



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 Cao<sup>1</sup>, Qizhong Wu<sup>1</sup>, Lingling Wang<sup>2</sup>, Nan Wang<sup>2</sup>, Huaqiong Cheng<sup>1</sup>, Xiao**  
5 **Tang<sup>3</sup>, Dongqing Li<sup>1</sup>, and Lanning Wang<sup>1</sup>**

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

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

10 <sup>3</sup>State 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](mailto:wqizhong@bnu.edu.cn)); Lingling  
15 Wang([928216422@qq.com](mailto:928216422@qq.com)); Lanning Wang ([wangln@bnu.edu.cn](mailto:wangln@bnu.edu.cn))  
16

17 **Abstract.** With semiconductor technology gradually approaching its physical and thermal limits,  
18 Graphics processing unit (GPU) is becoming an attractive solution in many scientific applications  
19 due to their high performance. This paper presents an application of GPU accelerators in air quality  
20 model. We endeavor to demonstrate an approach that runs a PPM solver of horizontal advection  
21 (HADVPPM) for air quality model CAMx on GPU clusters. Specifically, we first convert the  
22 HADVPPM to a new Compute Unified Device Architecture C (CUDA C) code to make it  
23 computable on the GPU (GPU-HADVPPM). Then, a series of optimization measures are taken,  
24 including reducing the CPU-GPU communication frequency, increasing the size of data  
25 computation on GPU, optimizing the GPU memory access, and using thread and block indices in  
26 order to improve the overall computing performance of CAMx model coupled with GPU-  
27 HADVPPM (named as CAMx-CUDA model). Finally, a heterogeneous, hybrid programming  
28 paradigm is presented and utilized with the GPU-HADVPPM on GPU clusters with Message  
29 Passing Interface (MPI) and CUDA. Offline experiment results show that running GPU-  
30 HADVPPM on one NVIDIA Tesla K40m and NVIDIA Tesla V100 GPU can achieve up to 845.4x  
31 and 1113.6x acceleration. By implementing a series of optimization schemes, the CAMx-CUDA  
32 model resulted in a 29.0x and 128.4x improvement in computational efficiency using a GPU



33 accelerator card on a K40m and V100 cluster, respectively. In terms of the single-module  
34 computational efficiency of GPU-HADVPPM, it can achieve 1.3x and 19.4x speedup on NVIDIA  
35 Tesla K40m GPU and NVIDA Tesla V100 GPU respectively. The multi-GPU acceleration algorithm  
36 enables 4.5x speedup with 8 CPU cores and 8 GPU accelerators on V100 cluster.

## 37 **1. Introduction**

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

56 The GPU has proven successful in weather models such as Non-Hydrostatic  
57 Icosahedral Model (NIM; Govett et al., 2017), Global/Regional Assimilation and  
58 Prediction System (GRAPES; Xiao et al., 2022), and Weather Research and Forecasting  
59 model (WRF; Huang et al., 2011; Huang et al., 2012; Mielikainen et al., 2012a;



60 Mielikainen et al., 2012b; Mielikainen et al., 2013a ; Mielikainen et al., 2013b; Price et  
61 al., 2014; Huang et al., 2015), ocean models such as LASG/IAP Climate System Ocean  
62 Model (LICOM; Jiang et al., 2019; Wang et al., 2021a) and Princeton Ocean Model  
63 (POM; Xu et al., 2015), and the Earth System Model of Chinese Academy of Sciences  
64 (CAS-EMS; Wang et al., 2021b ; Wang et al., 2021c).

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

86 Programming a GPU accelerator can be a hard and error-prone process that  
87 requires specially designed programming methods, there are three widely used methods  
88 for porting program to GPUs as described above. The first method uses the OpenACC



89 directive (<https://www.openacc.org/>, last access: 19 December 2022) which provides a  
90 set of high-level directives that enable C/C++ and Fortran programmers to utilize  
91 accelerators. The second method uses CUDA Fortran. CUDA Fortran is a software  
92 compiler which co-developed by the Portland Group (PGI) and NVIDIA, and tool chain  
93 for building performance optimized GPU-accelerated Fortran applications targeting the  
94 NVIDIA GPU platform (<https://developer.nvidia.com/cuda-fortran>, last access: 19  
95 December 2022). CUDA C involves rewriting the entire program using standard C  
96 programming language and low-level CUDA subroutines  
97 (<https://developer.nvidia.com/cuda-toolkit>, last access: 19 December 2022) to support  
98 the NVIDIA GPU accelerator. Compared to the other two technologies, CUDA C  
99 porting scheme is the most complex, but its computational performance is the highest  
100 (Mielikainen et al., 2012b; Wahib and Maruyama, 2013; Xu et al., 2015).

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

## 110 **2. The CAMx model and experiments**

### 111 **2.1. Model description**

112 CAMx model is a state-of-the air quality model developed by Ramboll Environ  
113 (<https://www.camx.com/>, last access: 19 December 2022). CAMx version 6.10 (CAMx  
114 V6.10; ENVIRON, 2014) is chosen in this study, it simulates the emission, dispersion,  
115 chemical reaction, and removal of pollutants by marching the Eulerian continuity



116 equation forward in time for each chemical species on a system of nested three-  
 117 dimensional grids. The Eulerian continuity equation is expressed mathematically in  
 118 terrain-following height coordinates as formula (1):

$$\begin{aligned}
 119 \quad \frac{\partial c_i}{\partial t} = & -\nabla_H \cdot V_H c_i + \left[ \frac{\partial(c_i \eta)}{\partial z} - c_i \frac{\partial^2 h}{\partial z \partial t} \right] + \nabla \cdot \rho K \nabla (c_i / \rho) \\
 120 \quad & + \left. \frac{\partial c_i}{\partial t} \right|_{Emission} + \left. \frac{\partial c_i}{\partial t} \right|_{Chemistry} + \left. \frac{\partial c_i}{\partial t} \right|_{Removal} \quad (1)
 \end{aligned}$$

$$121 \quad \nabla_H \cdot \rho V_H = \frac{m^2}{A_{yz}} \frac{\partial}{\partial x} \left( \frac{u A_{yz} \rho}{m} \right) + \frac{m^2}{A_{xz}} \frac{\partial}{\partial y} \left( \frac{v A_{xz} \rho}{m} \right) \quad (2)$$

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

129 In the Fortran code implementation of HADVPPM scheme, the CAMx main  
 130 program calls the emistns program, which mainly performs the physical processes such  
 131 as emission, diffusion, advection and dry/wet deposition of pollutants. And then, the  
 132 horizontal advection program is invoked by emistns program to solve the horizontal  
 133 advection equation by using the HADVPPM scheme.

## 134 2.2. Benchmark performance experiments

135 The first step of the porting is to test the performance of CAMx benchmark version  
 136 and identify the hotspots of the model. On the Intel x86 CPU platform, we launch two  
 137 processes concurrently to run the CAMx and take advantage of the Intel Trace Analyzer  
 138 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)  
 139 [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  
 140 Intel VTune

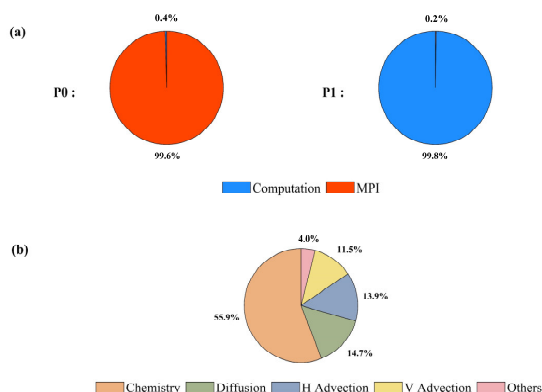


141 Profiler(VTune;<https://www.intel.com/content/www/us/en/develop/documentation/vtu>  
142 [ne-help/top.html](https://www.intel.com/content/www/us/en/develop/documentation/vtu), last access: 19 December 2022) performance analysis tools to collect  
143 performance information during CAMx operation.

