



1 **Enhancing the advection module performance in the EPICC-Model**
2 **V1.6.0 via GPU-HADVPPM4HIP V1.0 coupling and GPU-optimized**
3 **strategies**

4 **Kai Cao¹, Qizhong Wu², Xiao Tang^{1,3}, Jinxi Li¹, Xueshun Chen^{1,3}, Huansheng Chen¹,**
5 **Wending Wang¹, Huangjian Wu¹, Lei Kong¹, Jie Li^{1,3}, Jiang Zhu^{1,3}, and Zifa Wang^{1,3}**

6 ¹State Key Laboratory of Atmospheric Environment and Extreme Meteorology, Institute of
7 Atmospheric Physics, Chinese Academy of Sciences, Beijing 100029, China

8 ²College of Global Change and Earth System Science, Faculty of Geographical Science, Beijing
9 Normal University, Beijing 100875, China

10 ³College of Earth and Planetary Sciences, University of Chinese Academy of Sciences, Beijing,
11 100049, China

12

13 **Correspondence to:** Qizhong Wu (wqizhong@bnu.edu.cn); Xiao Tang (tangxiao@mail.iap.ac.cn)

14

15 **Abstract**

16 The rapid development of Graphics Processing Units (GPUs) has established new
17 computational paradigms for enhancing air quality modeling efficiency. In this study, the
18 heterogeneous-compute interface for portability (HIP) was implemented to parallel computing of
19 the piecewise parabolic method (PPM) advection solver (HADVPPM) on China's domestic GPU-
20 like accelerators (GPU-like), resulting in a GPU-accelerated version denoted as GPU-
21 HADVPPM4HIP V1.0. Computational performance was enhanced through three strategic
22 optimizations: reducing the central processing unit (CPU) and GPU (CPU-GPU) data transfer
23 frequency, thread-block coordinated indexing, and the Message Passing Interface (MPI) and HIP
24 ("MPI+HIP") hybrid parallelization across heterogeneous computing clusters. Following validation
25 of the GPU-HADVPPM4HIP V1.0 program's offline computational consistency and the pollutant
26 simulation performance of the Emission and atmospheric Processes Integrated and Coupled
27 Community version 1.6.0 (EPICC-Model V1.6.0) on the Earth System Numerical Simulation



28 Facility (EarthLab), comprehensive performance testing was conducted. Offline benchmark results
29 demonstrated that GPU-HADVPPM4HIP V1.0 achieved a maximum speedup of 556.5x on a GPU-
30 like using the compiler optimization option compared to the Fortran HADVPPM baseline compiled
31 option for a data size of 10^8 . Integrating GPU-HADVPPM4HIP V1.0 into EPIC-Model V1.6.0
32 yielded three distinct versions: the initial HIP-based version (HIP-Ori), a version optimized for CPU
33 and GPU communication frequency (HIP-Opt1), and a further-optimized version employing a
34 thread-block coordinated indexing strategy (HIP-Opt2). Compared to the HIP-Ori version,
35 HIP-Opt1 achieved a model-level computational efficiency improvement of 17.0x. Building upon
36 HIP-Opt1, HIP-Opt2 delivered an additional 1.5x enhancement in computational efficiency. At the
37 module level, including CPU and GPU data transfer overhead, the GPU implementation improves
38 computational efficiency of the advection module by 39.3%; when communication cost is excluded,
39 the advection module attains a $20.5\times$ acceleration relative to its CPU counterpart. This coupling
40 establishes a foundational framework for adapting air quality models to GPU-like architectures and
41 identifies critical optimization pathways. Moreover, the methodology provides essential technical
42 support for achieving full-model GPU implementation of the EPIC-Model, addressing both
43 current computational constraints and future demands for high-resolution air quality simulations.

44 **1. Introduction**

45 Air pollution, a source of fine particulate matter in both urban and rural areas, is associated
46 with an elevated risk of strokes, heart diseases, lung cancer, and acute and chronic respiratory
47 diseases (Atkinson et al., 2010; Kim et al., 2015; Liu et al., 2016; Milton and White, 2020). The air
48 quality forecasting system centered on the air quality model plays a critical role in the timely
49 dissemination of forecasting alerts and early warning information to the public. The accuracy of air
50 quality forecasting is jointly constrained by the spatial resolution of input datasets, including
51 emission inventories, terrain, and meteorological parameters (Gupta et al., 2015; Georgiou et al.,
52 2022). High-resolution model configurations have demonstrated improvements in the accuracy of
53 air quality forecasting (Georgiou et al., 2018; Podrascanin et al., 2019; Adani et al., 2022; Gao and
54 Zhou, 2024). However, current operational forecasting systems predominantly employ horizontal
55 resolutions ranging from several to tens of kilometers (Wu et al., 2014; Guevara et al., 2021; Tang



56 et al., 2022; Gao and Zhou, 2024), which inadequately address the requirements for urban-scale
57 high-resolution forecasting and precision management.

58 Computational demands emerge as a critical limiting factor for high-resolution air quality
59 modeling. On the one hand, doubling the horizontal resolution quadruples the number of
60 computational grids. On the other hand, maintaining numerical integration stability necessitates
61 proportional reduction in temporal integration steps (Georgiou et al., 2022). These combined effects
62 result in exponential growth of computational workload with increasing resolution. It is estimated
63 that when the horizontal resolution of the air quality model is increased by 18 times, the
64 computational load of the model increases by 300 times (Thompson and Selin 2012).

65 Enhancement of computational efficiency in air quality modeling has been predominantly
66 achieved through hardware-based acceleration strategies. Wang et al. (2017) ported the Global
67 Nested Air Quality Prediction Modeling System (GNAQPMS) to the second-generation Intel Xeon
68 Phi processor (KNL), achieving a 3.5x computational acceleration via MPI and OpenMP hybrid
69 parallelization, vectorization optimization, memory access pattern refinement, thread-local storage
70 reduction, and global communication optimization. The gas-phase chemistry module is widely
71 recognized as the dominant computational bottleneck in air quality models, typically accounting for
72 over 40% of total simulation time (Elbern, 1997; Linford et al., 2011; Wang et al., 2017; Cao et al.,
73 2023). To address this limitation, Wang et al. (2019) developed the MP CBM-Z mechanism by
74 implementing vectorized computation techniques within the CBM-Z framework. Leveraging Single
75 Instruction Multiple Data (SIMD) architecture, their approach enabled multi-point parallel
76 computation for gas-phase chemistry, achieving a 4.9x acceleration in the chemistry module and a
77 2.22x overall speedup for the entire NAQPMS model when deployed on Intel Xeon Gold 6132
78 CPUs.

79 In recent years, GPUs have emerged as transformative accelerators in artificial intelligence and
80 high-performance computing, driven by their massive parallel computing capabilities. In December
81 2024, the 64th TOP500 list of supercomputers revealed that the El Capitan system has achieved the
82 top spot, becoming the third exascale computing system following Frontier and Aurora, with an
83 HPL score of 1.742 EFLOP/s (Top500, 2024). This computational supremacy primarily originates
84 from its AMD Instinct MI300A GPU accelerators, each containing 14,592 stream processors and



85 delivering a double-precision floating-point performance of 61.3 TFLOP/s. Remarkably, the
86 computational efficiency of a single MI300A GPU exceeds 1.8 times the peak performance of the
87 Earth-Simulator supercomputer (CPU-based architecture) in Japan, which is Top1 supercomputer
88 in 2003.

89 The formidable computational capacity of GPUs has opened new directions for enhancing the
90 computational efficiency of air quality models. Alvanos and Christoudias (2017) developed a
91 software package for the global atmospheric chemistry model ECHAM/MESy Atmospheric
92 Chemistry (EMAC), enabling automated generates CUDA kernels to numerically integrate
93 atmospheric chemical kinetics by the Kinetic PreProcessor (KPP, Damian et al., 2002). Subsequent
94 memory optimization and thread management strategies achieved a 20.4x acceleration for the
95 chemistry module on NVIDIA P100 GPUs. In parallel efforts, Sun et al. (2018) implemented
96 CUDA-based optimization for the second-order Rosenbrock chemical solver (Sandu et al., 1997)
97 within the CAM4-Chem global chemistry-climate model. Through strategic enhancements in fully
98 interleaved memory layout, CUDA streams, and constant memory, they achieved an 11.7x speedup
99 for computation alone and a 3.8x speedup when the data transfer between the CPU and GPU is
100 considered on the NVIDIA Tesla K20X GPU. Notably, Quevedo et al. (2025) adapted the third-
101 order Rosenbrock solver in the CMAQ model by converting Fortran code to CUDA Fortran,
102 evaluating its performance across three chemical mechanisms: RACM2, CB6R5, and SAPRC07.
103 Comparative analysis revealed 51%, 50%, and 35% computational efficiency gains on NVIDIA
104 RTX 2080 Ti GPUs, respectively, while maintaining numerical consistency with CPU-based
105 benchmarks. Through code refactoring from Fortran to standard C and the HIP programming
106 technology, Cao et al. (2025) successfully parallelized the fourth-order Rosenbrock solver on GPU-
107 like architecture. Concurrently, the total model elapsed time was reduced by 46.9%. Regarding
108 another hotspot module in air quality models-the advection module, Cao et al. (2023, 2024)
109 implemented GPU-accelerated adaptations of the CAMx model's advection module using CUDA
110 and HIP heterogeneous technologies, respectively, and the optimized advection module achieved
111 maximum speedups of 80.2x on NVIDIA Tesla V100 GPUs and 28.9x on GPU-like.

112 Following Cao et al.'s (2025) successful implementation of parallel computing for the gas-
113 phase chemistry module in the EPIC-Model on GPU-like accelerators, the computational time



114 proportion of this module was significantly reduced. Consequently, the advection module has
115 emerged as a computational hotspot with comparable time consumption to the optimized chemistry
116 module. To address this shift, this study focuses on enhancing the advection module performance in
117 EPICC-Model V1.6.0 via GPU-HADVPPM4HIP V1.0 coupling and GPU-optimized strategies.
118 Sect. 2 details the EPICC-Model's computational framework, baseline performance tests, and the
119 heterogeneous computing platform employed in this research. Sect. 3 elaborates on the optimization
120 framework specifically designed for the computational characteristics of the EPICC-Model
121 advection module. Sect. 4 presents experimental results, including simulation performance and
122 computational performance analysis.

