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

**Kai Cao[1], Qizhong Wu[1], Lingling Wang[2], Nan Wang[2], Huaqiong Cheng[1], Xiao Tang[3], Dongqing Li[1], and Lanning Wang[1]**

[1]College of Global Change and Earth System Science, Beijing Normal University, Beijing 100875, China

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

[3]State Key Laboratory of Atmospheric Boundary Layer Physics and Atmospheric Chemistry, Institute of Atmospheric Physics, Chinese Academy of Science, Beijing 100029, China

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

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

1

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

## 1. Introduction

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

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

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

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

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

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

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

114   **2.1.   Model description**

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

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

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

$$+ \frac{\partial c_i}{\partial t} \bigg|_{Emission} + \frac{\partial c_i}{\partial t} \bigg|_{Chemistry} + \frac{\partial c_i}{\partial t} \bigg|_{Removal} \tag{1}$$

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

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

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

## 2.2. Benchmark performance experiments

The first step of the porting is to test the performance of CAMx benchmark version and identify the hotspots of the model. On the Intel x86 CPU platform, we launch two processes concurrently to run the CAMx and take advantage of the Intel Trace Analyzer
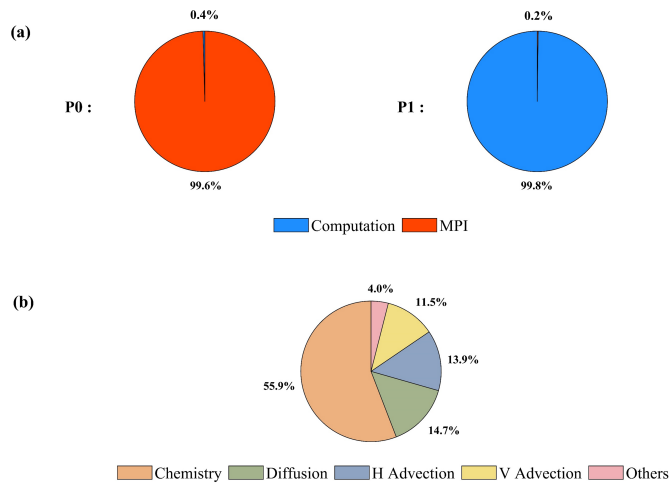
141  Collector(ITAC; https://www.intel.com/content/www/us/en/docs/trace-analyzer-
142  collector/get-started-guide/2021-4/overview.html, last access: 19 December 2022) and
143  Intel                                                                          VTune
144  Profiler(VTune;https://www.intel.com/content/www/us/en/develop/documentation/vtu
145  ne-help/top.html, last access: 19 December 2022) performance analysis tools to collect
146  performance information during CAMx operation.

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

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

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

170 WRF) also contains an advection module, so this study's heterogeneous porting method

171 and experience can be used for reference. Therefore, a GPU acceleration version of the

172 HADVPPM scheme, namely GPU-HADVPPM, is built to improve CAMx

173 performance.



174

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

178
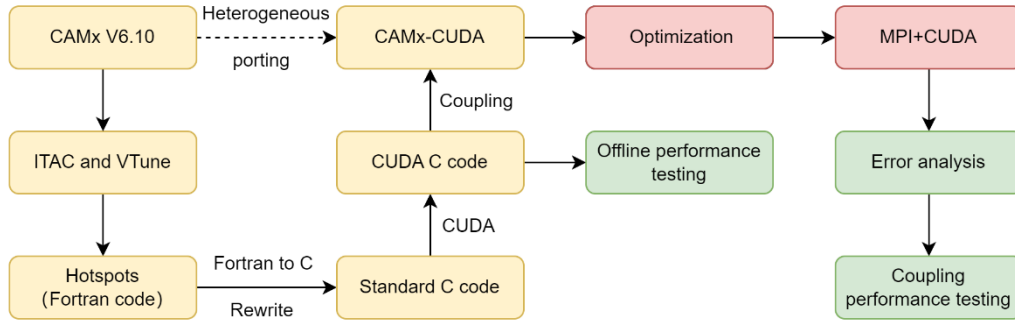
### 2.3. Porting scheme introduction

180 The heterogeneous scheme of CAMx-CUDA is shown in Figure 2. The second

181 time-consuming hadvppm program in the CAMx model was selected to implement the

182 heterogeneous porting. In order to map the hadvppm program to the GPU, the Fortran

183 code was converted to standard C code. Then, CUDA programing language, which was

184 tailor-made for NVIDIA, was added to convert the standard C code into CUDA C for

185 data-parallel execution on GPU, as GPU-HADVPPM. It prepared the input data for

186 GPU-HADVPPM by constructing random numbers and tested its offline performance

187 on the GPU platform.

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

189 optimized according to the characteristics of GPU architecture to improve the overall

190 computational efficiency on CPU-GPU heterogeneous platform. And then, the multi-

191 CPU core and multi-GPU card acceleration algorithm was adopted to improve the

192 parallel extensibility of heterogeneous computing. Finally, the coupling performance

193 test is implemented after verifying the different CAMx model simulation results.



194

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

196 **2.4. Hardware components and software environment of the testing system**

197 The experiments are conducted on two GPU clusters: K40m and V100.

198 hardware components and software environment of the two clusters are listed in Table

199 1. The K40m cluster is equipped with two 2.5GHz 16-core Intel Xeon E5-2682 v4 CPU

200 processors and one NVIDIA Tesla K40m GPU card on each node. The NVIDIA Tesla

201 K40m GPU has 2880 CUDA cores with 12GB of memory. The V100 cluster contains

202 two 2.7GHz 24-core Intel Xeon Platinum 8168 processors and eight NVIDIA Tesla

203 V100 GPU cards with 5120 CUDA cores and 16GB memory on each card.

204 **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 |

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

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

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

### 3.1. Mapping HADVPPM scheme to GPU

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

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

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

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

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

230 those in Fortran.

231      (4) Some built-in functions in Fortran are not available in standard C and need to

232 be defined in standard C macro definitions.

233      (5) For multidimensional arrays, Fortran and standard C follow column-major and

234 row-major order in-memory read and write, respectively;

235      (6) Array subscripts in Fortran and standard C are indexed from any integer and 0,

236 respectively.