144 The general MPI performance can be reported by the ITAC tool, and MPI load  
145 balance information, computation and communication profiling of each process is  
146 shown as Fig. 1a. During the running process of CAMx model, Process 0 (P0) spends  
147 99.6% of the time on the MPI\_Barrier function and only 0.4% of the time on  
148 computation, while Process 1(P1) spends 99.8% of its time computation and only 0.2%  
149 of its time receiving messages from P0. It is indicated that the parallel design of CAMx  
150 model adopts Master-Slave mode, P0 is responsible for inputting and outputting data  
151 and calling the MPI\_Barrier function to synchronize the process, so there is a lot of  
152 MPI waiting time. The other processes are responsible for computation.

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

165



166

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

170 **Table 1.** The top five most time-consuming programs on the P1 (Total runtime is 325.1 seconds).

	Program	Module	CPU time (s)
<b>Top1</b>	<i>ebirate.f</i>	Chemistry	49.4
<b>Top2</b>	<i>hadvppm.f</i>	Horizontal advection	35.6
<b>Top3</b>	<i>tridiag.f</i>	Vertical advection, Diffusion	28.0
<b>Top4</b>	<i>diffus.f</i>	Diffusion	26.9
<b>Top5</b>	<i>ebisolv.f</i>	Chemistry	26.2

171

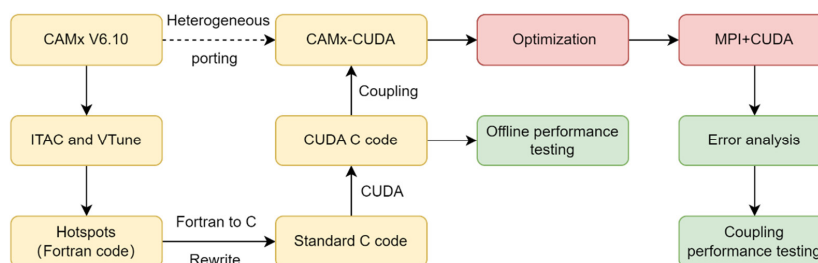
### 172 2.3. Porting scheme introduction

173 The heterogeneous scheme of CAMx-CUDA is shown in Figure 2. The second  
 174 time-consuming program *hadvppm* in CAMx model, was selected to implement the  
 175 heterogeneous porting. In order to map the *hadvppm* program to the GPU, the Fortran  
 176 code of *hadvppm* program is converted to standard C code. Then, CUDA programing  
 177 language which is tailor-made for NVIDIA was added to convert the standard C code



178 into CUDA C for data-parallel execution on GPU, as GPU-HADVPPM. It prepares the  
 179 input data for GPU-HADVPPM by constructing random numbers, and tests its offline  
 180 performance on GPU platform.

181 After coupling GPU-HADVPPM to CAMx model, the advection module code was  
 182 optimized according to the characteristics of GPU architecture to improve the overall  
 183 computational efficiency on CPU-GPU heterogeneous platform. And then, the multi-  
 184 CPU core and multi-GPU card acceleration algorithm was adopted to improve the  
 185 parallel extensibility of heterogeneous computing. Finally, the coupling performance  
 186 test is implemented after verifying the different CAMx model simulation results.



187  
 188 **Figure 2.** Heterogeneous porting scheme of CAMx-CUDA model.

189 **2.4. Hardware components and software environment of the testing system**

190 The experiments are conducted on two GPU clusters: K40m and V100.  
 191 hardware components and software environment of the two clusters are listed in Table  
 192 2. The K40m cluster is equipped with two 2.5GHz 16-core Intel Xeon E5-2682 v4 CPU  
 193 processors and one NVIDIA Tesla K40m GPU card on each node. The NVIDIA Tesla  
 194 K40m GPU has 2880 CUDA cores with 12GB of memory. The V100 cluster contains  
 195 two 2.7GHz 24-core Intel Xeon Platinum 8168 processors and eight NVIDIA Tesla  
 196 V100 GPU cards with 5120 CUDA cores and 16GB memory on each card.

197 **Table 2.** Configurations of GPU cluster.

Hardware components	
CPU	GPU





<b>K40m cluster</b>	Intel Xeon E5-2682 v4 CPU @2.5GHz, 16 cores	NVIDIA Tesla K40m, 2880 CUDA cores, 12GB memory
<b>V100 cluster</b>	Intel Xeon Platinum 8168 CPU @2.7 GHz, 24 cores	NVIDIA Tesla V100, 5120 CUDA cores, 16GB memory
<b>Software environment</b>		
	<b>Compiler and MPI</b>	<b>Programming Model</b>
<b>K40m cluster</b>	Intel-2021.4.0	CUDA-10.2
<b>V100 cluster</b>	Intel-2019.1.144	CUDA-10.0

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

### 209 **3. Porting and optimization of CAMx advection module on heterogeneous** 210 **platform**

#### 211 **3.1. Mapping HADVPPM scheme to GPU**

##### 212 **3.1.1. Manual code translation from Fortran to standard C**

213 As the CAMx V6.10 code was written in Fortran 90, we rewrote the hadvppm  
 214 program from Fortran to CUDA C. As an intermediate conversion step, we refactor the  
 215 original Fortran code using standard C. During the refactoring, some considerations are  
 216 listed in Table 3:

217 (1) The subroutine name refactored with standard C must be followed by an



218 underscore identifier, which can only be recognized when Fortran calls.

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

222 (3) Variable precision types defined in standard C must be strictly consistent with  
 223 those in Fortran.

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

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

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

230 **Table 3.** Some considerations during Fortran to C refactoring.

	<b>Fortran code</b>	<b>C code</b>
<b>Function name</b>	<i>subroutine hadvppm()</i>	<i>void hadvppm()</i>
<b>Parameter passing</b>	<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>
<b>Variable precision</b>	<i>real(kind=8) x</i>	<i>double x</i>
<b>Built-in functions</b>	<i>max</i>	<i>#define Max(a, b) ((a)&gt;(b)?(a):(b))</i>
<b>Memory read and write for multidimensional array</b>	Column-major	Row-major
<b>Array subscript index</b>	Starting from any integer	Starting from 0

231

232 **3.1.2. Converting standard C code into CUDA C**

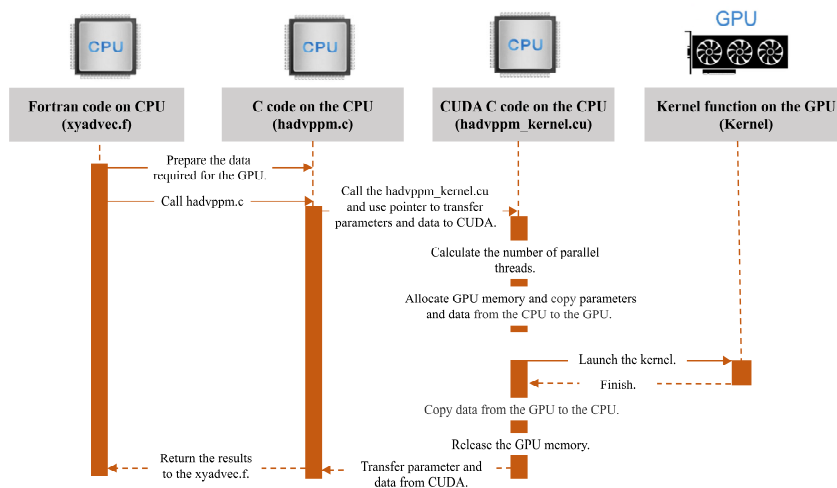
233 After refactoring the Fortran code of the hadvppm program with standard C,



234 CUDA was used to convert the C code into CUDA C to make it computable on the  
235 GPU. A standard C program using CUDA extensions distributes a large number of  
236 copies of the kernel functions into available multiprocessors and executes them  
237 simultaneously on the GPU.