123 **2. The EPICC-Model and experiments**

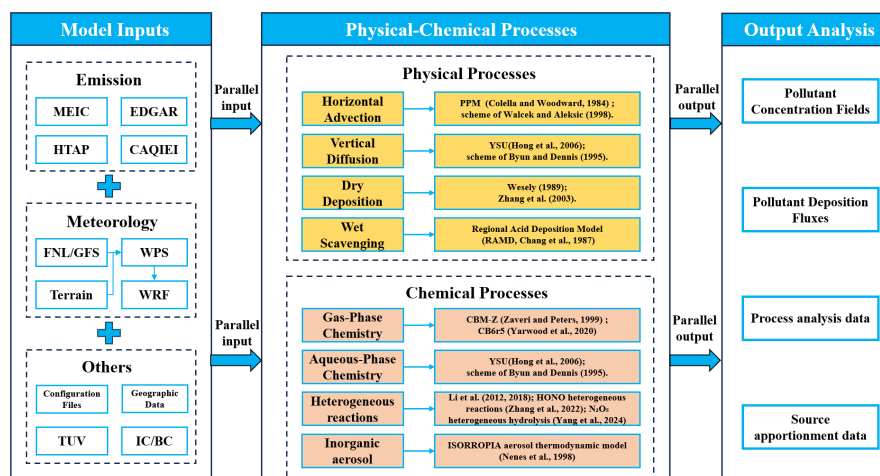
124 **2.1. The framework of the EPICC-Model**

125 The emission and atmospheric processes integrated and coupled community model version
126 V1.6.0 (EPICC-Model V1.6.0; EPICC-Model Working Group, 2025; Wang et al., 2025) is an air
127 quality modeling system specifically designed for air pollution complex in China (Zhu et al., 2023),
128 and developed by the Institute of Atmospheric Physics, Chinese Academy of Sciences based on the
129 Earth System Numerical Simulation Facility (EarthLab, Chai et al., 2021). The model framework is
130 fundamentally based on the species continuity equation and is used to simulate the complex physical
131 and chemical processes of pollutants in the atmosphere. These processes include emissions,
132 advection, diffusion, aerosol processes, gas-phase chemistry, and deposition. The EPICC-Model
133 V1.6.0 adopts a modular architecture developed using Fortran programming language, a high-
134 performance computing language specifically designed for scientific applications. The model code
135 is open-source and shared (EPICC-Model Working Group, 2026). This open-source code repository
136 enables rapid integration of novel mechanisms and modules proposed by diverse research groups,
137 thereby enhancing collaborative development efficiency.

138 The computational framework and workflow of the EPICC-Model V1.6.0 are illustrated in
139 Figure 1. The system primarily comprises three components: model inputs, physical-chemical
140 processes, and outputs analysis. Model input data include emissions, meteorological data, and other
141 datasets. Emissions inventories such as the Multi-resolution Emission Inventory for China (MEIC,



142 Li et al., 2017), the Emissions Database for Global Atmospheric Research (EDGAR, Crippa et al.,
143 2024), the HTAP (Crippa et al., 2023), and the Inversed Emission Inventory for Chinese Air Quality
144 (CAQIEI, Kong et al., 2024) can be utilized. Meteorological data are predominantly derived from
145 simulations generated by the mesoscale Weather Research and Forecasting (WRF) model. Other
146 datasets encompass configuration files, terrain data, TUV photolysis data, as well as initial
147 conditions (IC) and boundary conditions (BC). Physical-chemical processes primarily include
148 horizontal advection, vertical diffusion, dry deposition, wet scavenging, gas-phase chemistry,
149 aqueous-phase chemistry, heterogeneous reactions, inorganic aerosol thermodynamics, etc. For the
150 vertical diffusion module, either the scheme of Byun and Dennis (1995) or the YSU scheme (Hong
151 et al., 2006) can be selected to calculate the turbulent vertical diffusion coefficient. The dry
152 deposition module can employ either the scheme of Wesely (1989) or Zhang et al. (2003) to compute
153 deposition velocities. The gas-phase chemistry module offers the option to utilize either the CBM-
154 Z (Zaveri and Peters, 1999) or CB6r5 (Yarwood et al., 2020) chemical mechanisms. For
155 heterogeneous reactions, the model defaults to the scheme of Li et al. (2012). Additionally, it
156 integrates mechanisms for HONO heterogeneous chemical reactions (Zhang et al., 2022), sulfate
157 heterogeneous chemical reactions (Li et al., 2018), and N₂O₅ heterogeneous hydrolysis (Yang et al.,
158 2024). Inorganic aerosol is simulated using the ISORROPIA aerosol thermodynamic model (Nenes
159 et al., 1998). The aqueous-phase chemistry module originates from the Regional Acid Deposition
160 Model (RADM, Chang et al., 1987). Regarding model output analysis, the EPICC-Model can
161 generate pollutant concentration fields, pollutant deposition fluxes, process analysis data, and source
162 apportionment data. For a comprehensive technical description of the model architecture and
163 implementation details, refer to the EPICC-Model Working Group (2025).



164

165 **Figure 1.** The computational framework and workflow of the EPICC-Model V1.6.0. In the section
 166 of physical-chemical processes, yellow represents the physical module, and orange indicates the
 167 chemical module.

168 For the horizontal advection module that is the focus in this study, two high-precision
 169 numerical schemes are available, the positive-definite mass-conservative differencing scheme
 170 (Walcek and Aleksic, 1998) and the piecewise parabolic method (PPM, Colella and Woodward,
 171 1984). The PPM scheme, an extension of high-order Godunov’s method, operates by partitioning
 172 the integration domain into subregions and approximating solutions using parabolic functions.
 173 Renowned for its numerical precision and robustness in complex fluid dynamics, this classic
 174 algorithm has been widely adopted in atmospheric chemistry models including the latest CMAQ
 175 and CAMx (Appel, et al., 2021; Emery, et al., 2024). Within the EPICC-Model V1.6.0 framework,
 176 the advection module sequentially executes transport processes in the *x*-direction and *y*-direction. It
 177 employs species-specific PPM solvers (HADVPPM subroutine) for gaseous species, inorganic
 178 aerosols, organic aerosols, dust, and sea salt. Our previous studies have demonstrated a significant
 179 acceleration performance of PPM solver in the CAMx model through HIP heterogeneous
 180 programming technologies for GPU-like (Cao et al., 2024).

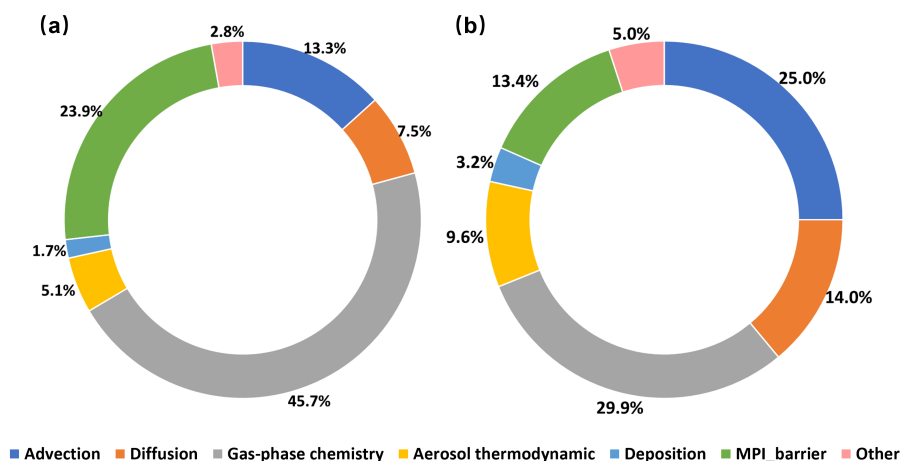
181 2.2. Benchmark performance of testing

182 As mentioned above, Cao et al. (2025) implemented the fourth-order Rosenbrock solver for
 183 the gas-phase chemistry module in the EPICC-Model, employing the CBM-Z chemical mechanism.



184 Through code restructure from Fortran to standard C programming language and implementation of
 185 the HIP heterogeneous programming framework, the computational efficiency of the gas-phase
 186 chemistry module improved by 2.88 times when accounting for data transfer between CPUs and the
 187 GPU-like accelerators.

188 Figure 2. illustrates the changes in the computational time proportion among modules before
 189 and after the heterogeneous parallel implementation of the gas-phase chemistry module on GPU-
 190 like accelerators, as reported by Cao et al. (2025). Specifically, Figure 2(a) shows the time
 191 proportion of each module in the original Fortran version, while Figure 2(b) presents the
 192 corresponding proportion after the gas-phase chemistry module was ported for parallel computing
 193 on GPU-like. As shown in Figure 2, the implementation of parallel computing for gas-phase
 194 chemistry modules on GPU-like achieved significant efficiency improvements, reducing its
 195 computational time proportion from 45.7% to 29.9%. Notably, the computational time proportion
 196 of MPI_Barrier synchronization function has decreased from 23.9% to 13.4%. A critical observation
 197 emerged regarding the advection module, whose computational time proportion increased from 13.3%
 198 to 25.0%, establishing it as a new performance bottleneck comparable to the optimized chemistry
 199 module. This performance shift necessitates subsequent optimization efforts focusing on
 200 heterogeneous porting and parallel acceleration of the PPM scheme for GPU-like architectures
 201 within the EPICC-Model framework, aiming to enhance the computational efficiency of the
 202 advection module.



203 ■ Advection ■ Diffusion ■ Gas-phase chemistry ■ Aerosol thermodynamic ■ Deposition ■ MPI_barrier ■ Other
 204 **Figure 2.** The computational time proportion among modules (a) for the Fortran version and (b)



205 after implementing the gas-phase chemistry module on GPU-like accelerators for parallel
206 computation.

207 **2.3. Hardware platform and software environment of experiments**

208 All performance testing of the EPICC-Model V1.6.0 and heterogeneous adaptation and
209 optimization studies of the advection module on GPU-like accelerators were conducted at the Earth
210 System Numerical Simulation Facility (EarthLab, Chai et al., 2021). Jointly developed by the
211 Institute of Atmospheric Physics, Chinese Academy of Sciences and collaborating institutions, this
212 platform, specifically designed for earth system modeling and high-resolution regional
213 environmental simulation, employs a CPU and GPU heterogeneous architecture. Detailed hardware
214 components and software environment are presented in Table 1. The Chinese domestic CPUs and
215 GPU-like accelerators used in this studying are the first-generation versions. Each GPU-like node
216 contains two China's domestic CPU processors and two GPU-like accelerators (Cao et al., 2024)
217 interconnected via PCIe 4.0 buses. The software stack employs Intel OneAPI 2021.3.0 toolkit for
218 CPU code compilation and dtk-23.04.1 toolkit for GPU-like code compilation, ensuring full
219 compatibility with heterogeneous computing paradigms.

220 **Table 1** The hardware components and software environment for the dedicated accelerator node on
221 the EarthLab.