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

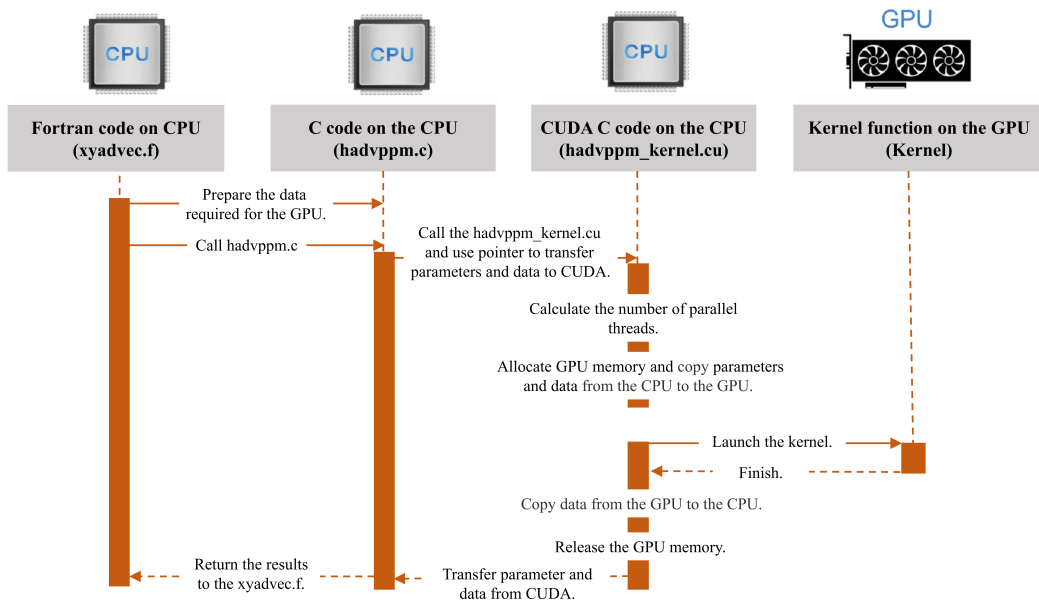| | **Fortran code** | **C code** |
|---|---|---|
| **Function name** | *subroutine hadvppm()* | *void hadvppm()* |
| **Parameter passing** | *hadvppm(nn,dt,dx,con,vel,area,areav, flxarr,mynn)* | *hadvppm(int *nn,float *dt, float *dx, float *con, float *vel, float *area, float *areav, float *flxarr,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

238

### 3.1.2.  Converting standard C code into CUDA C

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

241 CUDA was used to convert the C code into CUDA C to make it computable on the

242 GPU. A standard C program using CUDA extensions distributes a large number of

243 copies of the kernel functions into available multiprocessors and executes them

244 simultaneously on the GPU.

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

246    mentioned in Sect.2.1, xyadvec program calls the hadvppm program to solve the

247    horizontal advection function. Since the rewritten CUDA program cannot be called

248    directly by Fortran program (xyadvec.f), we add an intermediate subroutine

249    (hadvppm.c) as an interface to transfer the parameters and data required for GPU

250    computing from xyadvec Fortran program to hadvppm_kernel CUDA C program.

251    A CUDA program automatically uses numerous threads on GPU to execute kernel

252    functions. Therefore, the hadvppm_kernel CUDA C program first calculates the

253    number of parallel threads according to the array dimension. And then allocate GPU

254    memory, and copy parameters and data from the CPU to the GPU. As the CUDA

255    program launches a large number of parallel threads to execute kernel functions

256    simultaneously, the computation results will be copied from the GPU back to the CPU.

257    Finally, the GPU memory is released, and data computed on the GPU is returned to the

258    xyadvec program via hadvppm C program.



259

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

261    heterogeneous platform.

262    **3.2.    Coupling and optimization of GPU-HADVPPM scheme on a single GPU**

263    After the hadvppm program was rewritten with standard C and CUDA, the

264    implementation process of HADVPPM scheme is loaded from the CPU to the GPU.

11

265 And then, we coupled the GPU-HADVPPM to CAMx model. For ease of description,

266 we will refer to this original heterogeneous version of CAMx as CAMx-CUDA V1.0.

267 In the CAMx-CUDA V1.0, four external loops are nested when hadvppm C program is

268 called by the xyadvec program. It will result in the widespread data transfers from the

269 CPU to the GPU over the PCIe bus within a time step, making the computation of the

270 CAMx-CUDA V1.0 inefficient.

271 Therefore, we optimize the xyadvec Fortran program to significantly reduce the

272 frequency of data transmission between CPU and GPU, increase the amount of data

273 computation on GPU, and improve the total computing efficiency of the CAMx on

274 CPU-GPU heterogeneous platforms. In the original CAMx-CUDA V1.0, four external

275 loops outside of hadvppm C program and several one-dimensional arrays are computed

276 before calling hadvppm C program. Then the CPU will frequently launch the GPU and

277 transfer data to it within a time step. When the code optimization is completed, three or

278 four-dimensional arrays required for GPU computation within a time step will be sorted

279 before calling the hadvppm C program, and then the CPU will package and transfer the

280 arrays to the GPU in batches. The example of xyadvec Fortran program optimization

281 was shown in Figure S1.

282 The details of four different versions are shown in Table 3. In the CAMx-CUDA

283 V1.0, the Fortran code of the HADVPPM scheme was rewritten using standard C and

284 CUDA, and the xyadvec program is not optimized. The dimensions of the c1d variable

285 array transmitted to GPU in the X and Y directions are 157 and 145 in this case,

286 respectively. In CAMx-CUDA V1.1 and CAMx-CUDA V1.2, the c1d variable

287 transmitted from CPU to GPU are expanded to two (about 23,000 numbers) and four

288 dimensions (about 27.4 million numbers) by optimizing the xyadvec Fortran program

289 and hadvppm_kernel CUDA C program, respectively.

290 The order in which data is accessed in GPU memory affects the computational

291 efficiency of the code. In the CAMx-CUDA V1.3 of the Table 4, we further optimized

292 the order in which data is accessed in GPU memory based on the order in which it is

293 stored in memory, and eliminated unnecessary assignment loops that were added due

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

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

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

313  **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 *xyadvec.f* was not optimized. | 157 and 145 in the x direction and y direction for the c1d variable, respectively. |
| **CAMx-CUDA V1.1** | Optimize *xyadec.f* and *hadvppm_kernel.cu* 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. |

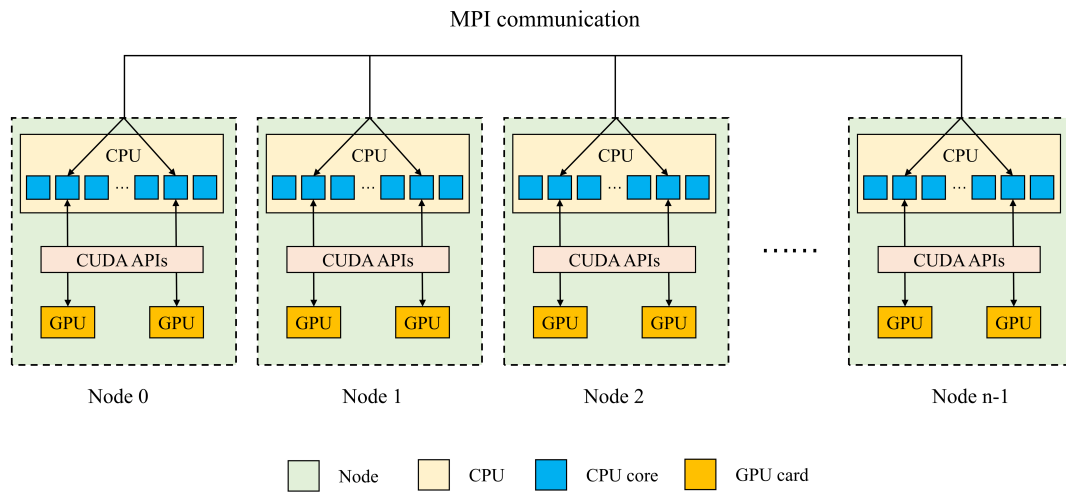| | Based on the CAMx-CUDA V1.2, the order of GPU memory access is optimized and unnecessary assignment loops are eliminated. | 157×145×14×86, about 27.4 million numbers for the c4d variable. |
|---|---|---|
| **CAMx-CUDA V1.3** | | |
| **CAMx-CUDA V1.4** | Based on the CAMx-CUDA V1.3, using thread and 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. |

314

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

315

316      Generally, super-large clusters have thousands of compute nodes. The current

317 CAMx V6.10, implemented by adopting MPI communication technology, typically

318 runs on dozens of compute nodes. Once the GPU-HADVPPM is coupled into the

319 CAMx, it also has to run on multiple compute nodes which equipped one or more GPUs

320 on each node. To make full use of multi-core and multi-GPU supercomputers and

321 further improve the overall computational performance of the CAMx-CUDA, we adopt

322 a parallel architecture with an MPI+CUDA hybrid paradigm, that is, the collaborative

323 computing strategy of multiple CPU cores and multiple GPU cards is adopted during

324 the operation of CAMx-CUDA model. Adopt this strategy, the GPU-HADVPPM can

325 run on multiple GPUs, the Fortran code of other modules in CAMx-CUDA model can

326 run on multiple CPU cores.