238 Figure 3 shows the implementation process of the GPU-HADVPPM. As  
239 mentioned in Sect.2.1, xyadvec program calls the hadvppm program to solve the  
240 horizontal advection function. Since the rewritten CUDA program cannot be called  
241 directly by Fortran program (xyadvec.f), we add an intermediate subroutine  
242 (hadvppm.c) as an interface to transfer the parameters and data required for GPU  
243 computing from xyadvec Fortran program to hadvppm\_kernel CUDA C program.

244 A CUDA program automatically uses numerous threads on GPU to execute kernel  
245 functions. Therefore, the hadvppm\_kernel CUDA C program first calculates the  
246 number of parallel threads according to the array dimension. And then allocate GPU  
247 memory, and copy parameters and data from the CPU to the GPU. As the CUDA  
248 program launches a large number of parallel threads to execute kernel functions  
249 simultaneously, the computation results will be copied from the GPU back to the CPU.  
250 Finally, the GPU memory is released, and data computed on the GPU is returned to the  
251 xyadvec program via hadvppm C program.



252

253 **Figure 3.** The calling and computation process of the GPU-HADVPPM on the CPU-GPU



254 heterogeneous platform.

### 255 **3.2. Coupling and optimization of GPU-HADVPPM scheme on a single GPU**

256 After the hadvppm program was rewritten with standard C and CUDA, the  
257 implementation process of HADVPPM scheme is loaded from the CPU to the GPU.  
258 And then, we coupled the GPU-HADVPPM to CAMx model. For ease of description,  
259 we will refer to this original heterogeneous version of CAMx as CAMx-CUDA V1.0.  
260 In the CAMx-CUDA V1.0, four external loops are nested when hadvppm C program is  
261 called by the xyadvec program. It will result in the widespread data transfers from the  
262 CPU to the GPU over the PCIe bus within a time step, making the computation of the  
263 CAMx-CUDA V1.0 inefficient.

264 Therefore, we optimize the xyadvec Fortran program to significantly reduce the  
265 frequency of data transmission between CPU and GPU, increase the amount of data  
266 computation on GPU, and improve the total computing efficiency of the CAMx on  
267 CPU-GPU heterogeneous platforms. In the original CAMx-CUDA V1.0, four external  
268 loops outside of hadvppm C program and several one-dimensional arrays are computed  
269 before calling hadvppm C program. Then the CPU will frequently launch the GPU and  
270 transfer data to it within a time step. When the code optimization is completed, three or  
271 four-dimensional arrays required for GPU computation within a time step will be sorted  
272 before calling the hadvppm C program, and then the CPU will package and transfer the  
273 arrays to the GPU in batches. The example of xyadvec Fortran program optimization  
274 was shown in Figure S1.

275 The details of four different versions are shown in Table 4. In the CAMx-CUDA  
276 V1.0, the Fortran code of the HADVPPM scheme was rewritten using standard C and  
277 CUDA, and the xyadvec program is not optimized. The dimensions of the c1d variable  
278 array transmitted to GPU in the X and Y directions are 157 and 145 in this case,  
279 respectively. In CAMx-CUDA V1.1 and CAMx-CUDA V1.2, the c1d variable  
280 transmitted from CPU to GPU are expanded to two (about 23,000 numbers) and four  
281 dimensions (about 27.4 million numbers) by optimizing the xyadvec Fortran program



282 and hadvppm\_kernel CUDA C program, respectively.

283 The order in which data is accessed in GPU memory affects the computational  
 284 efficiency of the code. In the CAMx-CUDA V1.3 of the Table 4, we further optimized  
 285 the order in which data is accessed in GPU memory based on the order in which it is  
 286 stored in memory, and eliminated unnecessary assignment loops that were added due  
 287 to the difference in memory read order between Fortran and C.

288 As described in Sect.2.4, a thread is the basic unit of parallelism in CUDA  
 289 programming. The structure of threads is organized into a three-level hierarchy. The  
 290 highest level is a grid, which consists of three-dimensional thread blocks. The second  
 291 level is a block, which also consists of three-dimensional threads. Built-in CUDA  
 292 variable *threadIdx.x* determines a unique thread "ID" number inside a thread block.  
 293 Similarly, built-in variable *blockIdx.x* and *blockIdx.y* determine which block to execute  
 294 on, and the size of the block is determined by using the built-in variable *blockDim.x*.  
 295 For the two-dimensional horizontal grid points, many threads and blocks can be  
 296 organized so that each CUDA thread computes the results for different spatial positions  
 297 simultaneously.

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

306 **Table 4.** 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 <i>cld</i> variable, respectively.
CAMx-CUDA V1.1	Optimize <i>xyadec.f</i> and	157×145,



	<i>hadvppm_kernel.cu</i> to expand the dimension of the array transmitted to the GPU from 1-dimensional to 2-dimensional.	about 23,000 numbers for the c2d variable.
<b>CAMx-CUDA V1.2</b>	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.
<b>CAMx-CUDA V1.3</b>	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.
<b>CAMx-CUDA V1.4</b>	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.

307

308 **3.3. MPI+CUDA acceleration algorithm of CAMx-CUDA on multiple GPUs**

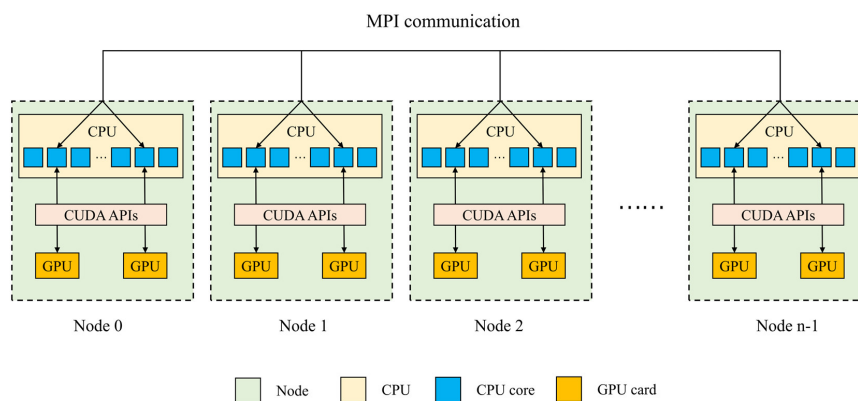
309 Generally, super-large clusters have thousands of compute nodes. The current  
 310 CAMx V6.10, implemented by adopting MPI communication technology, typically  
 311 runs on dozens of compute nodes. Once the GPU-HADVPPM is coupled into the  
 312 CAMx, it also has to run on multiple compute nodes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

320 As is shown in Figure 4., after the simulated region is subdivided by MPI, a CPU  
 321 core is responsible for the computation of a subregion. In order to improve the total  
 322 computational performance of the CAMx-CUDA model, we further used the NVIDIA  
 323 CUDA library to obtain the number of GPUs per node, and then used MPI process ID  
 324 and remainder function to determine the GPU ID to be launched by each node. Finally,



325 we used NVIDIA CUDA library cudaSetDevice to configure a GPU card for each CPU  
326 core.

327 According to the benchmark performance experiments, the parallel design of  
328 CAMx adopts Master-Slave mode, P0 is responsible for inputting and outputting data.  
329 If two processes (P0 and P1) were launched, only the P1 and its configured GPU  
330 participate in integration.