	CPU	GPU
Hardware components	two of China's domestic CPU processors, 2.0 GHz, 32 cores	two of GPU-like accelerators, 3840 computing units, 16 GB memory
Software environment	Intel OneAPI 2021.3.0 toolkit	dtk-23.04.1

222 Compared to CPU processors, GPU-like accelerators demonstrate superior capability in
223 launching massive thread-level parallelism. Similar to the NVIDIA GPU architectures (NVIDIA,
224 2020), these GPU-like accelerators employ a three-level parallelism hierarchy comprising grids,
225 blocks, and threads, which collaboratively execute parallel computations through coordinated
226 indexing. Specifically, a computational grid is partitioned into multiple thread blocks with three-
227 dimensional coordinates, each thread block containing an array of three-dimensional indexed
228 threads. As the fundamental execution unit, individual threads perform concrete computational tasks,
229 each possessing a unique index ID that precisely determines its spatial position within the thread



230 block hierarchy. Consequently, the design of hierarchical indexing schemes coordinating blocks and
231 threads constitutes a critical challenge in achieving efficient parallel computation for three-
232 dimensional numerical modeling grids.

233 Analogous to the AMD's ROCm software stack (AMD, 2023), the dtk-23.04.1 toolkit (Cao et
234 al., 2024) includes programming models, tools, compilers, libraries, and runtimes for artificial
235 intelligence and high-performance computing applications on GPU-like accelerators. Mirroring
236 ROCm's design paradigm, dtk-23.04.1 adopts the HIP programming language as its application
237 programming interface (API). This implementation leverages the Single-Instruction Multiple-
238 Thread (SIMT) execution model to effectively manage and coordinate massive thread parallelism
239 on GPU-like accelerators.

240 **3. Implementation details**

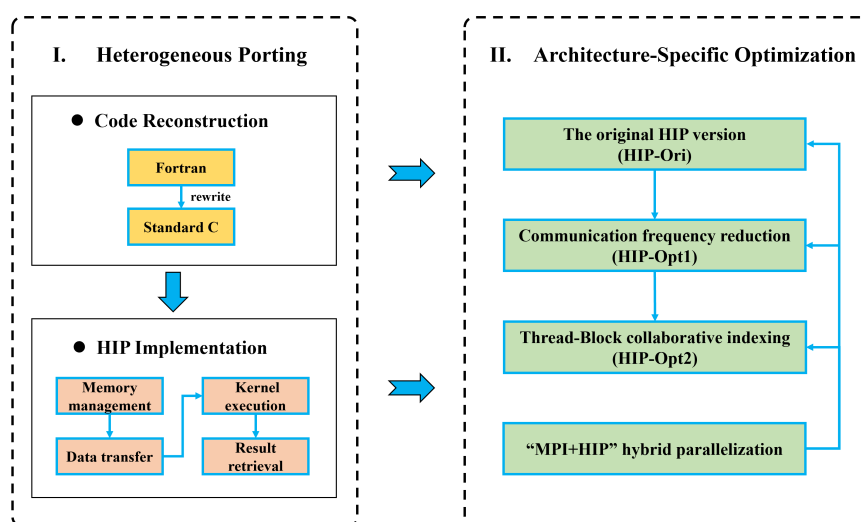
241 **3.1. Description of the heterogeneous porting and optimization scheme**

242 The heterogeneous porting and parallel optimization schemes of this study are illustrated in
243 Figure 3. Similar to the heterogeneous porting approach for the advection module of the air quality
244 model CAMx on GPU-like accelerators (Cao et al. 2024), the first step involved the porting and
245 adaptation of the HADVPPM advection solver from the EPIC-Model V1.6.0 to GPU-like
246 accelerators. Firstly, the Fortran code of the HADVPPM subroutine was reconstructed using
247 standard C programming language, followed by implementing parallel computing on GPU-like
248 accelerators through the HIP API. Similar to CUDA program execution on NVIDIA GPUs, the
249 implementation of GPU-HADVPPM4HIP V1.0 on GPU-like accelerators follows four key steps:
250 (1) Device memory allocation via the hipMalloc interface, (2) Data transfer from CPU to GPU-like
251 accelerator through hipMemcpy operations, (3) Parallel computation using kernel launching
252 (hipLaunchKernelGGL) with thread-index-based parallel processing after successful data
253 transmission, and (4) Final data retrieval from GPU back to CPU through hipMemcpy operations.

254 Following the implementation of GPU-HADVPPM4HIP V1.0 parallel computing on GPU-
255 like accelerators, the second phase involves architecture-specific parallel optimizations tailored for
256 GPU-like characteristics. Three optimization strategies were sequentially implemented to fully
257 exploit the SIMT vectorization parallelism of GPU-like accelerators, thereby enhancing the



258 computational performance of the EPICC-Model advection module. These strategies include: (1)
 259 reducing the frequency of communication between the CPU and GPU, (2) collaborative indexing
 260 between threads and blocks, and (3) hybrid parallelization of “MPI+HIP”. For systematic reference,
 261 three progressively optimized configurations were designated, namely HIP-Ori, HIP-Opt1, and
 262 HIP-Opt2. The HIP-Ori is baseline implementation after GPU-HADVPPM4HIP V1.0 integration
 263 into EPICC-Model without optimizations. The HIP-Opt1 is the version implementing reduced CPU-
 264 GPU communication frequency. The HIP-Opt2 is the enhanced version incorporating collaborative
 265 thread-block indexing. The hybrid “MPI+HIP” parallelization strategy was implemented across all
 266 three heterogeneous versions to enhance parallel scalability of the EPICC-Model V1.6.0 on the
 267 EarthLab.



268

269 **Figure 3.** Heterogeneous porting and parallel optimization scheme of advection module in the
 270 EPICC-Model V1.6.0.

271

272 3.2. HIP-Opt1: Communication frequency reduction

273 Influenced by the evolutionary trajectory of high-performance computing, most geoscientific
 274 numerical models, including the EPICC-Model, are predominantly coded in Fortran and designed
 275 for general-purpose CPU architectures. These models typically execute computations through grid-
 276 wise loop iterations. Taking the HADVPPM subroutine in the EPICC-Model as an example, its



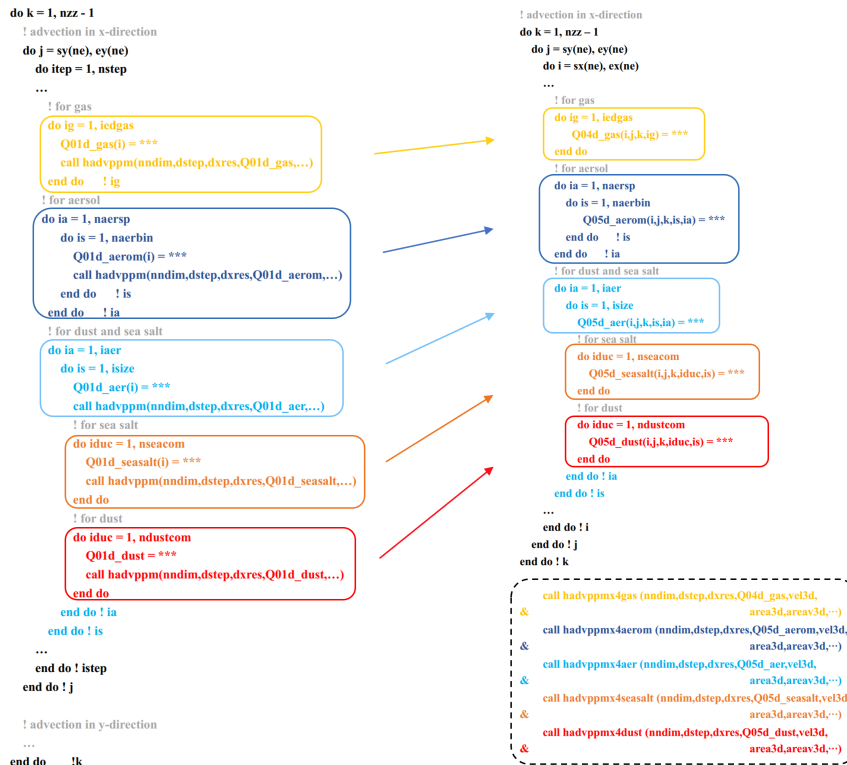
277 computational kernel is structured with triple nested loops progressing from innermost to outermost:
278 species loop (loop_species), latitudinal grid loop (loop_j), and vertical grid loop (loop_k). The
279 EPICC-Model innovatively categorizes atmospheric pollutants within the species loop into five
280 distinct classes: gaseous species (74 chemical constituents), inorganic aerosols (14 species
281 subdivided into coarse- and fine-mode particle sizes), organic aerosols, dust aerosols (8 species
282 across 4 size bins), and sea salt aerosols (11 species across 4 size bins).

283 The EPICC-Model calculates the advection process through looping of chemical-specified
284 variables, which has relatively low computational efficiency. A benchmark test using a two nested
285 grid configuration (horizontal grids: d01=228×165, d02=465×300) indicated that each timestep
286 requires approximately 9.8 million calls of the HADVPPM subroutine for advection processes in
287 both *x*-direction and *y*-direction. Consequently, the HIP-Ori version, generated by integrating GPU-
288 HADVPPM4HIP V1.0 into the EPICC-Model, incurs 9.8 million CPU-GPU data transfers per
289 timestep. Benchmark tests demonstrated that the computational time for 1-hour integration
290 increased from 1,015.0 seconds in the original Fortran version to 20,400.3 seconds under HIP-Ori
291 version, with frequent CPU-GPU communication identified as a primary performance bottleneck.
292 To address this critical bottleneck, our optimization framework prioritizes architectural redesign of
293 the advection module's loop hierarchy, strategically reducing communication frequency while
294 increasing data transfer sizes to better exploit GPU computational capacity.

295 Figure 4 illustrates the code-level implementation of CPU-GPU communication optimization,
296 with panel (a) depicting the HIP-Ori baseline and panel (b) presenting the optimized HIP-Opt1
297 version. In the two nested-domain case, the HIP-Ori configuration required approximately 4.9
298 million GPU calls for *x*-direction advection alone. To mitigate this computational overhead, we
299 restructured the advection module's loop architecture and expanded array dimensionality. The HIP-
300 Opt1 optimization framework implements multidimensional array restructuring, beginning with the
301 dimensional expansion of concentration variables from their original 1D representations (Q01d) in
302 HIP-Ori to 4D/5D configurations (Q04d/Q05d), while auxiliary parameters such as grid area
303 adjustment vector (area) and interfacial area adjustment vector (areav) are similarly upgraded from
304 1D to 3D structures. Prior to GPU execution, these variables undergo systematic multidimensional
305 reorganization, as demonstrated by the transformation of gaseous concentration variables from



306 Q01d_gas(i) to Q04d_gas(i, j, k, species) in Figures 4(a)-(b). This architectural redesign enables
 307 complete x-direction advection computation through a single GPU call per pollutant category.
 308 Consequently, the total GPU calls for both x and y direction advection decrease from approximately
 309 9.8 million in HIP-Ori to 10 in HIP-Opt1, achieved through one GPU call per spatial dimension
 310 across five pollutant categories, thereby optimizing computational efficiency through batched
 311 multidimensional data processing.