327      As is shown in Figure 4., after the simulated region is subdivided by MPI, a CPU

328 core is responsible for the computation of a subregion. In order to improve the total

329 computational performance of the CAMx-CUDA model, we further used the NVIDIA

330 CUDA library to obtain the number of GPUs per node, and then used MPI process ID

331 and remainder function to determine the GPU ID to be launched by each node. Finally,

332 we used NVIDIA CUDA library cudaSetDevice to configure a GPU card for each CPU

333 core.

334      According to the benchmark performance experiments, the parallel design of

335 CAMx adopts Master-Slave mode, P0 is responsible for inputting and outputting data.

336 If two processes (P0 and P1) were launched, only the P1 and its configured GPU

337    participate in integration.



MPI communication

Node 0          Node 1          Node 2                    Node n-1

Node    CPU    CPU core    GPU card

338

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

## 4.    Experimental results

The validation and evaluation of porting the HADVPPM scheme from the CPU to
the GPU platform were conducted using offline and coupling performance experiments.
First, we validated 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. Sect.4.2 and Sect.4.4, the CAMx version
of the HADVPPM scheme written by Fortran language, standard C, and CUDA C, is
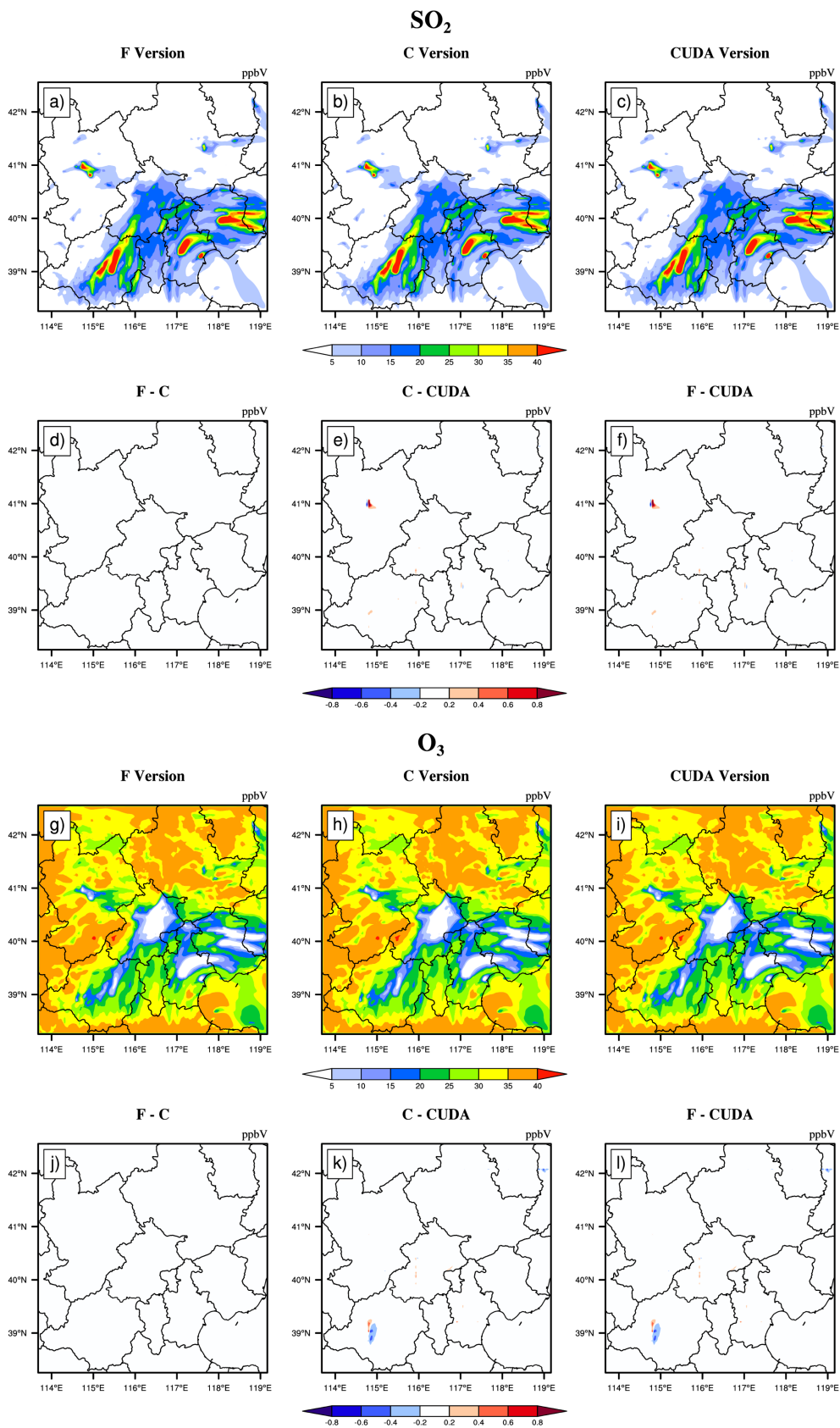named F, C, and CUDA C, respectively.