331

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

#### 334 4. Experimental results

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



#### 343 4.1. Experimental setup

344 The test case is a 48h simulation covering the Beijing, Tianjin and part region of  
345 Hebei province. The horizontal resolution is 3km with  $145 \times 157$  grid boxes. The  
346 model adopted 14 vertical layers. The simulation started at 12:00 UTC, 01 November  
347 2020, and ended at 12:00 UTC, 03 November 2020. The meteorological fields driving  
348 the CAMx model were provided by the Weather Research and Forecasting (WRF;  
349 Skamarock et al., 2008) model. The Sparse Matrix Operator Kernel Emission (SMOKE;  
350 Houyoux and Vukovich, 1999) version 2.4 model is used to provide gridded emission  
351 data for the CAMx model. The emission inventories (Sun et al., 2022) include the  
352 regional emissions in East Asia that were obtained from the Transport and Chemical  
353 Evolution over the Pacific (TRACE-P; Streets et al., 2003; Streets et al., 2006) project,  
354 30-min spatial resolution Intercontinental Chemical Transport Experiment-Phase B  
355 (INTEX-B; Zhang et al., 2009) and the updated regional emission inventories in North  
356 China. The physical and chemical numerical methods selected during CAMx model  
357 integration are listed in Table S2.

#### 358 4.2. Error analysis

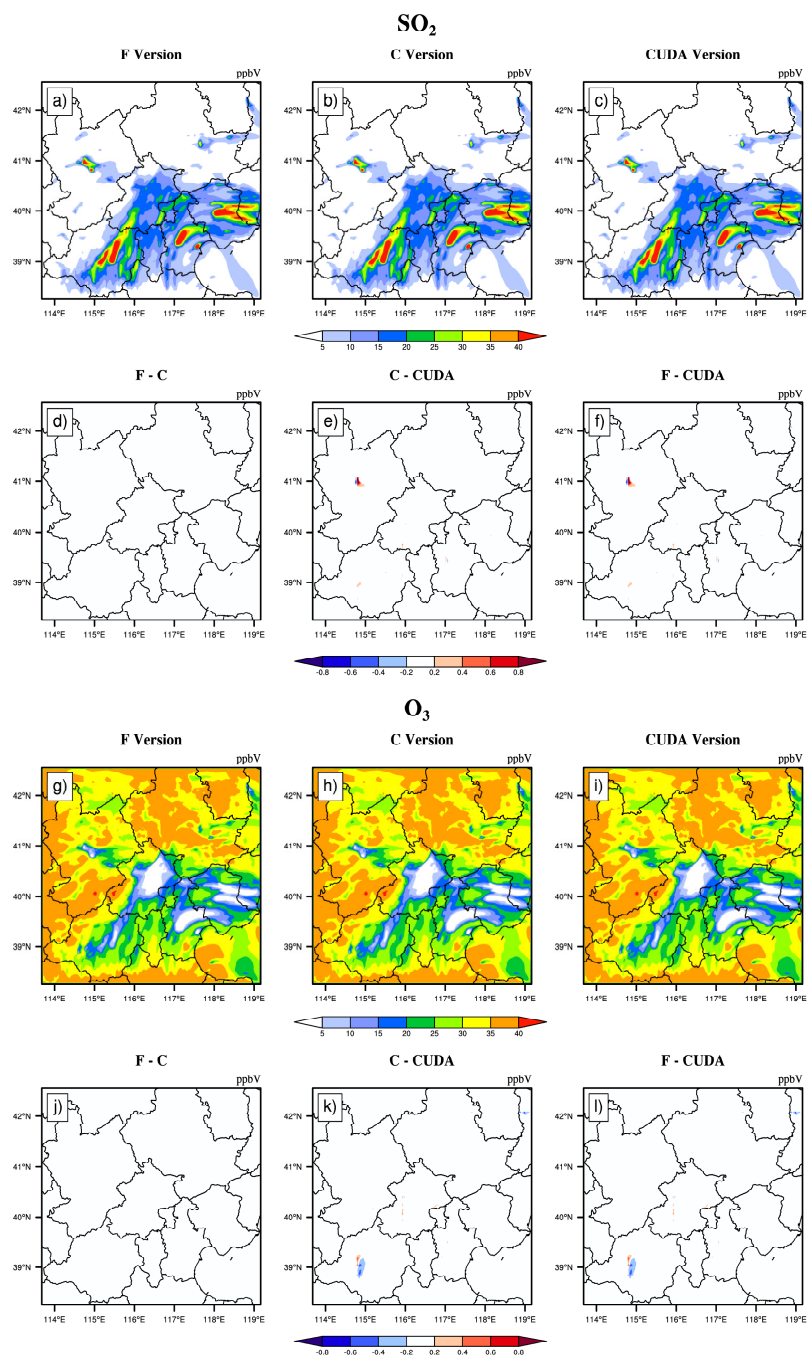
359 The hourly concentration of different CAMx simulations (Fortran, C, and CUDA  
360 C versions) are compared to verify the usefulness of the CUDA C version of CAMx for  
361 the numerical precision of scientific usage. Here, we chose six major species, i.e., SO<sub>2</sub>,  
362 O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub> after 48h integration to verify the results.

363 Due to the differences in programming languages and hardware, the simulation  
364 results are affected during the porting process. Figure 5~7 present the spatial  
365 distribution of SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub>, as well as the absolute errors (AEs)  
366 of their concentrations from different CAMx versions. The species' spatial patterns of  
367 the three CAMx versions are visually very similar. Especially between the Fortran and  
368 C versions, the AEs in all grid boxes are in the range of  $\pm 0.01$  ppbV (the unit of PSO<sub>4</sub>  
369 is  $\mu\text{g} \cdot \text{m}^{-3}$ ). During the porting process, the primary error comes from converting





370 standard C to CUDA C. In general, for SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub>, the AEs in the  
371 majority of grid boxes are in the range of  $\pm 0.8$  ppbV or  $\mu\text{g} \cdot \text{m}^{-3}$  between the  
372 standard C and CUDA C versions; for CO, because its background concentration is  
373 higher, the AEs of standard C and CUDA C versions are outside that range which falls  
374 into the range of -8 and 8 ppbV in some grid boxes and shows more obvious AEs than  
375 other species.

