be directly inputs to another CAMx module on GPU. Therefore, the role of I/O is greatly diminished once all of CAMx model have been converted to run on GPUs. In the future, other CAMx modules can be considered to adopt the scheme given in this paper to carry out heterogeneous porting. ↵

727

728 *Code and data availability.* The source codes of the CAMx version 6.10 are available
729 at <https://camx-wp.azurewebsites.net/download/source/> (last access: 24 March 2023,
730 ENVIRON,2022). The dataset related to this paper and CAMx-CUDA codes are
731 available online via ZENODO (<http://doi.org/10.5281/zenodo.7765218>; Cao et
732 al.,2023).

733

734 *Author contributions.* KC conducted the simulation and prepared the materials. QZW,
735 LLW, and LNW planned and organized the project. KC, QZW and XT refactored and
736 optimized code. LLW, NW, HQC, and DQL collected and prepared the data for
737 simulation. KC, QZW, XT, and LNW took part in the discussion.

738

739 *Competing interests.* The authors declare that they have no conflict of interest.

740

741 **Acknowledgements.** The National Key R&D Program of China (2020YFA0607804
742 & 2017YFC0209805), and National Supercomputing Center in Zhengzhou Innovation
743 Ecosystem Construction Technology Special Program (Grant No.201400210700), and
744 Beijing Advanced Innovation Program for Land Surface funded this work. The research
745 is support by the High Performance Scientific Computing Center (HSCC) of Beijing
746 Normal University and the National Supercomputing Center in Zhengzhou.

747

748 **Reference**

749 Bleichrodt, F., Bisseling, R. H., and Dijkstra, H. A.: Accelerating a barotropic ocean
750 model using a GPU, *Ocean Modelling*, 41, 16-21, 10.1016/j.ocemod.2011.10.001,
751 2012.

752 Cao, K., Wu, Q., Wang, L., Wang, N., Cheng, H., Tang, X., Li, D., and Wang, L.: The
753 dataset of the manuscript "GPU-HADVPPM V1.0: high-efficient parallel GPU

754 design of the Piecewise Parabolic Method (PPM) for horizontal advection in air
755 quality model (CAMx V6.10)", ZENODO,
756 <https://doi.org/10.5281/zenodo.7765218>, 2023.

757 Colella, P. and Woodward, P. R.: The Piecewise Parabolic Method (PPM) for gas-
758 dynamical simulations, *Journal of Computational Physics*, 54, 174-201,
759 [https://doi.org/10.1016/0021-9991\(84\)90143-8](https://doi.org/10.1016/0021-9991(84)90143-8), 1984.

760 ENVIRON: User Guide for Comprehensive Air Quality Model with Extensions
761 Version 6.1, available at: [https://camx-wp.azurewebsites.net/Files/CAMxUsers](https://camx-wp.azurewebsites.net/Files/CAMxUsersGuide_v6.10.pdf)
762 [Guide v6.10.pdf](https://camx-wp.azurewebsites.net/Files/CAMxUsersGuide_v6.10.pdf) (last access: 19 December 2022), 2014

763 Govett, M., Rosinski, J., Middlecoff, J., Henderson, T., Lee, J., MacDonald, A., Wang,
764 N., Madden, P., Schramm, J., and Duarte, A.: Parallelization and Performance of
765 the NIM Weather Model on CPU, GPU, and MIC Processors, *Bulletin of the*
766 *American Meteorological Society*, 98, 2201-2213, 10.1175/bams-d-15-00278.1,
767 2017.

768 Houyoux, M. R. and Vukovich, J. M.: Updates to the Sparse Matrix Operator Kernel
769 Emissions (SMOKE) Modeling System and Integration with Models-3,

770 Huang, B., Mielikainen, J., Plaza, A. J., Huang, B., Huang, A. H. L., and Goldberg, M.
771 D.: GPU acceleration of WRF WSM5 microphysics, *High-Performance*
772 *Computing in Remote Sensing*, 10.1117/12.901826, 2011.

773 Huang, B., Huang, M., Mielikainen, J., Huang, B., Huang, H. L. A., Goldberg, M. D.,
774 and Plaza, A. J.: On the acceleration of Eta Ferrier Cloud Microphysics Scheme in
775 the Weather Research and Forecasting (WRF) model using a GPU, *High-*
776 *Performance Computing in Remote Sensing II*, 10.1117/12.976908, 2012.