### 4.1.    Experimental setup

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

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

## 4.2. Error analysis

The hourly concentration of different CAMx simulations (Fortran, C, and CUDA C versions) are compared to verify the usefulness of the CUDA C version of CAMx for the numerical precision of scientific usage. Here, we chose six major species, i.e., $SO_2$, $O_3$, $NO_2$, CO, $H_2O_2$ and $PSO_4$ after 48h integration to verify the results. Due to the differences in programming languages and hardware, the simulation results are affected during the porting process. Figure 5~7 present the spatial distribution of $SO_2$, $O_3$, $NO_2$, CO, $H_2O_2$ and $PSO_4$, as well as the absolute errors (AEs) of their concentrations from different CAMx versions. The species' spatial patterns of the three CAMx versions are visually very similar. Especially between the Fortran and C versions, the AEs in all grid boxes are in the range of ±0.01 ppbV (the unit of $PSO_4$ is $\mu g \cdot m^{-3}$). During the porting process, the primary error comes from converting standard C to CUDA C, and the main reason was related to the hardware difference between the CPU and GPU. Due to the slight difference in data operation and accuracy between CPU and GPU (NVIDIA,2023), the concentration variable of hadvppm program appears to have minimal negative values (about $-10^{-9} \sim -10^{-4}$) when integrating on GPU. In order to allow the program to continue running, we forcibly replace these negative values with
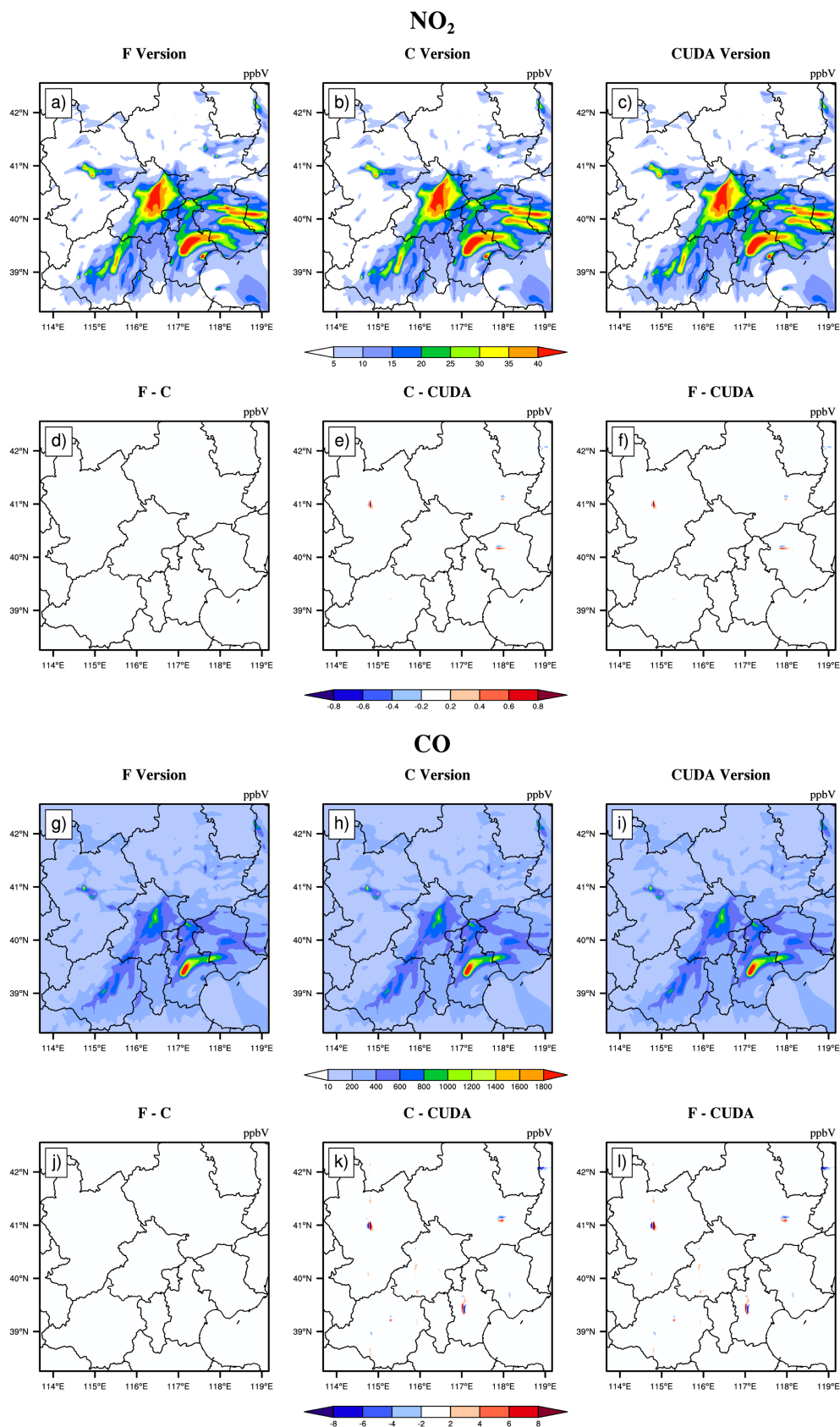
382     $10^{-9}$. It is because these negative values are replaced by positive values that the

383     simulation results are biased. In general, for $SO_2$, $O_3$, $NO_2$, $H_2O_2$ and $PSO_4$, the AEs in

384     the majority of grid boxes are in the range of $\pm 0.8$ ppbV or $\mu g \cdot m^{-3}$ between the

385     standard C and CUDA C versions; for CO, because its background concentration is

386     higher, the AEs of standard C and CUDA C versions are outside that range which falls

387     into the range of -8 and 8 ppbV in some grid boxes and shows more obvious AEs than

388     other species.

**Figure 5.** SO$_2$ and O$_3$ concentrations outputted by CAMx model for Fortran, standard C, and CUDA

391     C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard C

392     versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output

393     concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output

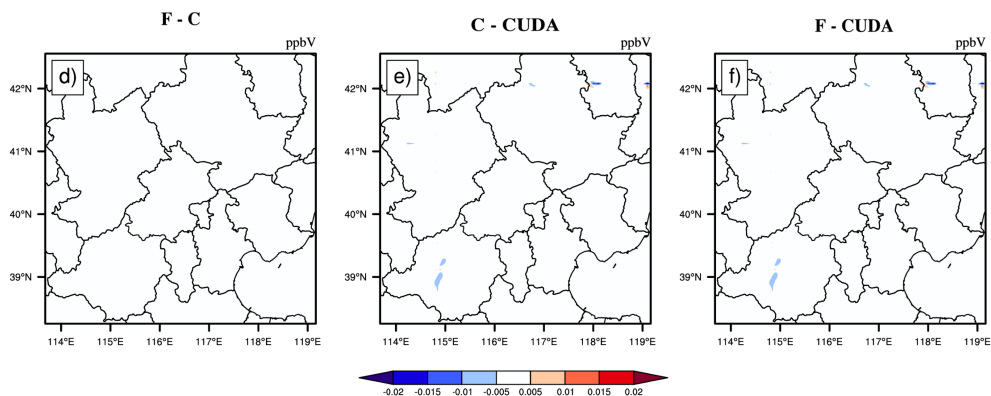394     concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output