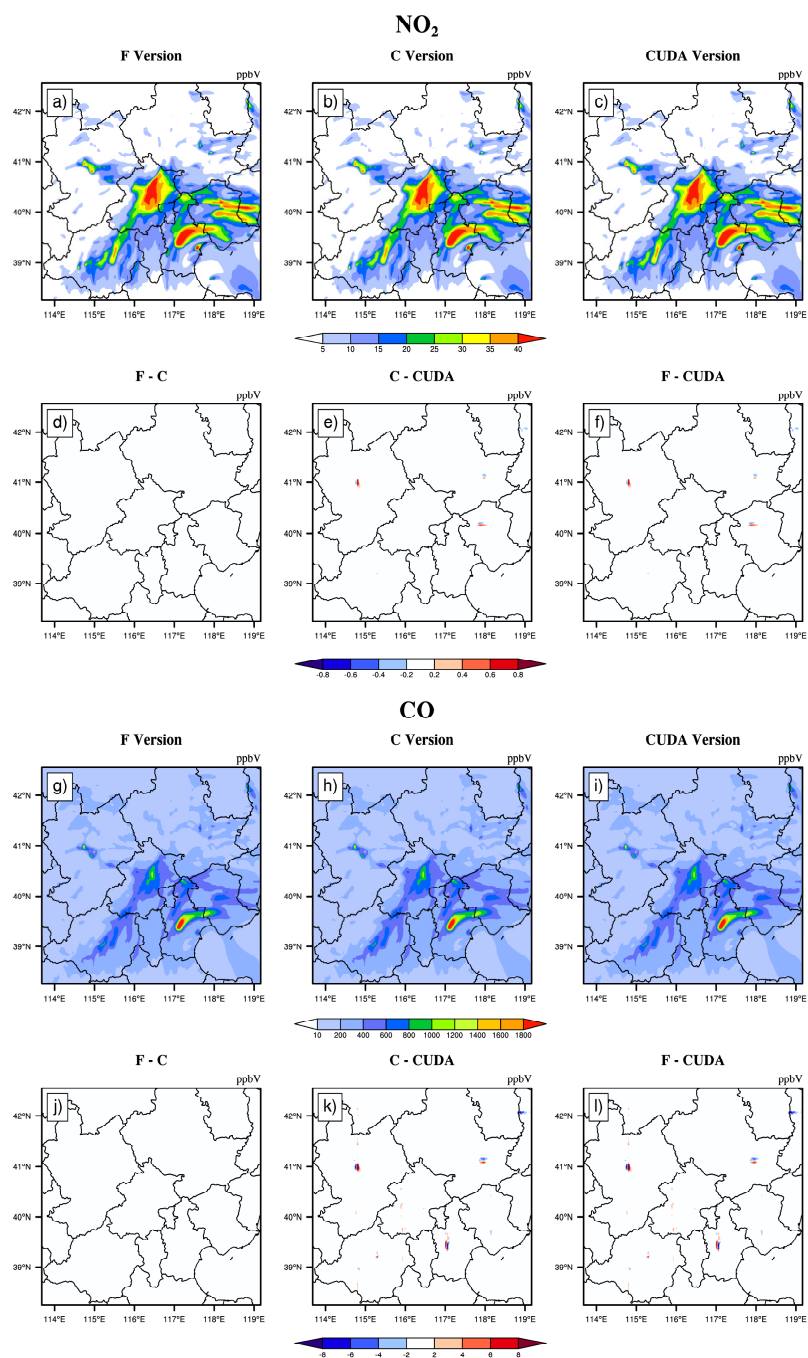


376

377 **Figure 5.** SO<sub>2</sub> and O<sub>3</sub> concentrations outputted by CAMx model for Fortran, standard C, and CUDA



378 C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard C  
379 versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output  
380 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output  
381 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output  
382 concentration differences of Fortran and CUDA C versions.

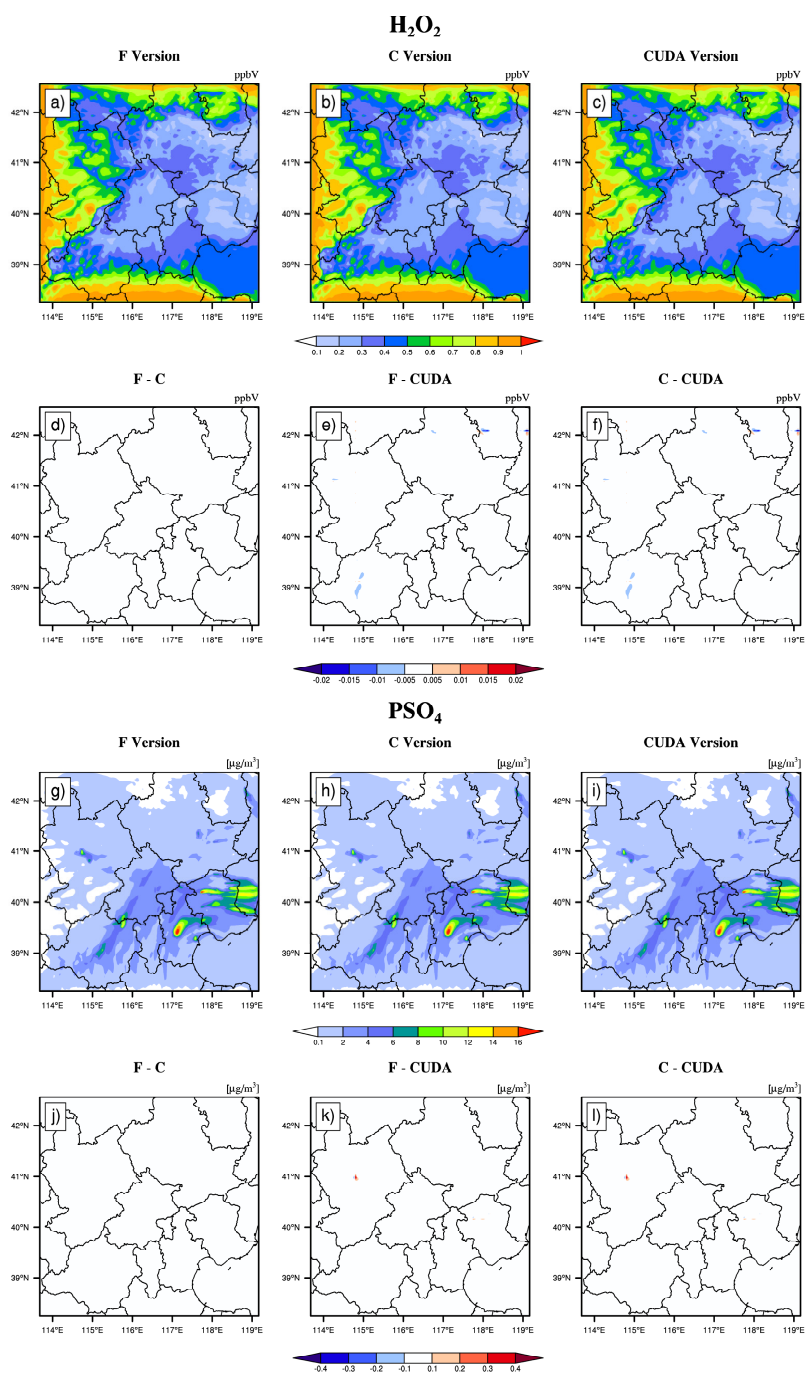


383

384 **Figure 6.** NO<sub>2</sub> and CO concentrations outputted by CAMx model for Fortran, standard C, and



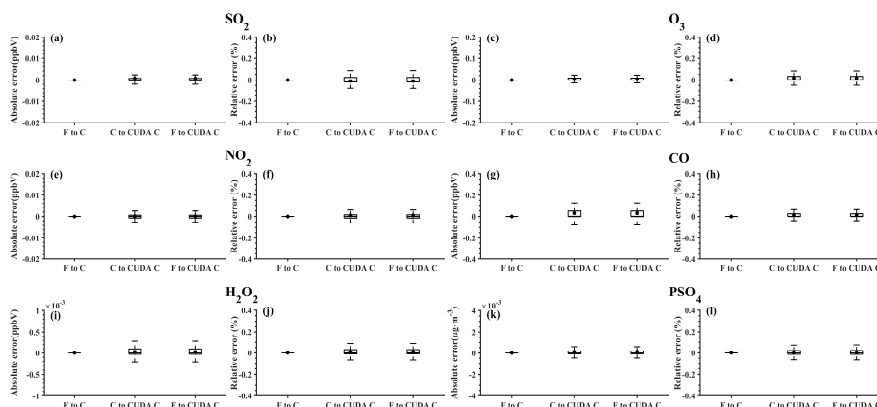
385    CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard  
386    C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output  
387    concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output  
388    concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output  
389    concentration differences of Fortran and CUDA C versions.





391 **Figure 7.** H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub> concentrations outputted by CAMx model for Fortran, standard C, and  
 392 CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard  
 393 C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output  
 394 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output  
 395 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output  
 396 concentration differences of Fortran and CUDA C versions.

397 Figure 8. shows the boxplot of AEs and relative error (REs) in all grid boxes for  
 398 the six species during the porting process. As described above, the AEs and REs  
 399 introduced by the Fortran to standard C code refactoring process are significantly small,  
 400 and the primary error comes from converting standard C to CUDA C. Statistically, the  
 401 average of AEs (REs) of SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub> were -0.0009 ppbV (-  
 402 0.01%), 0.0004 ppbV (-0.004%), 0.0005 ppbV (0.008%), 0.03 ppbV (0.01%),  
 403  $2.1 \times 10^{-5}$  ppbV (-0.01%) and 0.0002  $\mu\text{g} \cdot \text{m}^{-3}$  (0.0023%), respectively between  
 404 the Fortran and CUDA C versions.