312 (a) Baseline code of the HIP-ori. (b) Optimized code of the HIP-Opt1.

313 **Figure 4.** The code-level implementation of CPU-GPU communication optimization. Panel (a) is
 314 the baseline Fortran code of the HIP-ori Panel (b) is the optimized Fortran code of the HIP-Opt1.

315 3.3. HIP-Opt2: Thread and block coordinated indexing

316 The architectural advantage of GPU-like accelerators manifests in their capacity to support
 317 massive thread concurrency for parallel computing. To leverage this capability, the coordinated
 318 thread-block indexing methodology which proposed by Cao et al. (2023) was implemented,
 319 whereby in which each grid cell in the two-dimensional horizontal plane is assigned a dedicated



320 thread. Specifically, blocks were configured based on the meridional grid dimension, with each
321 block allocated threads corresponding to the zonal grid count. This hierarchical parallelization
322 strategy achieves comprehensive full parallel processing of across the two-dimensional planar grid
323 structure through coordinated thread-block resource allocation. Furthermore, given that GPUs are
324 well-suited for large-scale matrix parallel computations without data dependencies, prior to
325 implementing parallel computation on a two-dimensional grid using thread and block coordinated
326 indexing, we decoupled the iterative computations present in the advection module by introducing
327 intermediate variables. This ensures computational independence at each step, meaning that the
328 computation of the next step does not depend on the results of the previous step, thereby fully
329 leveraging the multi-thread parallel computing capability of the GPU.

330 **3.4. “MPI+HIP” hybrid parallelization**

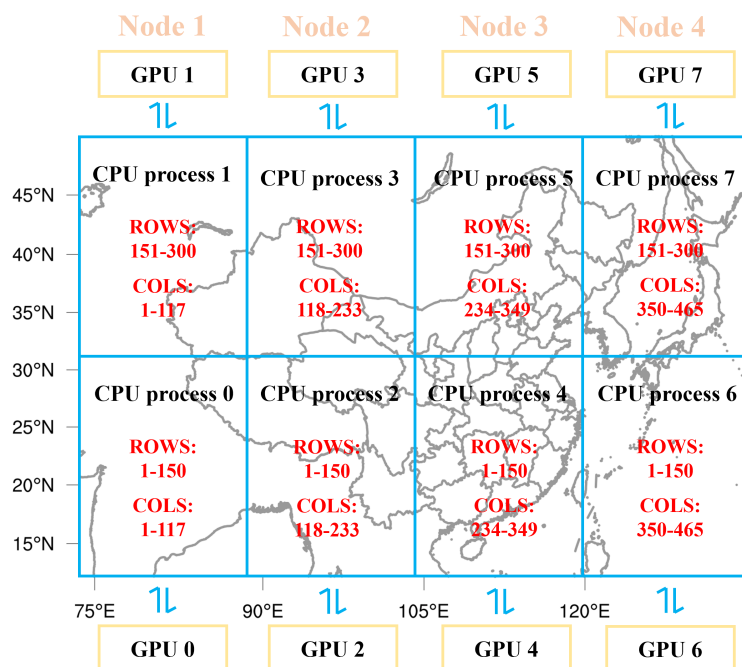
331 The recent advancements in GPU technology have solidified the dominance of "CPU+GPU"
332 heterogeneous architectures in global high-performance computing, with 9 out of the top 10
333 supercomputers in 2024 utilizing heterogeneous configurations (Top 100, 2024). These large-scale
334 heterogeneous clusters typically deploy one or multiple GPU accelerators per compute node.
335 Aligned with this trend, EarthLab employs heterogeneous architecture in its dedicated computing
336 nodes, integrating China's domestic CPUs with GPU-like accelerators. To maximize GPU
337 utilization and enhance the parallel scalability of the EPICC-Model on heterogeneous clusters, an
338 “MPI+HIP” hybrid parallelization scheme was designed tailored for EarthLab, inspired by the
339 “MPI+CUDA” approach proposed by Cao et al. (2023) for the CAMx model.

340 In the original Fortran version of EPICC-Model, the entire simulation domain is decomposed
341 into subdomains using MPI, with each CPU process responsible for computations, including those
342 of the advection module, within one subdomain. In this study, after MPI domain decomposition, we
343 offloaded the advection module computations originally performed on the CPU to China's domestic
344 GPUs through heterogeneous porting, thereby enabling parallel execution of the advection module
345 on GPUs. Considering future high-resolution applications with large-scale data, assigning a single
346 GPU to multiple CPU processes could lead to data-transfer contention and potential GPU memory
347 overflow. Therefore, we adopted an “MPI+HIP” hybrid parallelization approach, in which each
348 participating CPU process is assigned a dedicated GPU accelerator. This design expands the parallel



349 computing capacity of the model on heterogeneous clusters and makes full use of GPU resources.
 350 The implementation can be summarized in three main steps: (1) Obtain the MPI process rank
 351 information, as well as the number and indices of GPUs available on each compute node; (2) Based
 352 on the MPI process rank and using the remainder function in standard C, determine the number and
 353 indices of GPUs to be launched; (3) Map each MPI process to a specific GPU index, thereby
 354 realizing the “one CPU process – one GPU” configuration in the MPI+HIP hybrid parallel scheme.

355 As illustrated in Figure 5, taking the d02 domain configured in Section 4.1 with 8 CPU
 356 processes and 8 GPU-like accelerators as an example: the EPICCC model decomposes the simulation
 357 domain into 8 subdomains using the MPI software standard, with each CPU process handling its
 358 assigned subdomain. During the execution of the advection module, the “MPI+HIP” hybrid parallel
 359 scheme allocates one GPU accelerator to each CPU process. Computational tasks originally
 360 performed by the CPU are offloaded to the corresponding GPU; after the advection computations
 361 are completed, the results are returned to the CPU.



362
 363 **Figure 5.** An example schematic of domain subdivision and mapping to CPU processes, where each
 364 CPU process is equipped with a GPU-like accelerator.



365 4. Experimental results

366 4.1. Experimental setup

367 The centre of the simulation domain is located at (35 °N, 105 °E) and its two true latitude lines
368 are 25 °N and 45 °N, respectively. The EPIC-Model V1.6.0 employs a two-nested configuration,
369 the first domain (d01) covers East Asia with a 45 km × 45 km horizontal resolution on 228×165
370 grid cells. The lower left corner of the second domain has its starting positions in the grid of the first
371 domain as 36 and 39, respectively, and the second domain focuses on China with a 15 km × 15 km
372 horizontal resolution on 465×300 grid cells. In vertical, 20 terrain-following layers are configured
373 with the height of the top layer set to 20 km and six layers below 1 km. The numerical simulation
374 was conducted from 00:00 UTC July 1 to 23:00 UTC on 31 July, 2021, spanning a total duration of
375 744 hours. The initial 168-hour period was allocated for model spin-up. The MEIC emission
376 inventory (Li et al., 2017) was adopted as the emission input, and baseline year is 2019. The
377 numerical schemes selected during the EPIC-Model V1.6.0 simulations are listed in Table 2.
378 Furthermore, the PM_{2.5} and O₃ observations are from the China National Environmental Monitoring
379 Centre, which provides hourly PM_{2.5} and O₃ observations for eight cities in China. The station
380 information, including station name and its latitude and longitude, is listed in Table 3.
381 Meteorological inputs are generated using the Weather Research and Forecasting (WRF, Skamarock
382 et al., 2008) model, a mesoscale numerical weather prediction system, and is widely adopted in both
383 theoretical research and operational forecasting. This study employed the WRF version 3.9.1, and
384 the model domain configurations maintains identical nesting architecture and spatial coverage as
385 the EPIC-Model.

386 **Table 2.** The physical and chemical numerical schemes selected during EPIC-Model V1.6.0
387 simulation.

Process	Numerical schemes
Horizontal advection	PPM (Colella and Woodward, 1984)
Vertical diffusion	YSU scheme (Hong et al., 2006)
Gas-phase chemistry	CBM-Z (Zaveri and Peters, 1999)
Aqueous-phase chemistry and wet deposition	RADM (Chang et al., 1987)
Dry deposition	Scheme of Wesely (1989)
Inorganic aerosol thermodynamic partitioning	ISORROPIA v1.7 (Nenes et al., 1998)



388

389 **Table 3.** The names and latitude-longitude information of the PM_{2.5} and O₃ observation stations.

Name	Latitude (°N)	Longitude (°E)
Beijing	40.2865	116.1700
Taiyuan	37.7394	112.5583
Hangzhou	30.3058	120.3480
Hefei	31.7386	117.2780
Fuzhou	25.9664	119.5189
Qingdao	36.1032	120.3664
Guangzhou	23.5538	113.5890
Kunming	24.9786	102.7997

390

391 **4.2. Simulation performance analysis**

392 **4.2.1. Error analysis of the GPU-HADVPPM4HIP**

393 As elaborated in Sect. 3.1, the implementation of the HADVPPM Fortran code on GPU-like
394 comprises two principal phases. Initially, the Fortran code undergoes reconstruction using standard
395 C programming language through a Fortran-to-C conversion process. Subsequently, the standard C-
396 version HADVPPM program is adapted to GPU-like accelerators through C-to-HIP expansion
397 employing the HIP heterogeneous programming technology. Following successful GPU adaptation,
398 the offline precision verification was conducted to compare computational accuracy among three
399 implementations, the original Fortran source code, restructured standard C code, and HIP-
400 accelerated code. To ensure input consistency, a dedicated Fortran program was developed to
401 generate identical input datasets, including 100 double-precision floating-point numbers, for all
402 three implementations. Each implementation executed a complete advection integration
403 computation, with subsequent output recording and analysis. Therefore, the absolute errors (AE)
404 and relative errors (RE) presented in Table 4 represent the average values computed from the 100
405 double-precision floating-point results produced by the original Fortran code, restructured standard
406 C code, and HIP-accelerated code of the HADVPPM program after performing the advection
407 solution computation under the given input conditions.

408 Notably, both Fortran and C implementations were executed on China's domestic CPUs and
409 compiled using the Intel OneAPI 2021.3.0 toolkit, while GPU-HADVPPM4HIP was compiled with



410 the dtk-23.04.1 toolkit for GPU execution. For enhanced analytical rigor, we further implemented
411 compilation with -O0 and -O3 optimization flags. The -O0 flag maintains default compilation
412 settings without code optimization, whereas -O3 flag enables more aggressive loop and memory-
413 access optimizations, such as scalar replacement, loop unrolling (Intel Software, 2018).