777 Huang, M., Huang, B., Chang, Y.-L., Mielikainen, J., Huang, H.-L. A., and Goldberg,
778 M. D.: Efficient Parallel GPU Design on WRF Five-Layer Thermal Diffusion
779 Scheme, *IEEE Journal of Selected Topics in Applied Earth Observations and*
780 *Remote Sensing*, 8, 2249-2259, 10.1109/jstars.2015.2422268, 2015.

781 Huang, M., Huang, B., Mielikainen, J., Huang, H. L. A., Goldberg, M. D., and Mehta,
782 A.: Further Improvement on GPU-Based Parallel Implementation of WRF 5-Layer

783 Thermal Diffusion Scheme, 2013 International Conference on Parallel and
784 Distributed Systems, 10.1109/icpads.2013.126, 2013.

785 Jiang, J., Lin, P., Wang, J., Liu, H., Chi, X., Hao, H., Wang, Y., Wang, W., and Zhang,
786 L.: Porting LASG/ IAP Climate System Ocean Model to Gpus Using OpenAcc,
787 IEEE Access, 7, 154490-154501, 10.1109/access.2019.2932443, 2019.

788 Mielikainen, J., Huang, B., Huang, H.-L. A., and Goldberg, M. D.: GPU Acceleration
789 of the Updated Goddard Shortwave Radiation Scheme in the Weather Research
790 and Forecasting (WRF) Model, IEEE Journal of Selected Topics in Applied Earth
791 Observations and Remote Sensing, 5, 555-562, 10.1109/jstars.2012.2186119,
792 2012a.

793 Mielikainen, J., Huang, B., Huang, H.-L. A., and Goldberg, M. D.: GPU
794 Implementation of Stony Brook University 5-Class Cloud Microphysics Scheme
795 in the WRF, IEEE Journal of Selected Topics in Applied Earth Observations and
796 Remote Sensing, 5, 625-633, 10.1109/jstars.2011.2175707, 2012b.

797 Mielikainen, J., Huang, B., Huang, H. L. A., Goldberg, M. D., and Mehta, A.: Speeding
798 Up the Computation of WRF Double-Moment 6-Class Microphysics Scheme with
799 GPU, Journal of Atmospheric and Oceanic Technology, 30, 2896-2906,
800 10.1175/jtech-d-12-00218.1, 2013a.

801 Mielikainen, J., Huang, B., Wang, J., Allen Huang, H. L., and Goldberg, M. D.:
802 Compute unified device architecture (CUDA)-based parallelization of WRF
803 Kessler cloud microphysics scheme, Computers & Geosciences, 52, 292-299,
804 10.1016/j.cageo.2012.10.006, 2013b.

805 NVIDIA: CUDA C++ Programming Guide Version 10.2, available at:
806 [https://docs.nvidia.com/cuda/archive/10.2/pdf/CUDA_C_Programming_Guide.p](https://docs.nvidia.com/cuda/archive/10.2/pdf/CUDA_C_Programming_Guide.pdf)
807 [df](https://docs.nvidia.com/cuda/archive/10.2/pdf/CUDA_C_Programming_Guide.pdf) (last access: 19 December 2022), 2020

808 [NVIDIA: Floating Point and IEEE 754 Compliance for NVIDIA GPUs. Release 12.1,](https://docs.nvidia.com/cuda/floating-point/#differences-from-x86)
809 [available at: https://docs.nvidia.com/cuda/floating-point/#differences-from-x86](https://docs.nvidia.com/cuda/floating-point/#differences-from-x86)
810 [\(last access: 18 May 2023\), 2023.](https://docs.nvidia.com/cuda/floating-point/#differences-from-x86)

811 Odman, M. and Ingram, C.: Multiscale Air Quality Simulation Platform (MAQSIP):

812 Source Code Documentation and Validation, 1996.

813 Price, E., Mielikainen, J., Huang, M., Huang, B., Huang, H.-L. A., and Lee, T.: GPU-
814 Accelerated Longwave Radiation Scheme of the Rapid Radiative Transfer Model
815 for General Circulation Models (RRTMG), IEEE Journal of Selected Topics in
816 Applied Earth Observations and Remote Sensing, 7, 3660-3667,
817 10.1109/jstars.2014.2315771, 2014.