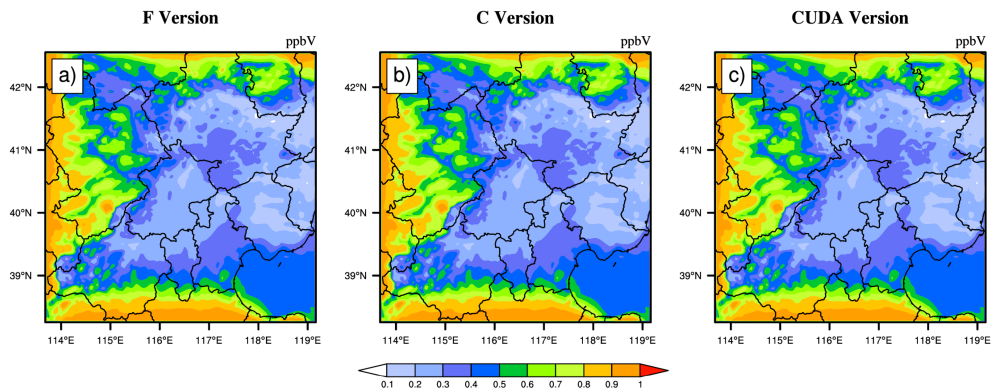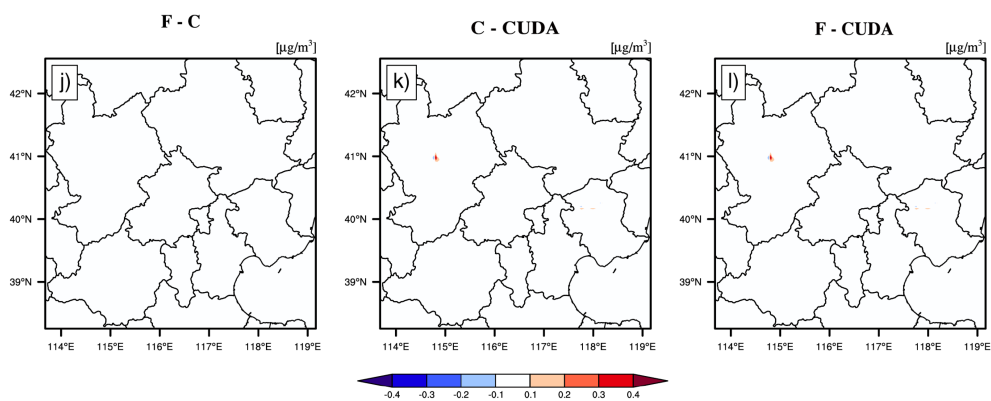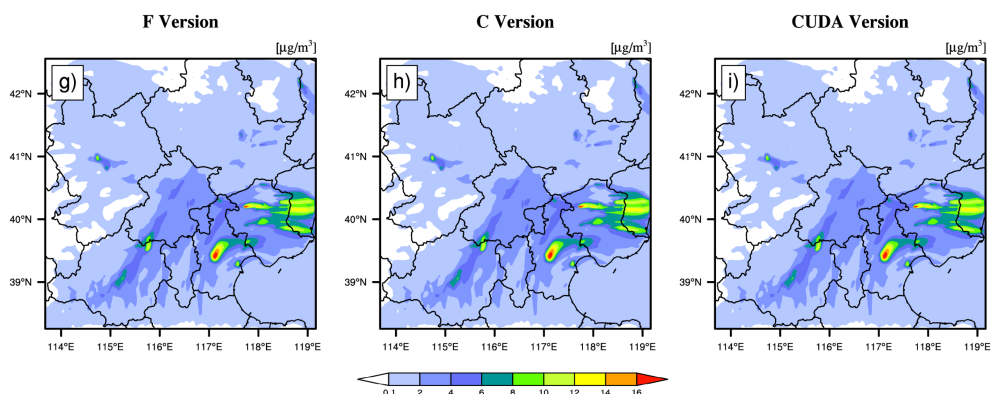395     concentration differences of Fortran and CUDA C versions.

**NO₂**

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

398     CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard

399     C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output

400     concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output

401     concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output

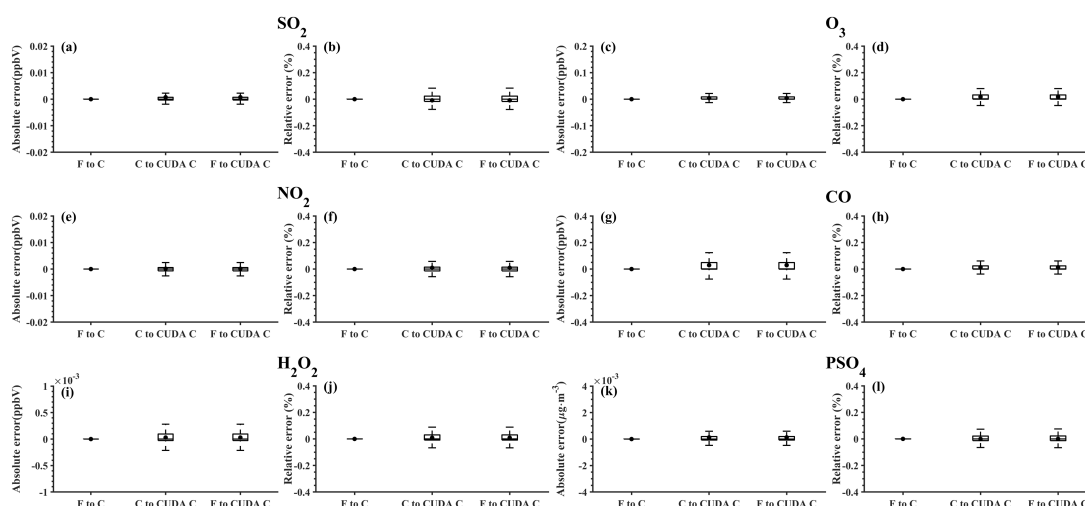402     concentration differences of Fortran and CUDA C versions.

# H$_2$O$_2$

**F Version**      **C Version**      **CUDA Version**



**F - C**      **C - CUDA**      **F - CUDA**



# PSO$_4$

**F Version**      **C Version**      **CUDA Version**



**F - C**      **C - CUDA**      **F - CUDA**



403

404 **Figure 7.** $H_2O_2$ and $PSO_4$ concentrations outputted by CAMx model for Fortran, standard C, and

405 CUDA C versions. Panels (a) and (g) are from Fortran versions. Panels (b) and (h) are from standard

406 C versions. Panels (c) and (i) are from CUDA C versions. Panels (d) and (j) are the output

407 concentration differences of Fortran and standard C versions. Panels (e) and (k) are the output

408 concentration differences of standard C and CUDA C versions. Panels (f) and (l) are the output

409 concentration differences of Fortran and CUDA C versions.

410      Figure 8. shows the boxplot of AEs and relative error (REs) in all grid boxes for

411 the six species during the porting process. As described above, the AEs and REs

412 introduced by the Fortran to standard C code refactoring process are significantly small,

413 and the primary error comes from converting standard C to CUDA C. Statistically, the

414 average of AEs (REs) of $SO_2$, $O_3$, $NO_2$, CO, $H_2O_2$ and $PSO_4$ were -0.0009 ppbV (-

415 0.01%), 0.0004 ppbV (-0.004%), 0.0005 ppbV (0.008%), 0.03 ppbV (0.01%),

416 $2.1 \times 10^{-5}$ ppbV (-0.01%) and 0.0002 $\mu g \cdot m^{-3}$ (0.0023%), respectively between

417 the Fortran and CUDA C versions. In terms of time series, the regionally averaged time

418 series of the three versions are almost consistent (as is shown in Figure S2), and the

419 maximum AEs for the above six species are 0.001ppbV, 0.005 ppbV, 0.002 ppbV,

420 0.03ppbV, 0.0001 ppbV and 0.0002 $\mu g \cdot m^{-3}$, respectively, between the Fortran and

421 CUDA C versions.

422



423 **Figure 8.** The distributions of absolute errors and relative errors for $SO_2$, $O_3$, $NO_2$, CO, $H_2O_2$ and

424 $PSO_4$ in all of the grid boxes after 48 hours of integration.