414 Table 4 presents the mean absolute error (AE) and relative errors (RE) in computational
415 precision between the Fortran and standard C implementations (F-to-C), standard C and HIP
416 implementations (C-to-HIP), and Fortran and HIP implementations (F-to-HIP) of the HADVPPM
417 program under two compilation configurations. Notably, the -O3 optimization flag prioritizes
418 computational performance through code optimization at the expense of precision degradation.
419 Consequently, the AE and RE values for all three porting processes (F-to-C, C-to-HIP, and F-to-HIP)
420 under the -O0 configuration flag are systematically lower than those under -O3 flag. For instance,
421 between the Fortran and HIP implementations, the AE and RE under -O0 flag are measured at
422 1.4×10^{-7} and $4.3 \times 10^{-7}\%$, respectively. However, when compiled with -O3 flag, these errors
423 increase significantly to 3.1×10^{-7} and $2.5 \times 10^{-6}\%$, respectively. Similar error escalation
424 patterns are observed in the F-to-C and C-to-HIP comparisons under -O3 optimization flag. This
425 phenomenon aligns with expectations, as aggressive compiler optimizations (e.g., loop unrolling
426 and memory access reorganization) may introduce numerical instability through altered operation
427 sequences and reduced intermediate precision preservation.

428 Furthermore, it is noteworthy that across both -O0 and -O3 compilation configurations, the AE
429 and RE values of the F-to-C process consistently exceed those of the C-to-HIP process. This
430 indicates that the computational errors introduced during the heterogeneous porting of the
431 HADVPPM Fortran code from domestic CPUs to GPU-like accelerators predominantly originate
432 from the Fortran-to-C transcoding phase. For instance, under the -O0 configuration, the AE and RE
433 for F-to-C are measured at 1.5×10^{-7} and $5.1 \times 10^{-7}\%$, respectively, whereas those for C-to-HIP
434 are significantly smaller, at -9.5×10^{-9} and $-8.0 \times 10^{-8}\%$, respectively. This discrepancy arises
435 from inherent differences between Fortran and C in programming paradigms and data precision
436 management. The Fortran-to-C code restructuring involves two fundamentally distinct
437 programming languages, differing in object-oriented design philosophies and numerical
438 representation conventions, thereby introducing computational inaccuracies. In contrast, the HIP



439 programming model, analogous to NVIDIA CUDA, is inherently an extension of standard C. As
 440 detailed in Sect. 3.1, HIP achieves GPU compatibility by augmenting standard C programming
 441 language with critical GPU-specific functionalities, such as memory allocation and data transfer
 442 operations. Since HIP code essentially constitutes an enhanced standard C framework, the C-to-HIP
 443 adaptation introduces minimal computational bias, resulting in markedly lower error compared to
 444 the F-to-C transformation.

445 **Table 4.** Comparative of mean AE and RE across compilation configurations for F-to-C, C-to-HIP,
 446 and F-to-HIP Processes.

	AE			RE (%)		
	F-to-C	C-to-HIP	F-to-HIP	F-to-C	C-to-HIP	F-to-HIP
-O0	1.5×10^{-7}	-9.5×10^{-9}	1.4×10^{-7}	5.1×10^{-7}	-8.0×10^{-8}	4.3×10^{-7}
-O3	5.4×10^{-7}	-2.4×10^{-7}	3.1×10^{-7}	2.8×10^{-6}	-2.9×10^{-7}	2.5×10^{-6}

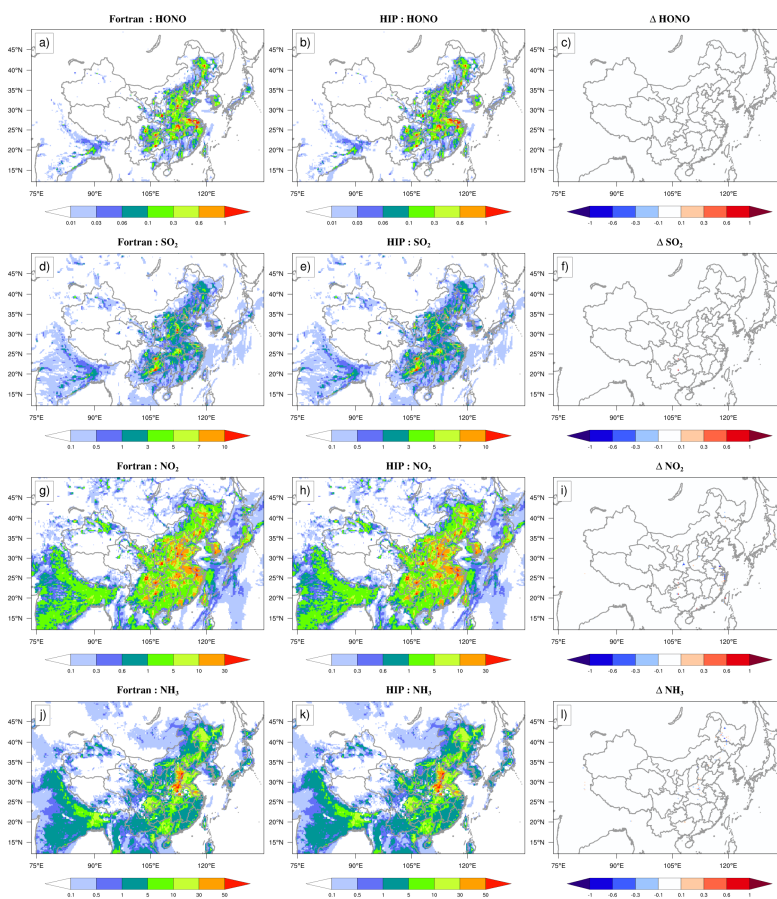
447 By coupling GPU-HADVPPM4HIP V1.0 into the EPIC-Model and implementing
 448 optimizations—such as reducing the communication frequency between CPUs and domestic GPUs,
 449 adopting thread and block coordinated indexing, and employing "MPI+HIP" hybrid
 450 parallelization—we developed the HIP-Opt2 version. Using identical input data and model
 451 configurations, the differences in simulation results between the HIP-Opt2 version and the original
 452 Fortran version were compared in a real-world case study. The model input data and configurations
 453 were set as described in Section 4.1, with the simulation period covering 24 hours from 00:00 to
 454 23:00 UTC on July 1, 2021.

455 Figures 6 and 7 present the simulated concentrations of gases (e.g., HONO, SO₂, NO₂, NH₃)
 456 and aerosols (e.g., BC, PM_{2.5}, ASO₄, ANH₄) after a 24-hour integration by both the HIP-Opt2 and
 457 the original Fortran versions, along with the Absolute Errors (AEs) between the two model versions.
 458 As visually evident from the figures, the results from HIP-Opt2 after the 24-hour integration are in
 459 close agreement with those from the original Fortran version. For the vast majority of grid points,
 460 the AEs for gases and aerosols remain within ± 0.1 ppbV or $\pm 0.1 \mu\text{g}\cdot\text{m}^{-3}$.

461 To further evaluate the suitability of HIP-Opt2 for scientific research, we followed the
 462 methodologies of Wang et al. (2021) and Cao et al. (2024) by introducing two metrics: the root mean
 463 square error (RMSE) and the standard deviation (std). The ratio of RMSE to std was calculated to
 464 quantify the scientific applicability of HIP-Opt2. Here, RMSE represents the error between HIP-



465 Opt2 and the original Fortran version for different species, while std denotes the standard deviation
 466 of each species in the Fortran version. Taking NO_2 as an example, a very small RMSE/std ratio
 467 indicates that the computational deviation introduced by heterogeneous parallel acceleration is
 468 negligible compared to the inherent spatial variability of NO_2 . This implies that such minor
 469 computational errors do not compromise the model's utility in scientific research. Table 5 lists the
 470 RMSE, std, and their ratios for the aforementioned gases and aerosols. The RMSE/std ratios for
 471 these species range from $10^{-5}\%$ to $10^{-2}\%$. Specifically, BC exhibits the smallest ratio at $4.8 \times 10^{-5}\%$,
 472 while ASO_4 shows the largest ratio, yet only at $7.0 \times 10^{-2}\%$. The RMSE/std ratios for all species in
 473 HIP-Opt2 are comparable to those reported by Cao et al. (2024), demonstrating that the simulation
 474 results from the HIP-Opt2 version are fully suitable for scientific applications.

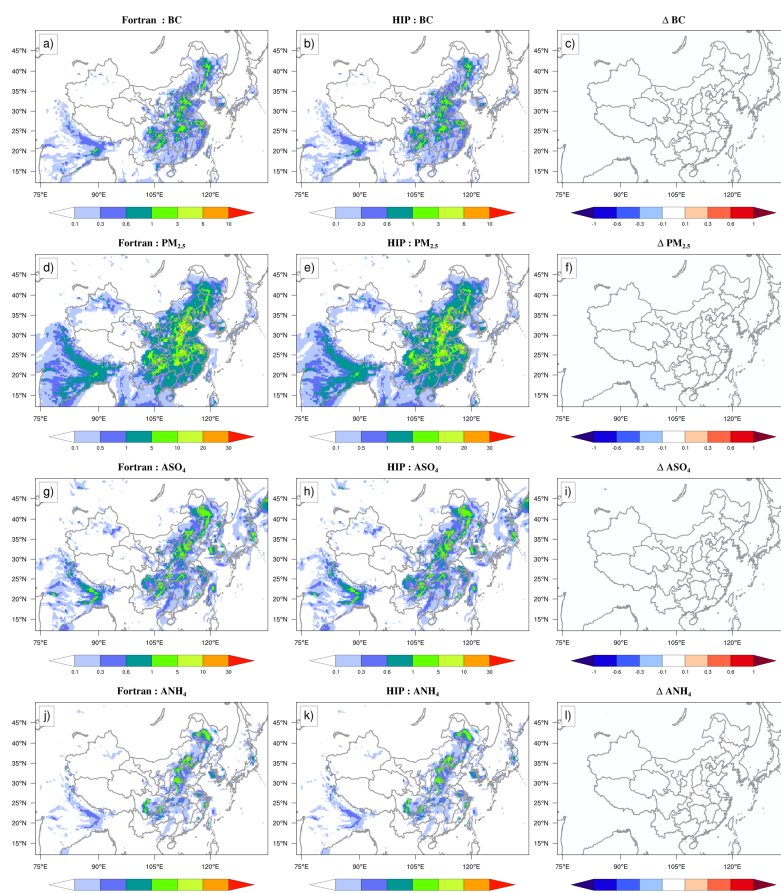


475

476 **Figure 6.** HONO, SO_2 , NO_2 , and NH_3 concentrations outputted by the EPICC-Model for the Fortran



477 and HIP-Opt2. versions. Panels (a), (d), (g), and (j) are from the Fortran version. Panels(b), (e), (h),
 478 and (k) are from the HIP-Opt2. version. Panels (c), (f), (i), and (l) are the absolute errors (AEs)
 479 between the Fortran and HIP-Opt2. Versions.