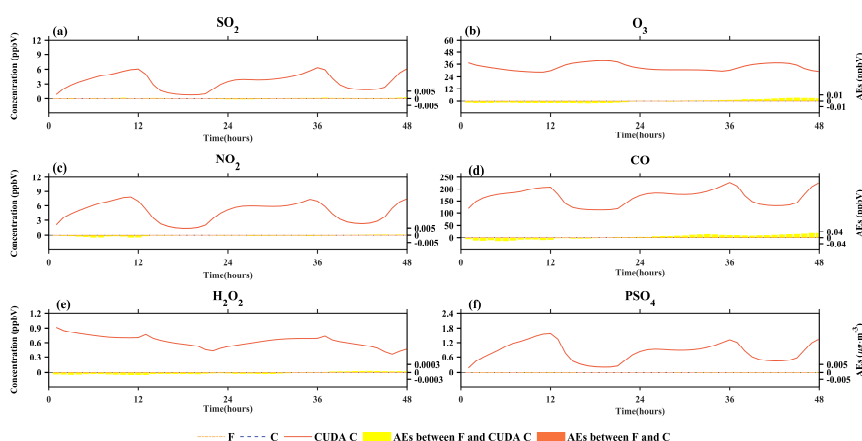


405 **Figure 8.** The distributions of absolute errors and relative errors for SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub> and  
 406 PSO<sub>4</sub> in all of the grid boxes after 48 hours of integration.

408 Figure 9. presents the regionally averaged time series and AEs of SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>,  
 409 CO, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub>. The time series between different versions is almost consistent,  
 410 and the maximum AEs for above six species are 0.001ppbV, 0.005 ppbV, 0.002 ppbV,  
 411 0.03ppbV, 0.0001 ppbV and 0.0002  $\mu\text{g} \cdot \text{m}^{-3}$ , respectively between the Fortran and  
 412 CUDA C versions.



413 It is difficult to verify the scientific applicability of the results from CUDA C  
 414 version because the programming language and hardware are different between the  
 415 Fortran and CUDA C version. Here, we used the evaluation method of Wang et al.  
 416 (2021a) to compute the root mean square errors (RMSEs) of SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub>  
 417 and PSO<sub>4</sub> between the Fortran and CUDA C versions, which are 0.0007 ppbV, 0.001  
 418 ppbV, 0.0002 ppbV, 0.0005 ppbV, 0.00003 ppbV, and 0.0004  $\mu\text{g} \cdot \text{m}^{-3}$  respectively,  
 419 much smaller than the spatial variation of the whole region, which is 7.0 ppbV  
 420 (approximately 0.004%), 9.7 ppbV (approximately 0.003%), 7.4 ppbV (approximately  
 421 0.003%), 142.2 ppbV (approximately 0.006%), 0.2ppbV (approximately 0.015%) and  
 422 1.7  $\mu\text{g} \cdot \text{m}^{-3}$  (approximately 0.004%). It is indicated that the bias between CUDA C  
 423 and Fortran version of the above six species is negligible compared with their own  
 424 spatial changes, and the results of the CUDA C version are generally acceptable for  
 425 research.



426  
 427 **Figure 9.** Time series and AEs of SO<sub>2</sub>, O<sub>3</sub>, NO<sub>2</sub>, CO, H<sub>2</sub>O<sub>2</sub> and PSO<sub>4</sub> outputted by CAMx model  
 428 for Fortran, standard C, and CUDA C versions.

### 429 4.3. Offline performance comparison of GPU-HADVPPM

430 As described in the Sect. 4.2, we validate that the CAMx model result of the  
 431 CUDA C version can be generally acceptable for scientific research. We tested the  
 432 offline performance of the HADVPPM and GPU-HADVPPM scheme on 1 CPU core





433 and 1 GPU card, respectively. There are 7 variables input into the HADVPPM program,  
434 which are nn, dt, dx, con, vel, area, and area<sub>v</sub>, and their specific meanings are shown in  
435 Table S1.

436 Firstly, we use `random_number` function in Fortran to create random single-  
437 precision floating-point numbers of different sizes for the above 7 variables, and then  
438 transmit these random numbers to the `hadvppm` Fortran program and `hadvppm_kernel`  
439 CUDA C program for computation, respectively. Finally, test the offline performance  
440 of the HADVPPM and GPU-HADVPPM on the CPU and GPU platforms. During the  
441 offline performance experiments, we used two different CPUs and GPUs described in  
442 the Sect. 2.4., and the experimental results are shown in Figure 10.

443 On the CPU platform, the wall time of `hadvppm` Fortran program does not change  
444 significantly when the data size is less than 1000. With the increase in the data size, its  
445 wall time increases linearly. When the data size reaches  $10^7$ , the wall time of the  
446 `hadvppm` Fortran program on Intel Xeon E5-2682v4 and Intel Platinum 8168 CPU  
447 platforms is 1737.3ms and 1319.0ms, respectively. On the GPU platform, the  
448 reconstructed and extended CUDA C program implements parallel computation of  
449 multiple grid points by executing a large number of kernel function copies, so the  
450 computational efficiency of `hadvppm_kernel` CUDA C code on it is significantly  
451 improved. In the size of  $10^7$  random numbers, the `hadvppm_kernel` CUDA C program  
452 takes only 12.1ms and 1.6ms to complete the computation on the NVIDIA Tesla K40m  
453 and NVIDIA Tesla V100 GPU.

454 Figure 10. (b) shows the speedup of HADVPPM and GPU-HADVPPM on CPU  
455 platform and GPU platform under different data sizes. When mapping the HADVPPM  
456 scheme to GPU, the computational efficiency under different data size is not only  
457 significantly improved, but also the larger the data size, the more obvious the  
458 acceleration effect of the GPU-HADVPPM. For example, in the size of  $10^7$  random  
459 numbers, the GPU-HADVPPM achieved 1113.6x and 845.4x acceleration on the  
460 NVIDIA Tesla V100 GPU, respectively, compared to the two CPU platforms. Although  
461 the K40m GPU's single-card computing performance is slightly lower than that of the





482 CAMx-CUDA V1.0 is not optimized, it is extremely computationally inefficient when  
 483 starting two CPU processes and configuring a GPU card for P1. On the K40m and V100  
 484 cluster, it takes 10829 seconds and 37237 seconds respectively to complete 1-hour  
 485 simulation.

486 By optimizing the algorithm of xyadvec Fortran program and hadvppm\_kernel  
 487 CUDA C program, the frequency of data transmission between CPU and GPU was  
 488 decreased, and the overall computing efficiency was improved after GPU-HADVPPM  
 489 coupling to CAMx-CUDA model. In CAMx-CUDA V1.2, the frequency of data  
 490 transmission between CPU-GPU within one time step is reduced to 1, and the elapsed  
 491 time on the two heterogeneous clusters is 1207 seconds and 548 seconds, respectively,  
 492 and the speedup is 9.0x and 68.0x compared to the CAMx-V1.0.

493 GPU memory access order can directly affect the overall computational  
 494 efficiency of GPU-HADVPPM on the GPU. In CAMx-CUDA V1.3, we have optimized  
 495 the memory access order of hadvppm\_kernel CUDA C program on the GPU and  
 496 eliminated unnecessary assignment loops before kernel functions launched, which  
 497 further improved the CAMx-CUDA model computational efficiency, resulting in 12.7x  
 498 and 94.8x speedups.