425      Figure 9. presents the regionally averaged time series and AEs of $SO_2$, $O_3$, $NO_2$,

426 CO, H$_2$O$_2$ and PSO$_4$. The time series between different versions is almost consistent,

427 and the maximum AEs for above six species are 0.001ppbV, 0.005 ppbV, 0.002 ppbV,

428 0.03ppbV, 0.0001 ppbV and 0.0002 $\mu g \cdot m^{-3}$, respectively between the Fortran and

429 CUDA C versions.

430 It is difficult to verify the scientific applicability of the results from CUDA C

431 version because the programming language and hardware are different between the

432 Fortran and CUDA C version. Here, we used the evaluation method of Wang et al.

433 (2021a) to compute the root mean square errors (RMSEs) of SO$_2$, O$_3$, NO$_2$, CO, H$_2$O$_2$

434 and PSO$_4$ between the Fortran and CUDA C versions, which are 0.0007 ppbV, 0.001

435 ppbV, 0.0002 ppbV, 0.0005 ppbV, 0.00003 ppbV, and 0.0004 $\mu g \cdot m^{-3}$ respectively,

436 much smaller than the spatial variation of the whole region, which is 7.0 ppbV

437 (approximately 0.004%), 9.7 ppbV (approximately 0.003%), 7.4 ppbV (approximately

438 0.003%), 142.2 ppbV (approximately 0.006%), 0.2ppbV (approximately 0.015%) and

439 1.7 $\mu g \cdot m^{-3}$ (approximately 0.004%). It is indicated that the bias between CUDA C

440 and Fortran version of the above six species is negligible compared with their own

441 spatial changes, and the results of the CUDA C version are generally acceptable for

442 research.

443

444 **4.3. Offline performance comparison of GPU-HADVPPM**

445 As described in the Sect. 4.2, we validate that the CAMx model result of the

446 CUDA C version can be generally acceptable for scientific research. We tested the

447 offline performance of the HADVPPM and GPU-HADVPPM scheme on 1 CPU core

448 and 1 GPU card, respectively. There are 7 variables input into the HADVPPM program,

449 which are nn, dt, dx, con, vel, area, and areav, and their specific meanings are shown in

450 Table S1.

451 Firstly, we use random_number function in Fortran to create random single-

452 precision floating-point numbers of different sizes for the above 7 variables, and then

453 transmit these random numbers to the hadvppm Fortran program and hadvppm_kernel

454 CUDA C program for computation, respectively. Finally, test the offline performance
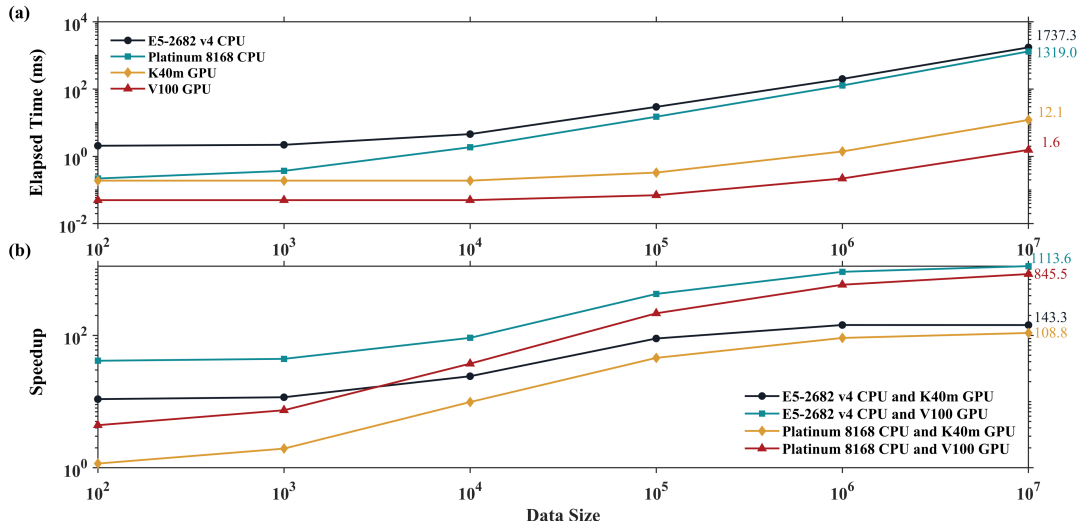
455  of the HADVPPM and GPU-HADVPPM on the CPU and GPU platforms. During the

456  offline performance experiments, we used two different CPUs and GPUs described in

457  the Sect. 2.4., and the experimental results are shown in Figure 9.

458      On the CPU platform, the wall time of hadvppm Fortran program does not change

459  significantly when the data size is less than 1000. With the increase in the data size, its

460  wall time increases linearly. When the data size reaches $10^7$, the wall time of the

461  hadvppm Fortran program on Intel Xeon E5-2682v4 and Intel Platinum 8168 CPU

462  platforms is 1737.3ms and 1319.0ms, respectively. On the GPU platform, the

463  reconstructed and extended CUDA C program implements parallel computation of

464  multiple grid points by executing a large number of kernel function copies, so the

465  computational efficiency of hadvppm_kernel CUDA C code on it is significantly

466  improved. In the size of $10^7$ random numbers, the hadvppm_kernel CUDA C program

467  takes only 12.1ms and 1.6ms to complete the computation on the NVIDIA Tesla K40m

468  and NVIDIA Tesla V100 GPU.

469      Figure 9. (b) shows the speedup of HADVPPM and GPU-HADVPPM on CPU

470  platform and GPU platform under different data sizes. When mapping the HADVPPM

471  scheme to GPU, the computational efficiency under different data size is not only

472  significantly improved, but also the larger the data size, the more obvious the

473  acceleration effect of the GPU-HADVPPM. For example, in the size of $10^7$ random

474  numbers, the GPU-HADVPPM achieved 1113.6x and 845.4x acceleration on the

475  NVIDIA Tesla V100 GPU, respectively, compared to the two CPU platforms. Although

476  the K40m GPU's single-card computing performance is slightly lower than that of the

477  V100 GPU, GPU-HADVPPM can also achieve up to 143.3x and 108.8x acceleration.

478      As described in Sect. 3.2, the thread is the most basic unit of GPU for parallel

479  computing. Each dimension of the three-dimensional block can contain a maximum

480  number of threads of 1024,1024, and 64, respectively. Each dimension of the three-

481  dimensional grid can contain a maximum number of blocks of $2^{31} - 1$, 65535, and

482  65535. It is theoretically possible to distribute a large number of copies of kernel

483  functions into tens of billions of threads for parallel computing without exceeding the

25

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



**Figure 9.** The offline performance of the HADVPPM and GPU-HADVPPM scheme on CPU and GPU. The unit of the wall times for the offline performance experiments is millisecond(ms).

## 4.4. Coupling performance comparison of GPU-HADVPPM with different GPU configurations

### 4.4.1. CAMx-CUDA on a single GPU

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

504    between CPU and GPU is reduced to 1 within one time-step, we further optimize the

505    GPU memory access order on GPU card, eliminate unnecessary assignment loops

506    before kernel functions launched and use thread and block indices.

507          Table 4. lists the total elapsed time for different versions of CAMx-CUDA model

508    during the optimization, as described in Section 3.2. Since the xyadvec program in the

509    CAMx-CUDA V1.0 is not optimized, it is extremely computationally inefficient when

510    starting two CPU processes and configuring a GPU card for P1. On the K40m and V100

511    cluster, it takes 10829 seconds and 37237 seconds respectively to complete 1-hour

512    simulation.