480
 481 **Figure 7.** BC, PM_{2.5}, ASO₄, and ANH₄ concentrations outputted by the EPIC-Model for the
 482 Fortran and HIP-Opt2. versions. Panels (a), (d), (g), and (j) are from the Fortran version. Panels(b),
 483 (e), (h), and (k) are from the HIP-Opt2. version. Panels (c), (f), (i), and (l) are the absolute errors
 484 (AEs) between the Fortran and HIP-Opt2 versions.

485

486 **Table 5.** The root mean square error (RMSE) between the Fortran and HIP-Opt2. versions, standard
 487 deviation (std) of the Fortran version, and the RMSE and std ratio.

	RMSE	std	RMSE/std (%)
--	------	-----	--------------



HONO (ppbV)	2.1×10^{-5}	0.1	2.1×10^{-2}
SO₂ (ppbV)	3.3×10^{-5}	0.7	4.5×10^{-3}
NO₂ (ppbV)	2.8×10^{-4}	3.4	8.3×10^{-3}
NH₃ (ppbV)	9.6×10^{-5}	4.1	2.4×10^{-3}
BC ($\mu\text{g} \cdot \text{m}^{-3}$)	9.7×10^{-8}	0.2	4.8×10^{-5}
PM_{2.5} ($\mu\text{g} \cdot \text{m}^{-3}$)	4.4×10^{-6}	1.9	2.3×10^{-4}
ASO₄ ($\mu\text{g} \cdot \text{m}^{-3}$)	1.7×10^{-4}	0.2	7.0×10^{-2}
ANH₄ ($\mu\text{g} \cdot \text{m}^{-3}$)	1.1×10^{-4}	0.2	5.8×10^{-2}

488

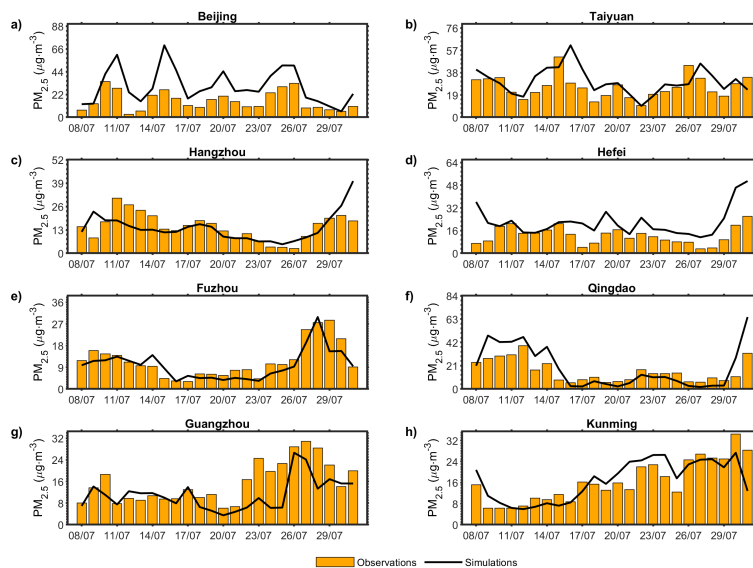
489 **4.2.2. Simulation performance verification of the EPICC-Model**

490 After parallelizing the HADVPPM program in the air quality model CAMx on GPU-like
491 accelerators, Cao et al. (2024) conducted comparative analyses of simulation results between
492 NVIDIA GPUs and domestic GPUs through offline and coupled testing approaches. Although both
493 experimental results indicated smaller computational errors introduced by China's domestic GPUs,
494 the study did not validate the discrepancies between CAMx simulation results and actual
495 observational data, particularly regarding model performance in real-case scenarios. To address this
496 gap, the current study integrates GPU-HADVPPM4HIP V1.0 into the EPICC-Model V1.6.0 and
497 performs one-month real-case simulations following the experimental configuration described in
498 Sect. 4.1. This serves dual purposes: firstly, to verify the computational stability of EPICC-Model
499 V1.6.0 in cross-architecture heterogeneous cluster environments using China's domestic hardware,
500 and secondly, to evaluate the model's pollutant simulation performance through observational
501 validation. For observational data, we collected PM_{2.5} and O₃ observations from major Chinese
502 cities including Beijing, Taiyuan, Hangzhou, Hefei, Fuzhou, Qingdao, Guangzhou, and Kunming,
503 implementing quality control procedures following the methodology established by Wu et al. (2018).
504 Regarding simulation data, we extracted model outputs from grid cells in the d02 domain
505 corresponding to the geographical coordinates of monitoring stations.

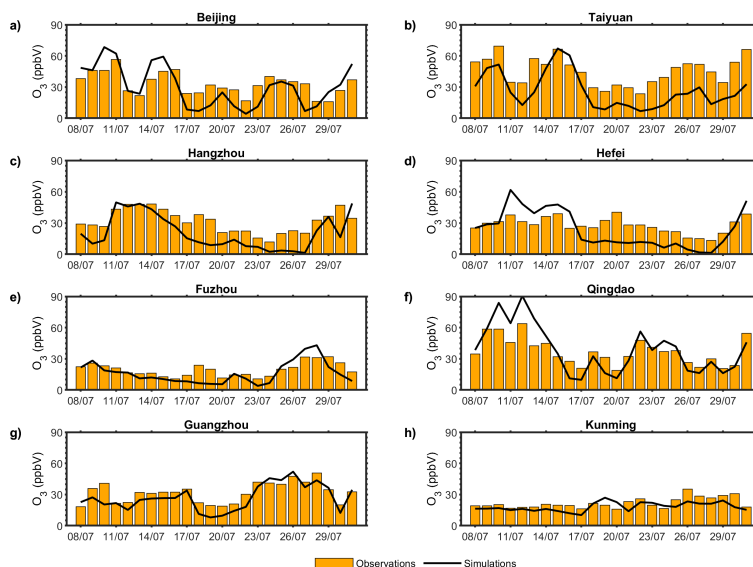
506 Figure 8 and Figure 9 along with Table 6, present the time series comparisons between daily
507 simulated and observed PM_{2.5} and O₃ concentrations, as well as relevant statistical metrics,
508 following the one-month simulation after coupling GPU-HADVPPM4HIP V1.0 to the EPICC-
509 Model V1.6.0. The formulas for calculating statistical parameters are detailed in the Supplementary
510 Materials. For daily PM_{2.5} concentrations, the EPICC-Model V1.6.0 demonstrated robust simulation



511 performance across most cities, with slight overestimations observed in Beijing and Hefei. Notably,
 512 the model effectively captured the temporal variations of PM_{2.5} in Fuzhou, achieving a root mean
 513 square error (RMSE) of $3.9 \mu\text{g} \cdot \text{m}^{-3}$ and a correlation coefficient (r) of 0.89. Across the eight cities,
 514 the mean RMSE and r between simulated and observed PM_{2.5} were $9.4 \mu\text{g} \cdot \text{m}^{-3}$ and 0.70,
 515 respectively. Regarding daily maximum 8-hour average (MDA8) O₃ concentrations, the model
 516 exhibited minor underestimations at certain times in Taiyuan, and Hefei but performed well in
 517 Beijing, Fuzhou, Qingdao, and Guangzhou. In Qingdao and Guangzhou, the correlation coefficient
 518 between simulated and observed O₃ reached 0.91 and 0.87, with RMSE values of 12.2 ppbV and
 519 8.7 ppbV, respectively. The mean RMSE and r for O₃ across eight cities were 12.1 ppbV and 0.76.
 520 These statistical results indicate that the EPIC-Model V1.6.0, integrated with GPU-
 521 HADVPPM4HIP V1.0, reasonably reproduces the spatiotemporal characteristics of PM_{2.5} and O₃
 522 concentrations on China's domestic heterogeneous computing clusters.



523
 524 **Figure 8.** Time series of daily observed and simulated PM_{2.5} concentrations in major cities of China
 525 on 8-31 July, 2021. Panels (a) - (h) represent Beijing, Taiyuan, Hangzhou, Hefei, Fuzhou, Qingdao,
 526 Guangzhou, and Kunming.



527

528 **Figure 9.** Time series of observed and simulated MAD8 O₃ concentrations in major cities of China
 529 on 8-31 July, 2021. Panels (a) - (h) represent Beijing, Taiyuan, Hangzhou, Hefei, Fuzhou, Qingdao,
 530 Guangzhou, and Kunming.

531

532 **Table 6.** The statistics of the PM_{2.5} and O₃ simulations of EPIC-Model over different eight cities.

	PM _{2.5}		O ₃	
	RMSE ($\mu\text{g} \cdot \text{m}^{-3}$)	r	RMSE (ppbV)	r
Beijing	17.3	0.84	13.2	0.77
Taiyuan	11.0	0.54	21.6	0.78
Hangzhou	7.4	0.50	14.8	0.80
Hefei	12.4	0.54	14.2	0.70
Fuzhou	3.9	0.87	7.5	0.77
Qingdao	11.0	0.90	12.2	0.91
Guangzhou	7.2	0.65	7.8	0.87
Kunming	5.3	0.76	6.0	0.46
Average	9.4	0.70	12.1	0.76

533

534 The underestimations or overestimations of PM_{2.5} and O₃ simulations observed in specific
 535 cities and periods are primarily attributable to two factors. First, the coarse model resolution—the
 536 horizontal resolution of the d02 domain in this experiment is 15 km—hindered the accurate
 537 representation of terrain features, meteorological variables, and spatial variations in emission
 538 sources. Second, discrepancies exist between the baseline year of the emission inventory and the
 539 simulation year. Specifically, the MEIC emission inventory used in this study is based on 2019 data,
 540 whereas the simulation year is 2021. For cities with stringent pollution control policies, such as
 Beijing, the 2019 MEIC inventory may inadequately reflect actual emission reductions achieved by



541 2021, particularly for pollutants under strict abatement measures. This discrepancy could lead to
542 overestimated simulated concentrations.