818 Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Barker, D.M., Duda, M. G.,
819 Huang, X. Y., Wang, W., and Powers, J. G.: A Description of the Advanced
820 Research WRF Version3 (No.NCAR/TN-475CSTR), University Corporation for
821 Atmospheric Research, <https://doi.org/10.5065/D68S4MVH>, NCAR, 2008.

822 Streets, D. G., Zhang, Q., Wang, L., He, K., Hao, J., Wu, Y., Tang, Y., and Carmichael,
823 G. R.: Revisiting China's CO emissions after the Transport and Chemical
824 Evolution over the Pacific (TRACE-P) mission: Synthesis of inventories,
825 atmospheric modeling, and observations, Journal of Geophysical Research:
826 Atmospheres, 111, <https://doi.org/10.1029/2006JD007118>, 2006.

827 Streets, D. G., Bond, T. C., Carmichael, G. R., Fernandes, S. D., Fu, Q., He, D., Klimont,
828 Z., Nelson, S. M., Tsai, N. Y., Wang, M. Q., Woo, J. H., and Yarber, K. F.: An
829 inventory of gaseous and primary aerosol emissions in Asia in the year 2000,
830 Journal of Geophysical Research: Atmospheres, 108,
831 <https://doi.org/10.1029/2002JD003093>, 2003.

832 Sun, Y., Wu, Q., Wang, L., Zhang, B., Yan, P., Wang, L., Cheng, H., Lv, M., Wang, N.,
833 and Ma, S.: Weather Reduced the Annual Heavy Pollution Days after 2016 in
834 Beijing, Sola, 18, 135-139, 10.2151/sola.2022-022, 2022.

835 Wahib, M. and Maruyama, N.: Highly optimized full GPU-acceleration of non-
836 hydrostatic weather model SCALE-LES, 2013 IEEE International Conference on
837 Cluster Computing (CLUSTER), 23-27 Sept. 2013, 1-8,
838 10.1109/CLUSTER.2013.6702667,

839 Wang, P., Jiang, J., Lin, P., Ding, M., Wei, J., Zhang, F., Zhao, L., Li, Y., Yu, Z., Zheng,
840 W., Yu, Y., Chi, X., and Liu, H.: The GPU version of LASG/IAP Climate System

841 Ocean Model version 3 (LICOM3) under the heterogeneous-compute interface for
842 portability (HIP) framework and its large-scale application, *Geosci. Model Dev.*,
843 14, 2781-2799, 10.5194/gmd-14-2781-2021, 2021a.

844 Wang, Y., Guo, M., Zhao, Y., and Jiang, J.: GPUs-RRTMG_LW: high-efficient and
845 scalable computing for a longwave radiative transfer model on multiple GPUs,
846 *The Journal of Supercomputing*, 77, 4698-4717, 10.1007/s11227-020-03451-3,
847 2021b.

848 Wang, Z., Wang, Y., Wang, X., Li, F., Zhou, C., Hu, H., and Jiang, J.: GPU-
849 RRTMG_SW: Accelerating a Shortwave Radiative Transfer Scheme on GPU,
850 *IEEE Access*, 9, 84231-84240, 10.1109/access.2021.3087507, 2016.

851 Xiao, H., Lu, Y., Huang, J., and Xue, W.: An MPI+OpenACC-based PRM scalar
852 advection scheme in the GRAPES model over a cluster with multiple CPUs and
853 GPUs, *Tsinghua Science and Technology*, 27, 164-173,
854 10.26599/TST.2020.9010026, 2022.

855 Xu, S., Huang, X., Oey, L. Y., Xu, F., Fu, H., Zhang, Y., and Yang, G.: POM.gpu-v1.0:
856 a GPU-based Princeton Ocean Model, *Geoscientific Model Development*, 8,
857 2815-2827, 10.5194/gmd-8-2815-2015, 2015.

858 Zhang, Q., Streets, D. G., Carmichael, G. R., He, K. B., Huo, H., Kannari, A., Klimont,
859 Z., Park, I. S., Reddy, S., Fu, J. S., Chen, D., Duan, L., Lei, Y., Wang, L. T., and
860 Yao, Z. L.: Asian emissions in 2006 for the NASA INTEX-B mission, *Atmos.*
861 *Chem. Phys.*, 9, 5131-5153, 10.5194/acp-9-5131-2009, 2009.