513          By optimizing the algorithm of xyadvec Fortran program and hadvppm_kernel

514    CUDA C program, the frequency of data transmission between CPU and GPU was

515    decreased, and the overall computing efficiency was improved after GPU-HADVPPM

516    coupling to CAMx-CUDA model. In CAMx-CUDA V1.2, the frequency of data

517    transmission between CPU-GPU within one time step is reduced to 1, and the elapsed

518    time on the two heterogeneous clusters is 1207 seconds and 548 seconds, respectively,

519    and the speedup is 9.0x and 68.0x compared to the CAMx-CUDA V1.0.

520          GPU memory access order can directly affect the overall computational

521    efficiency of GPU-HAVPPM on the GPU. In CAMx-CUDA V1.3, we have optimized

522    the memory access order of hadvppm_kernel CUDA C program on the GPU and

523    eliminated unnecessary assignment loops before kernel functions launched, which

524    further improved the CAMx-CUDA model computational efficiency, resulting in 12.7x
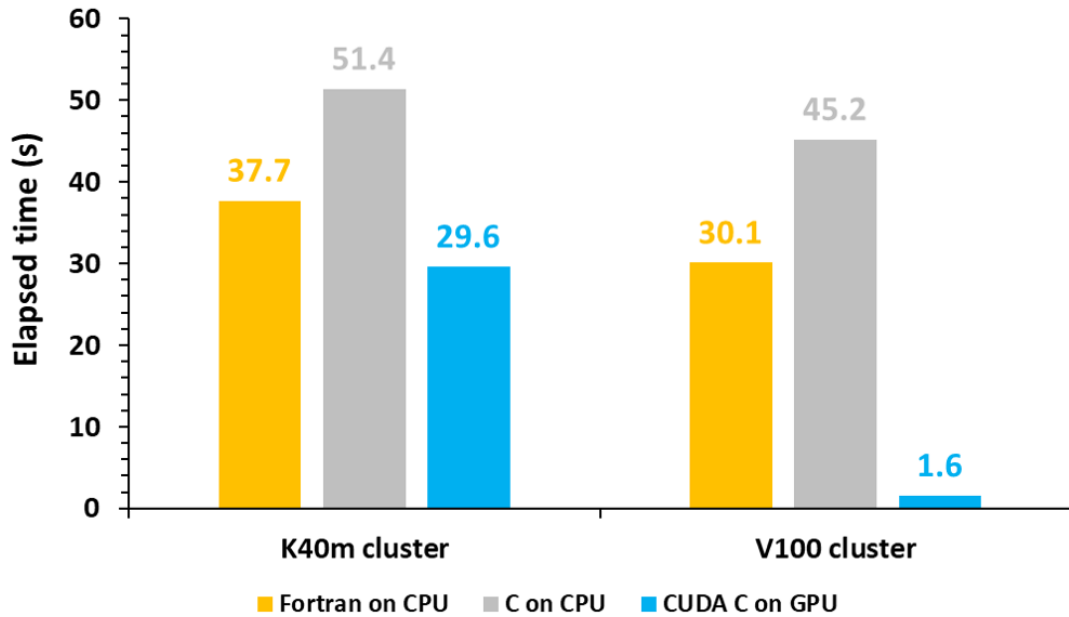
525    and 94.8x speedups.

526          Using thread and block indices to compute horizontal grid points simultaneous can

527    greatly improve the computational efficiency of GPU-HADVPPM and thus reduce the

528    overall elapsed time of CAMx-CUDA model. CAMx-CUDA V1.4 further reduces the

529    elapsed time by 378 seconds and 103 seconds respectively on K40m cluster and V100

530    cluster compared with CAMx-CUDA V1.3, and achieving up to 29.0x and 128.4x

531    speedup compared with CAMx-CUDA V1.0.

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

27

533　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 |

534　　　　In terms of the single module computational efficiency of HADVPPM and GPU-

535　HADVPPM, we further coupling test the computational performance of the Fortran

536　version HADVPPM on the CPU, C version HADVPPM on the CPU, and CUDA C

537　version GPU-HADVPPM in CAMx-CUDA V1.4 (GPU-HADVPPM V1.4) on the GPU,

538　using system_clock functions in the Fortran language and cudaEvent_t in CUDA

539　programming. The specific results are shown in Figure 10. On the K40m cluster, it takes

540　37.7 seconds and 51.4 seconds to launch the Intel Xeon E5-2682 v4 CPU to run Fortran

541　and C version HADVPPM, the C version is 26.7% slower than the Fortran version.

542　After the CUDA technology was used to convert the C code into CUDA C, the CUDA

543　C version took 29.6 seconds to launch an NVIDIA Telsa K40m GPU to run GPU-

544　HADVPPM V1.4, with 1.3x and 1.7x acceleration. On the V100 cluster, the Fortran,

545　the C, and the CUDA C version are computationally more efficient than those on the

546　K40m cluster. It takes 30.1 seconds and 45.2 seconds to launch Intel Xeon Platinum

547　8168 CPU to run Fortran and C version HADVPPM and 1.6 seconds to run the GPU-

548　HADVPPM V1.4 using an NVIDIA V100 GPU. The computational efficiency of the

549　CUDA C version is 18.8x and 28.3x higher than Fortran and C versions.
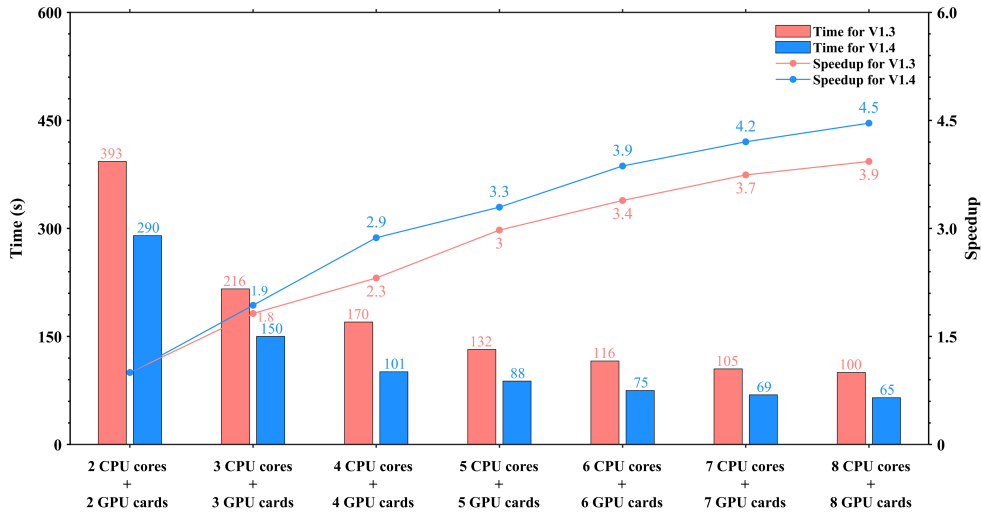
**Figure 10.** The elapsed time of the Fortran version HADVPPM on the CPU, the C version HADVPPM on the CPU, and CUDA C version GPU-HADVPPM V1.4 on the GPU. The unit is seconds (s).