543 **4.3. Computational performance analysis**

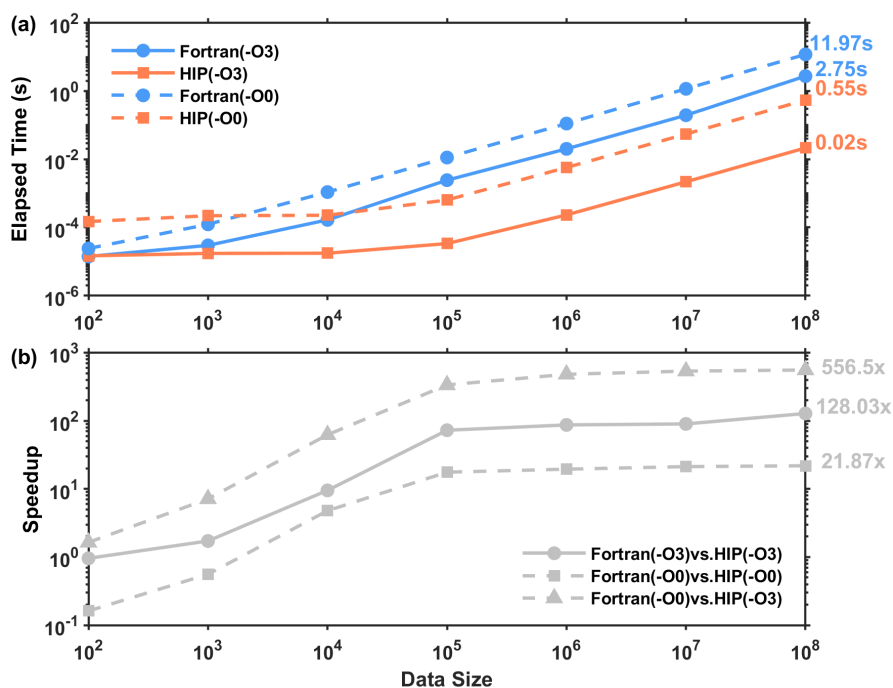
544 **4.3.1. Offline performance comparison**

545 A comparison of offline computational results indicates that the computational errors
546 introduced during the Fortran-to-HIP heterogeneous porting process under both compilation settings
547 are minimal. Specifically, the HADVPPM program exhibits small discrepancies on the order of 10^{-7}
548 when ported from CPU to GPU-like architectures. That is, differences in the computational results
549 only begin to appear at the seventh decimal place. Based on the verified consistency of offline results,
550 the computational performance of the HADVPPM program was further evaluated on domestic CPU
551 and GPU-like under different compilation configurations. To achieve this, Fortran-based test
552 programs were implemented to generate randomized input arrays with varying scales for both
553 Fortran and HIP versions of the HADVPPM program. These random arrays are one-dimensional
554 and contain 10^2 , 10^3 , 10^4 , 10^5 , 10^6 , 10^7 , and 10^8 random numbers. Figure 10 illustrates the
555 computational time and speedup ratios of the HADVPPM program across different compilation
556 options and data scales on domestic CPU and GPU-like. It is explicitly stated that the execution
557 time measurements for the HADVPPM program on the GPU-like exclusively account for the
558 computational time of the kernel functions on the device, while overheads such as GPU memory
559 allocation and host-device data transfer are excluded.

560 As illustrated in Figure 10(a), under both -O0 and -O3 compilation flags, the SIMD vectorized
561 parallel computing advantages of the GPU-like accelerators become prominent for large data scales
562 exceeding 10^4 , demonstrating significantly higher computational efficiency compared to domestic
563 CPU. At a data scale of 10^8 with the -O0 flag, the Fortran-based HADVPPM program required
564 approximately 11.97 seconds to complete computation on the CPU, while the HIP version on the
565 GPU-like achieved the same task in 0.55 seconds, yielding a speedup ratio of 21.87x. The -O3
566 compilation flag further enhances computational efficiency through automated code optimization,
567 albeit at a slight cost to numerical precision. At the same 10^8 data scale, the Fortran version on the
568 CPU required 2.75 seconds, whereas the HIP version on the GPU completed computations in 0.02
569 seconds, achieving a remarkable speedup of 128.03x. Notably, the HIP version compiled with -O3



570 exhibited 556.5x higher efficiency than the Fortran version compiled with -O0 flag. However, for
 571 smaller data scales such as 10^2 or 10^3 , the architectural advantages of the GPU-like diminish, with
 572 computational efficiency comparable to or even lower than that of a CPU. For instance, at a 10^2
 573 scale with the -O3 option, the GPU's performance matched the CPU (speedup ratio is approximately
 574 1). Under the -O0 option, the GPU's speedup dropped to 0.16x, indicating inferior efficiency relative
 575 to the CPU. These results underscore that the GPU-like is well-suited for large-scale matrix parallel
 576 computing tasks.



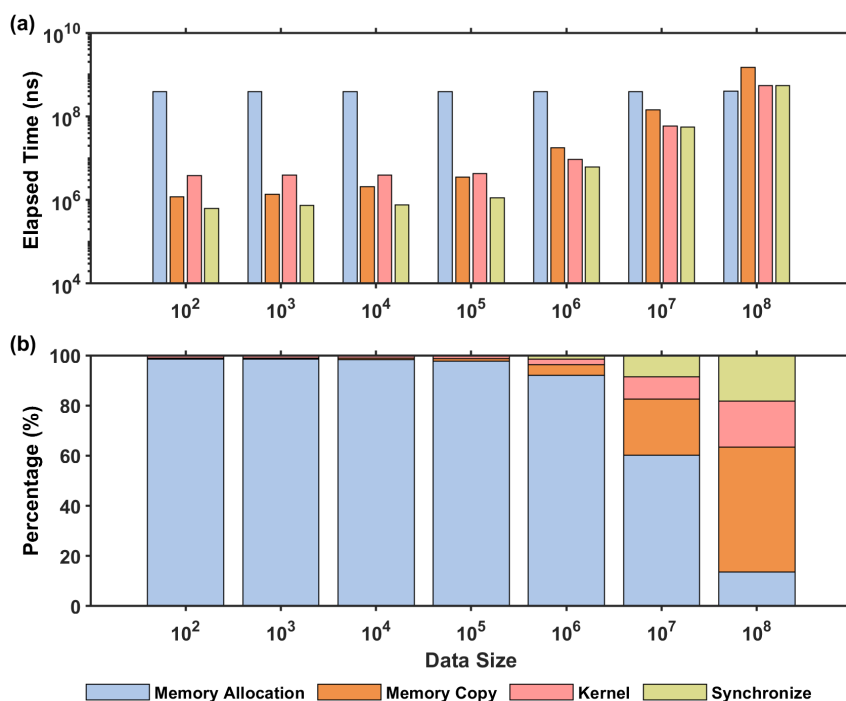
577
 578 **Figure 10.** The offline computational time (a) and speedup ratios (b) of the HADVPPM program
 579 across different compilation options and data scales on domestic CPU and GPU-like.

580 As described in Sect. 3.1, a complete heterogeneous computation of GPU-HADVPPM4HIP
 581 V1.0 on the GPU-like involves critical processes such as GPU memory allocation, data transfer
 582 between CPU and GPU, and kernel launching for parallel computing. To quantify the overheads of
 583 these processes, the time consumption and proportional contributions was analyzed under the -O0
 584 compilation flag and varying data scales. In Figure 11., the label of “Memory Allocation”, “Memory
 585 Copy”, “Kernel”, and “Synchronize” represent the processes of memory allocation on the GPU,



586 bidirectional data transfer between CPU and GPU, kernel launch and parallel computation, as well
587 as thread synchronization within the kernel, respectively.

588 As shown in Figure 11., for data scales smaller than 10^6 , memory allocation dominates the time
589 consumption, exceeding 90% of the total execution time and significantly surpassing the durations
590 of the other three processes. Notably, the time required for memory allocation remains
591 approximately 0.4 seconds regardless of increases in data scale. The dominance of memory
592 allocation at small data scales highlights its fixed overhead nature. This suggests that memory
593 allocation is largely independent of data volume. While negligible in large-scale computations, this
594 fixed cost becomes a critical bottleneck for small-scale tasks. When the data scale surpasses 10^5 , the
595 overhead of memory copy rises rapidly, with a growth rate higher than those of kernel execution
596 and synchronization processes. At a data scale of 10^8 , memory copy accounts for approximately 50%
597 of the total time and exhibits a tendency for further increase. Under this condition, the time
598 contributions of memory allocation, kernel execution, and synchronization are approximately 14%,
599 18%, and 18%, respectively. The rapid escalation of memory copy overhead underscores the
600 limitations of host-device data transfer bandwidth. The growth rate of memory copy time implies
601 that the data transfer between the CPU and GPU becomes one of the primary factors influencing the
602 computational performance of programs in heterogeneous cluster environments, and the I/O-bound
603 workloads often underutilize GPU compute capabilities. Future efforts could focus on reducing
604 communication overhead through strategies such as unified memory architecture and asynchronous
605 communication. Additionally, integrating mixed-precision methods (Vaña et al., 2017), converting
606 variables with minimal impact on simulation results from double-precision to single-precision,
607 could further enhance data transfer efficiency between CPUs and GPUs.



608

609 **Figure 11.** Time consumption (a) and proportional contributions (b) for memory allocation,
 610 memory copy, kernel, and synchronize process under the -O0 compilation flag and varying data
 611 scales.

612 4.3.2. Coupling performance comparison

613 Following the integration of GPU-HADVPPM4HIP V1.0 into the EPICC-Model V1.6.0,
 614 computational efficiency on the EarthLab was improved through communication optimization
 615 described in Sect. 3.2 and enhanced thread and block collaborative indexing detailed in Sect. 3.3.
 616 Furthermore, the hybrid parallelization scheme outlined in Sect. 3.4 was employed to extend the
 617 parallel computing scalability of the EPICC-Model V1.6.0 on the EarthLab. As introduced in Sect.
 618 3.1, three model versions were established: the baseline unoptimized version HIP-Ori, the
 619 communication-optimized version HIP-Opt1, and the collaboratively indexed version HIP-Opt2
 620 Figure 12(a) displays the average elapsed time required for a 28-hour simulation across these three
 621 versions. To ensure comparability, all tests adopted identical hardware configurations, including the
 622 MPI+HIP hybrid parallelization scheme with 10 CPU processes and 10 GPU-like accelerators, and
 623 were compiled using the -O3 optimization flag.