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

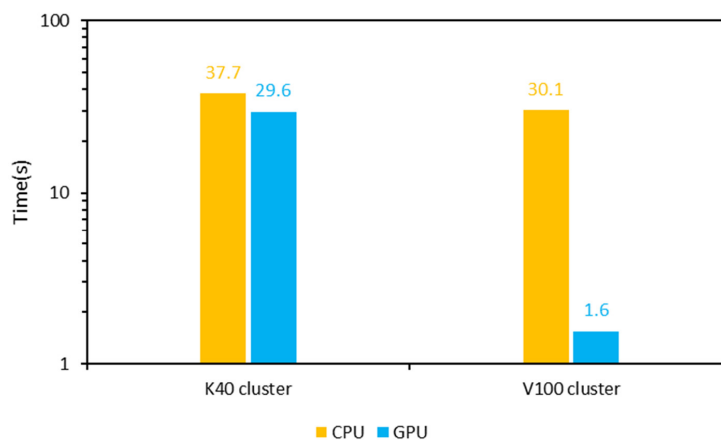
505 **Table 5.** Total elapsed time for different versions of CAMx-CUDA during the optimization. The  
 506 unit of elapsed time for experiments is seconds (s).

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

507 In terms of the single-module computational efficiency of HADVPPM and GPU-  
 508 HADVPPM, we further test their performance in CPU and GPU using system\_clock  
 509 functions in the Fortran language and cudaEvent\_t in CUDA programming. Figure 11.  
 510 show the elapsed time of HADVPPM and GPU-HADVPPM in CAMX-CUDA V1.4  
 511 (GPU-HADVPPM V1.4) on K40m cluster and V100 cluster. On K40m cluster, it takes  
 512 37.7 seconds and 29.6 seconds to launch the Intel Xeon E5-2682 v4 CPU and a NVIDIA  
 513 Tesla K40m GPU to run HADVPPM and GPU-HADVPPM, respectively, with 1.3x  
 514 acceleration. On the V100 cluster, it takes 30.6 seconds to launch the Intel Xeon  
 515 Platinum 8168 CPU to complete the HADVPPM operation, and only 1.6 seconds to run  
 516 the GPU-HADVPPM using a NVIDIA V100 GPU after porting, with a speedup of  
 517 19.4x.

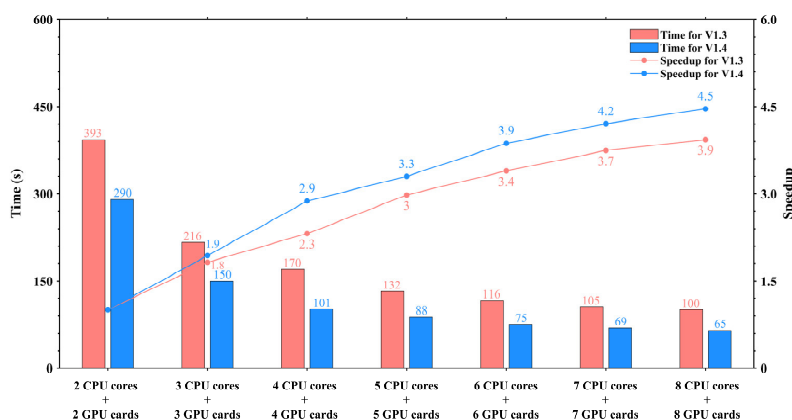


518  
 519 **Figure 11.** The elapsed time of HADVPPM and GPU-HADVPPM V1.4 on CPU and GPU. The  
 520 unit is seconds (s).



#### 521 4.4.2. CAMx-CUDA on multiple GPUs

522 To make full use of multi-core and multi-GPU in the heterogeneous cluster,  
523 MPI+CUDA acceleration algorithm was implemented to improve the total  
524 computational performance of the CAMx-CUDA model. Two different compile flags  
525 were implemented in this study before comparing the computational efficiency of  
526 CAMx-CUDA V1.3 and V1.4 on multiple GPUs, namely *-mieee-fp* and *-fp-model*  
527 *precise*. The *-mieee-fp* compile flag comes from the *Makefile* of the official CAMx  
528 version, which uses the IEEE standard to compare floating-point numbers. Its  
529 computation accuracy is higher, but the efficiency is slower. The *-fp-model precise*  
530 compile flag control the balance between precision and efficiency of floating-point  
531 calculations, and it can force the compiler to use the vectorization of some calculations  
532 under the value-safe. The experiment results show that *-fp model precise* compile flag  
533 is 41.4% faster than *-mieee-fp*, and the AEs of the simulation results are less than  
534  $\pm 0.05\text{ppbV}$  (Figure S2). Therefore, the *-fp model precise* compile flag is implemented  
535 when comparing the computational efficiency of CAMx-CUDA V1.3 and V1.4 on  
536 multiple GPU cards. Figure 12. shows the total elapsed time and speedup of CAMx-  
537 CUDA V1.3 and V1.4 on the V100 cluster. The total elapsed time decreases as the  
538 number of CPU cores and GPU cards increases. When starting 8 CPU cores and 8 GPU  
539 cards, the speedup of CAMx-CUDA V1.4 is increased from 3.9x to 4.5x compared with  
540 V1.3, and the computational efficiency is increased by 35.0%.



541

542 **Figure 12.** The total elapsed time and speedup of CAMx-CUDA V1.3 and V1.4 on multiple  
 543 GPUs. The unit of elapsed time for experiments is seconds (s).

544 **5. Conclusions and discussion**

545 GPU accelerators are playing an increasingly important role in high-performance  
 546 computing. In this study, a GPU acceleration version of the PPM solver (GPU-  
 547 HADVPPM) of horizontal advection for air quality model is developed, that can be run  
 548 on GPU accelerators using the standard C programming language and CUDA  
 549 technology. Offline performance experiments results show that K40m and V100 GPU  
 550 can achieve up to 845.4x and 1113.6x speedup, respectively, and the larger the data  
 551 input to the GPU, the more obvious the acceleration effect. After coupling GPU-  
 552 HADVPPM to CAMx model, a series of optimization measures are taken, including  
 553 reducing the CPU-GPU communication frequency, increasing the size of data  
 554 computation on GPU, optimizing the GPU memory access order, and using thread and  
 555 block indices to improve the overall computing performance of CAMx-CUDA model.  
 556 Using a single GPU card, the optimized CAMx-CUDA V1.4 model improves the  
 557 computing efficiency by 29.0x and 128.4x on the K40m cluster and the V100 cluster,  
 558 respectively. In terms of the single-module computational efficiency of GPU-



559 HADVPPM, it can achieve 1.3x and 19.4x speedup on NVIDIA Tesla K40m GPU and  
560 NVIDIA Tesla V100 GPU respectively. To make full use of multi-core and multi-GPU  
561 supercomputers and further improve the total computational performance of CAMx-  
562 CUDA model, a parallel architecture with an MPI+CUDA hybrid paradigm is presented.  
563 After implementing the acceleration algorithm, the total elapsed time decreases as the  
564 number of CPU cores and GPU cards increases, and it can achieve up to 4.5x speedup  
565 when launch 8 CPU cores and 8 GPU cards compared with 2 CPU cores and 2 GPU  
566 cards.

567 The communication bandwidth of data transfer is one of the main issues for  
568 restricting the computing performance of CUDA C codes on GPUs. This restriction not  
569 only holds true for GPU-HADVPPM, but also WRF module as well (Mielikainen et al.,  
570 2012b; Mielikainen et al., 2013b; Huang et al., 2013). Data transfer efficiency between  
571 CPU and GPU can be optimized.

572 The results of this offline performance experiment shows that the larger the  
573 amount of data transferred to the GPU, the more obvious the acceleration effect.  
574 However, the number of 3D grids points in the coupling test case in this paper is only  
575  $145 \times 157 \times 14$ , a larger simulation case can be used.

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

583

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



588 al.,2023).

589

590 **Acknowledgements.** The National Key R&D Program of China (2020YFA0607804  
591 & 2017YFC0209805), and National Supercomputing Center in Zhengzhou Innovation  
592 Ecosystem Construction Technology Special Program (Grant No.201400210700), and  
593 Beijing Advanced Innovation Program for Land Surface funded this work. The research  
594 is support by the High Performance Scientific Computing Center (HSCC) of Beijing  
595 Normal University and the National Supercomputing Center in Zhengzhou.