**4.4.2. CAMx-CUDA on multiple GPUs**

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

29

569    multiple GPU cards. Figure 11. shows the total elapsed time and speedup of CAMx-

570    CUDA V1.3 and V1.4 on the V100 cluster. The total elapsed time decreases as the

571    number of CPU cores and GPU cards increases. When starting 8 CPU cores and 8 GPU

572    cards, the speedup of CAMx-CUDA V1.4 is increased from 3.9x to 4.5x compared with

573    V1.3, and the computational efficiency is increased by 35.0%.



574

575    **Figure 11.** The total elapsed time and speedup of CAMx-CUDA V1.3 and V1.4 on multiple

576    GPUs. The unit of elapsed time for experiments is seconds (s).

577    **5.    Conclusions and discussion**

578    GPU accelerators are playing an increasingly important role in high-performance

579    computing. In this study, a GPU acceleration version of the PPM solver (GPU-

580    HADVPPM) of horizontal advection for air quality model is developed, that can be run

581    on GPU accelerators using the standard C programming language and CUDA

582    technology. Offline performance experiments results show that K40m and V100 GPU

583    can achieve up to 845.4x and 1113.6x speedup, respectively, and the larger the data

584    input to the GPU, the more obvious the acceleration effect. After coupling GPU-

585    HADVPPM to CAMx model, a series of optimization measures are taken, including

586    reducing the CPU-GPU communication frequency, increasing the size of data

587 computation on GPU, optimizing the GPU memory access order, and using thread and

588 block indices to improve the overall computing performance of CAMx-CUDA model.

589 Using a single GPU card, the optimized CAMx-CUDA V1.4 model improves the

590 computing efficiency by 29.0x and 128.4x on the K40m cluster and the V100 cluster,

591 respectively. In terms of the single-module computational efficiency of GPU-

592 HADVPPM, it can achieve 1.3x and 18.8x speedup on NVIDA Tesla K40m GPU and

593 NVIDA Tesla V100 GPU respectively. To make full use of multi-core and multi-GPU

594 supercomputers and further improve the total computational performance of CAMx-

595 CUDA model, a parallel architecture with an MPI+CUDA hybrid paradigm is presented.

596 After implementing the acceleration algorithm, the total elapsed time decreases as the

597 number of CPU cores and GPU cards increases, and it can achieve up to 4.5x speedup

598 when launch 8 CPU cores and 8 GPU cards compared with 2 CPU cores and 2 GPU

599 cards.

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

601    1) We currently implemented thread and block co-indexing to compute horizontal

602 grid points in parallel. Given the CAMx model 3-dimensional grid computing

603 characteristics, 3-dimensional thread and block co-indexing will be considered to

604 compute 3-dimensional grid points in parallel.

605    2) The communication bandwidth of data transfer is one of the main issues for

606 restricting the computing performance of CUDA C codes on GPUs. This restriction not

607 only holds true for GPU-HADVPPM, but also WRF module as well (Mielikainen et al.,

608 2012b; Mielikainen et al., 2013b; Huang et al., 2013). In this study, data transmission

609 efficiency between CPU and GPU is improved only by reducing communication

610 frequency. In the future, more technologies, such as pinned memory (Wang et al.,2016),

611 will be considered to solve the communication bottleneck between CPU and GPU.

612    3) In order to further improve the overall computational efficiency of the CAMx

613 model, the heterogeneous porting scheme proposed in this study will be considered to

614 carry out the heterogeneous porting of other CAMx modules in the future.

615

**Reference**

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

Cao, K., Wu, Q., Wang, L., Wang, N., Cheng, H., Tang, X., Li, D., and Wang, L.: The dataset of the manuscript "GPU-HADVPPM V1.0: high-efficient parallel GPU design of the Piecewise Parabolic Method (PPM) for horizontal advection in air

quality          model          (CAMx          V6.10)",          ZENODO,
https://doi.org/10.5281/zenodo.7765218, 2023.

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

ENVIRON: User Guide for Comprehensive Air Quality Model with Extensions
Version 6.1, available at: https://camx-wp.azurewebsites.net/Files/CAMxUsers
Guide_v6.10.pdf (last access: 19 December 2022), 2014

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

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

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

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

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

Huang, M., Huang, B., Mielikainen, J., Huang, H. L. A., Goldberg, M. D., and Mehta,
A.: Further Improvement on GPU-Based Parallel Implementation of WRF 5-Layer
Thermal Diffusion Scheme, 2013 International Conference on Parallel and

Distributed Systems, 10.1109/icpads.2013.126, 2013.

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

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

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

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

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

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

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

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

701 Price, E., Mielikainen, J., Huang, M., Huang, B., Huang, H.-L. A., and Lee, T.: GPU-
702     Accelerated Longwave Radiation Scheme of the Rapid Radiative Transfer Model
703     for General Circulation Models (RRTMG), IEEE Journal of Selected Topics in
704     Applied Earth Observations and Remote Sensing, 7, 3660-3667,
705     10.1109/jstars.2014.2315771, 2014.

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

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

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

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

723 Wahib, M. and Maruyama, N.: Highly optimized full GPU-acceleration of non-
724     hydrostatic weather model SCALE-LES, 2013 IEEE International Conference on
725     Cluster Computing (CLUSTER), 23-27 Sept. 2013, 1-8,
726     10.1109/CLUSTER.2013.6702667,

727 Wang, P., Jiang, J., Lin, P., Ding, M., Wei, J., Zhang, F., Zhao, L., Li, Y., Yu, Z., Zheng,
728     W., Yu, Y., Chi, X., and Liu, H.: The GPU version of LASG/IAP Climate System
729     Ocean Model version 3 (LICOM3) under the heterogeneous-compute interface for

730      portability (HIP) framework and its large-scale application, Geosci. Model Dev.,

731      14, 2781-2799, 10.5194/gmd-14-2781-2021, 2021a.

732  Wang, Y., Guo, M., Zhao, Y., and Jiang, J.: GPUs-RRTMG_LW: high-efficient and

733      scalable computing for a longwave radiative transfer model on multiple GPUs,

734      The Journal of Supercomputing, 77, 4698-4717, 10.1007/s11227-020-03451-3,

735      2021b.

736  Wang, Z., Wang, Y., Wang, X., Li, F., Zhou, C., Hu, H., and Jiang, J.: GPU-

737      RRTMG_SW: Accelerating a Shortwave Radiative Transfer Scheme on GPU,

738      IEEE Access, 9, 84231-84240, 10.1109/access.2021.3087507, 2016.

739  Xiao, H., Lu, Y., Huang, J., and Xue, W.: An MPI+OpenACC-based PRM scalar

740      advection scheme in the GRAPES model over a cluster with multiple CPUs and

741      GPUs, Tsinghua Science and Technology, 27, 164-173,

742      10.26599/TST.2020.9010026, 2022.

743  Xu, S., Huang, X., Oey, L. Y., Xu, F., Fu, H., Zhang, Y., and Yang, G.: POM.gpu-v1.0:

744      a GPU-based Princeton Ocean Model, Geoscientific Model Development, 8,

745      2815-2827, 10.5194/gmd-8-2815-2015, 2015.

746  Zhang, Q., Streets, D. G., Carmichael, G. R., He, K. B., Huo, H., Kannari, A., Klimont,

747      Z., Park, I. S., Reddy, S., Fu, J. S., Chen, D., Duan, L., Lei, Y., Wang, L. T., and

748      Yao, Z. L.: Asian emissions in 2006 for the NASA INTEX-B mission, Atmos.

749      Chem. Phys., 9, 5131-5153, 10.5194/acp-9-5131-2009, 2009.