596

## 597 **Reference**

598 Bleichrodt, F., Bisseling, R. H., and Dijkstra, H. A.: Accelerating a barotropic ocean  
599 model using a GPU, *Ocean Modelling*, 41, 16-21, 10.1016/j.ocemod.2011.10.001,  
600 2012.

601 Cao, K., Wu, Q., Wang, L., Wang, N., Cheng, H., Tang, X., Li, D., and Wang, L.: The  
602 dataset of the manuscript "GPU-HADVPPM V1.0: high-efficient parallel GPU  
603 design of the Piecewise Parabolic Method (PPM) for horizontal advection in air  
604 quality model (CAMx V6.10)", ZENODO,  
605 <https://doi.org/10.5281/zenodo.7765218>, 2023.

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

609 ENVIRON: User Guide for Comprehensive Air Quality Model with Extensions  
610 Version 6.1, available at: [https://camx-wp.azurewebsites.net/Files/CAMxUser](https://camx-wp.azurewebsites.net/Files/CAMxUser_sGuide_v6.10.pdf)  
611 [sGuide\\_v6.10.pdf](https://camx-wp.azurewebsites.net/Files/CAMxUser_sGuide_v6.10.pdf) (last access: 19 December 2022), 2014

612 Govett, M., Rosinski, J., Middlecoff, J., Henderson, T., Lee, J., MacDonald, A., Wang,  
613 N., Madden, P., Schramm, J., and Duarte, A.: Parallelization and Performance of  
614 the NIM Weather Model on CPU, GPU, and MIC Processors, *Bulletin of the*





615 American Meteorological Society, 98, 2201-2213, 10.1175/bams-d-15-00278.1,  
616 2017.

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

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

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

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

630 Huang, M., Huang, B., Mielikainen, J., Huang, H. L. A., Goldberg, M. D., and Mehta,  
631 A.: Further Improvement on GPU-Based Parallel Implementation of WRF 5-Layer  
632 Thermal Diffusion Scheme, 2013 International Conference on Parallel and  
633 Distributed Systems, 10.1109/icpads.2013.126, 2013.

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

637 Mielikainen, J., Huang, B., Huang, H.-L. A., and Goldberg, M. D.: GPU Acceleration  
638 of the Updated Goddard Shortwave Radiation Scheme in the Weather Research  
639 and Forecasting (WRF) Model, IEEE Journal of Selected Topics in Applied Earth  
640 Observations and Remote Sensing, 5, 555-562, 10.1109/jstars.2012.2186119,  
641 2012a.

642 Mielikainen, J., Huang, B., Huang, H.-L. A., and Goldberg, M. D.: GPU  
643 Implementation of Stony Brook University 5-Class Cloud Microphysics Scheme



644 in the WRF, IEEE Journal of Selected Topics in Applied Earth Observations and  
645 Remote Sensing, 5, 625-633, 10.1109/jstars.2011.2175707, 2012b.

646 Mielikainen, J., Huang, B., Huang, H. L. A., Goldberg, M. D., and Mehta, A.: Speeding  
647 Up the Computation of WRF Double-Moment 6-Class Microphysics Scheme with  
648 GPU, Journal of Atmospheric and Oceanic Technology, 30, 2896-2906,  
649 10.1175/jtech-d-12-00218.1, 2013a.

650 Mielikainen, J., Huang, B., Wang, J., Allen Huang, H. L., and Goldberg, M. D.:  
651 Compute unified device architecture (CUDA)-based parallelization of WRF  
652 Kessler cloud microphysics scheme, Computers & Geosciences, 52, 292-299,  
653 10.1016/j.cageo.2012.10.006, 2013b.

654 NVIDIA: CUDA C++ Programming Guide Version 10.2, available at:  
655 [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)  
656 [df](https://docs.nvidia.com/cuda/archive/10.2/pdf/CUDA_C_Programming_Guide.pdf) (last access: 19 December 2022), 2020

657 Odman, M. and Ingram, C.: Multiscale Air Quality Simulation Platform (MAQSIP):  
658 Source Code Documentation and Validation, 1996.

659 Price, E., Mielikainen, J., Huang, M., Huang, B., Huang, H.-L. A., and Lee, T.: GPU-  
660 Accelerated Longwave Radiation Scheme of the Rapid Radiative Transfer Model  
661 for General Circulation Models (RRTMG), IEEE Journal of Selected Topics in  
662 Applied Earth Observations and Remote Sensing, 7, 3660-3667,  
663 10.1109/jstars.2014.2315771, 2014.

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

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



- 673 Streets, D. G., Bond, T. C., Carmichael, G. R., Fernandes, S. D., Fu, Q., He, D., Klimont,  
674 Z., Nelson, S. M., Tsai, N. Y., Wang, M. Q., Woo, J. H., and Yarber, K. F.: An  
675 inventory of gaseous and primary aerosol emissions in Asia in the year 2000,  
676 *Journal of Geophysical Research: Atmospheres*, 108,  
677 <https://doi.org/10.1029/2002JD003093>, 2003.
- 678 Sun, Y., Wu, Q., Wang, L., Zhang, B., Yan, P., Wang, L., Cheng, H., Lv, M., Wang, N.,  
679 and Ma, S.: Weather Reduced the Annual Heavy Pollution Days after 2016 in  
680 Beijing, *Sola*, 18, 135-139, [10.2151/sola.2022-022](https://doi.org/10.2151/sola.2022-022), 2022.
- 681 Wahib, M. and Maruyama, N.: Highly optimized full GPU-acceleration of non-  
682 hydrostatic weather model SCALE-LES, 2013 IEEE International Conference on  
683 Cluster Computing (CLUSTER), 23-27 Sept. 2013, 1-8,  
684 [10.1109/CLUSTER.2013.6702667](https://doi.org/10.1109/CLUSTER.2013.6702667),
- 685 Wang, P., Jiang, J., Lin, P., Ding, M., Wei, J., Zhang, F., Zhao, L., Li, Y., Yu, Z., Zheng,  
686 W., Yu, Y., Chi, X., and Liu, H.: The GPU version of LASG/IAP Climate System  
687 Ocean Model version 3 (LICOM3) under the heterogeneous-compute interface for  
688 portability (HIP) framework and its large-scale application, *Geosci. Model Dev.*,  
689 14, 2781-2799, [10.5194/gmd-14-2781-2021](https://doi.org/10.5194/gmd-14-2781-2021), 2021a.
- 690 Wang, Y., Guo, M., Zhao, Y., and Jiang, J.: GPUs-RRTMG\_LW: high-efficient and  
691 scalable computing for a longwave radiative transfer model on multiple GPUs,  
692 *The Journal of Supercomputing*, 77, 4698-4717, [10.1007/s11227-020-03451-3](https://doi.org/10.1007/s11227-020-03451-3),  
693 2021b.
- 694 Wang, Z., Wang, Y., Wang, X., Li, F., Zhou, C., Hu, H., and Jiang, J.: GPU-  
695 RRTMG\_SW: Accelerating a Shortwave Radiative Transfer Scheme on GPU,  
696 *IEEE Access*, 9, 84231-84240, [10.1109/access.2021.3087507](https://doi.org/10.1109/access.2021.3087507), 2021c.
- 697 Xiao, H., Lu, Y., Huang, J., and Xue, W.: An MPI+OpenACC-based PRM scalar  
698 advection scheme in the GRAPES model over a cluster with multiple CPUs and  
699 GPUs, *Tsinghua Science and Technology*, 27, 164-173,  
700 [10.26599/TST.2020.9010026](https://doi.org/10.26599/TST.2020.9010026), 2022.
- 701 Xu, S., Huang, X., Oey, L. Y., Xu, F., Fu, H., Zhang, Y., and Yang, G.: POM.gpu-v1.0:



702 a GPU-based Princeton Ocean Model, *Geoscientific Model Development*, 8,  
703 2815-2827, 10.5194/gmd-8-2815-2015, 2015.

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