624 The GPU-HADVPPM4HIP V1.0 module features triple nested loops, and each invocation by



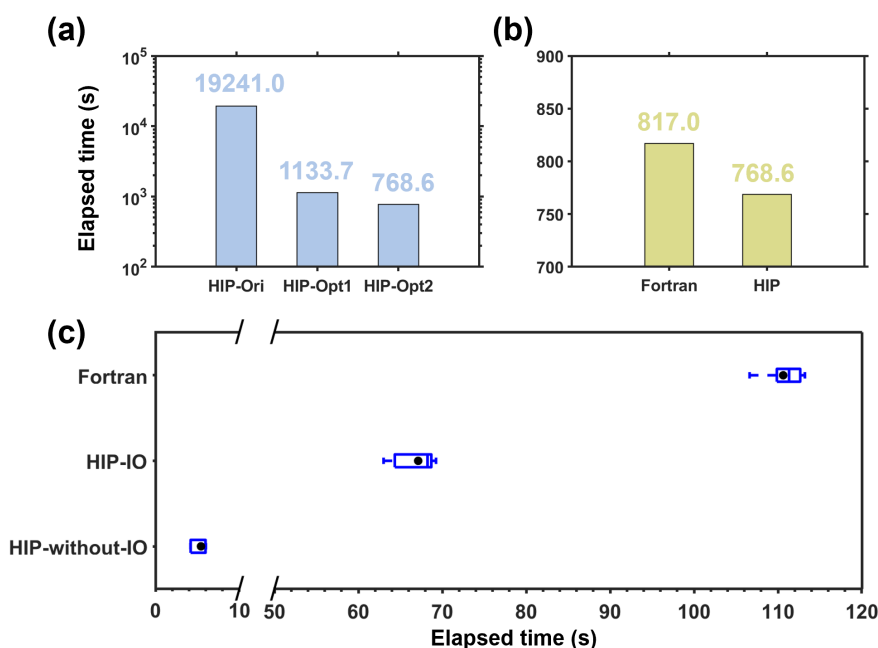
625 the EPICC-Model V1.6.0 necessitates data transfer between the CPU and GPU-like. Frequent CPU-
626 GPU data transfer severely compromised the computational efficiency of the EPICC-Model V1.6.0
627 on the EarthLab. In the HIP-Ori unoptimized version, completing a 28-hour simulation required
628 149.7 hours, with an average hourly elapsed time of 19,241.0 seconds, equivalent to approximately
629 5.3 hours, reflecting critically low efficiency. After implementing communication optimizations in
630 the HIP-Opt1 version, the communication frequency between the CPU and GPU-like was reduced
631 from roughly 9.8 million to 10. This optimization drastically decreased the total elapsed time from
632 149.7 hours for HIP-Ori to 8.8 hours for HIP-Opt1, while the average hourly elapsed time dropped
633 from 19,241.0 seconds to 1,133.7 seconds, achieving a 17.0x improvement in computational
634 efficiency. Subsequently, collaborative indexing of threads and blocks was applied to parallelize the
635 two-dimensional grid computations in the EPICC-Model V1.6.0. Compared to HIP-Opt1, the HIP-
636 Opt2 version further reduced the total elapsed time from 8.8 hours to 6.0 hours, with the average
637 hourly elapsed time decreasing from 1,133.7 seconds to 768.6 seconds. This enhancement delivered
638 an additional 1.5x efficiency gain. Cumulatively, these optimizations elevated the computational
639 efficiency of the EPICC-Model V1.6.0 on the EarthLab by approximately 25.0x.

640 By significantly improving the computational efficiency of the model on heterogeneous
641 clusters through reducing the communication frequency between the CPU and GPU, we aim to
642 emphasize the fundamental architectural differences between CPUs and GPUs. In recent years,
643 although the substantial increase in GPU computational performance has opened up new directions
644 for enhancing the efficiency of numerical models, the historical development of modeling
645 frameworks means that early numerical models were predominantly designed for CPU architectures.
646 Achieving parallel computation of numerical models on GPUs cannot be accomplished merely
647 through straightforward code translation or by relying solely on heterogeneous programming
648 interfaces. Instead, beyond adapting the model to GPUs, it is necessary to redesign the
649 computational workflow, loop structures, and logical organization in accordance with GPU
650 architectural characteristics, so as to fully leverage the powerful parallel computing capabilities of
651 GPUs.

652 Figure 12(b) further compares the computational performance between the original Fortran-
653 based version and the HIP-Opt2 version. In the figure, the label Fortran refers to the legacy CPU



654 cluster implementation of the EPIC-Model V1.6.0, while HIP represents the HIP-Opt2 version.
 655 Following the parallelization of the advection module on GPU-like accelerators, the EPIC-Model
 656 V1.6.0 demonstrated superior computational efficiency on heterogeneous clusters compared to its
 657 Fortran counterpart on conventional CPU clusters. The total elapsed time for a 28-hour simulation
 658 decreased from 2,287.4 seconds for the Fortran version to 2,152.0 seconds for the HIP version, with
 659 the average hourly elapsed time reduced from 817.0 seconds to 768.6 seconds.



660
 661 **Figure 12.** (a) The average hourly elapsed time required for a 28-hour simulation across HIP-Ori,
 662 HIP-Opt1, Opt2 versions; (b) the average hourly elapsed time required for a 28-hour simulation
 663 between the Fortran and HIP-Opt2 versions; (c) the hourly elapsed time required for a 28-hour
 664 simulation across the original Fortran-based HADVPPM program on CPUs, the GPU-
 665 HADVPPM4HIP V1.0 with CPU-GPU data transfer, and the GPU-HADVPPM4HIP V1.0 without
 666 data transfer. The black dots represent the average values, and the unit is seconds (s).

667 As demonstrated by the timing analysis of key processes in Sect. 4.3.1 for the offline
 668 heterogeneous computation of GPU-HADVPPM4HIP V1.0 on GPU-like accelerators, memory
 669 copying, specifically data transfer between the CPU and GPU, emerges as the dominant factor
 670 influencing parallel computational efficiency on heterogeneous clusters when handling large-scale



671 datasets. This overhead surpasses the time spent on kernel function parallelization. To quantify this
672 effect, the computational time of GPU-HADVPPM4HIP V1.0 was separately evaluated on GPU-
673 like accelerators under two scenarios: (1) including CPU-GPU data transfer and (2) excluding data
674 transfer (kernel-only computation). In Figure 12 (c), the y-axis labels Fortran, HIP-IO, and HIP-
675 without-IO correspond to the computational time of the original Fortran-based HADVPPM module
676 on CPUs, the GPU-HADVPPM4HIP V1.0 with CPU-GPU data transfer on GPU-like accelerators,
677 and the GPU-HADVPPM4HIP V1.0 without data transfer on GPU-like accelerators, respectively.
678 All tests were compiled with the -O3 optimization flag, and timing metrics were averaged over a
679 28-hour simulation, focusing on the hourly computational cost of the advection module.

680 The original Fortran-based advection module required an average of 110.6 seconds per
681 simulated hour on CPUs. After heterogeneous porting and parallel optimization to GPU-like
682 accelerators, the GPU-HADVPPM4HIP V1.0 with data transfer achieved an average time of 67.1
683 seconds per simulation hour, representing a 39.3% improvement in computational efficiency. When
684 excluding data transfer, the kernel-only GPU-HADVPPM4HIP V1.0 reduced the average time to
685 5.4 seconds per simulation hour, achieving a 20.5x acceleration compared to the Fortran version.
686 These results underscore that while GPU-like accelerators deliver substantial computational power,
687 the efficiency of CPU-GPU data transfer critically constrains overall performance on heterogeneous
688 clusters. To mitigate this bottleneck, future work will focus on retaining I/O operations on CPUs
689 while porting the entire physicochemical integration module (excluding I/O) to GPUs for parallel
690 computation. This strategy is expected to reduce the impact of inter-device data transfer and further
691 enhance scalability.

692 **5. Conclusions and discussion**

693 In recent years, the rapid advancement in the computational performance of GPUs has provided
694 novel approaches and hardware foundations for improving the computational efficiency of air
695 quality models. Building upon the heterogeneous porting and parallel optimization technology
696 system for air quality model, this study further implements parallel computing of the advection
697 module in the EPIC-Model air quality modeling system on GPU-like accelerators, validating the
698 feasibility of the heterogeneous porting framework. The study involves three key technical



699 improvements. The first is restructuring the original Fortran code of the advection module using
700 standard C language programming. The second is porting the advection module to GPU-like
701 accelerators through HIP heterogeneous programming technology, in addition, computational
702 efficiency was enhanced through optimizing CPU-GPU data transfer frequency reduction,
703 coordinated indexing of threads and blocks, and hybrid parallelization strategies. These
704 optimizations collectively improved both the computational performance of the advection module
705 and the parallel computing scalability of the EPICC-Model V1.6.0 on the EarthLab.

706 This study systematically conducted comparative efficiency analyses by the validation of
707 computational consistency in GPU-HADVPPM4HIP V1.0 through offline testing methodologies,
708 as well as the verification of EPICC-Model's pollutant simulation performance on the EarthLab.
709 Initial benchmarking compared the offline computational efficiency between the original Fortran-
710 based HADVPPM program on domestic CPUs and the GPU-HADVPPM4HIP V1.0
711 implementation. The results demonstrated that the -O3 compiler optimization flag significantly
712 enhanced GPU-HADVPPM4HIP's computational efficiency, with acceleration effects becoming
713 more pronounced at larger data scales. Specifically, at 10^8 data size configuration, GPU-
714 HADVPPM4HIP V1.0 achieved a maximum 556.5x speedup over the Fortran baseline using default
715 -O0 compilation, while maintaining a 128.0x speedup advantage even when both implementations
716 employed -O3 optimization. Further profiling of GPU-HADVPPM4HIP's heterogeneous
717 computation on GPU-like accelerators revealed critical characteristics: Memory copy operations
718 (i.e., CPU-GPU data transfers) exhibited elapsed time increases rapidly with data size increasing,
719 accounting for approximately 50% of total computation time at 10^8 data size with a continuing
720 upward trend. This observation underscores data transfer efficiency as a critical bottleneck for high-
721 resolution air quality simulations in heterogeneous computing environments.

722 Coupling GPU-HADVPPM4HIP V1.0 into EPICC-Model V1.6.0 with data transfer
723 optimizations and thread-block coordinated indexing strategies yielded system-level performance
724 improvements of 17.0x and 1.5x respectively on the EarthLab. The detailed module-level analysis
725 demonstrated that GPU-HADVPPM4HIP V1.0 achieved 39.3% efficiency enhancement over the
726 original Fortran advection module when accounting for CPU-GPU data transfer overheads,
727 escalating to a 20.5x acceleration when excluding data transfer costs. These findings quantitatively



728 validate the substantial impact of CPU-GPU data transfer efficiency on the operational performance
729 of air quality models in heterogeneous computing architectures.

730 There remains considerable room for improving the computational performance of GPU-
731 HADVPPM4HIP V1.0 within the EPICCC-Model in this study. First, constrained by the thermal
732 dissipation limits of transistors, the growth in computational performance of CPU processors has
733 slowed and is gradually approaching its physical limits. In recent years, however, GPU processors
734 have continued to achieve substantial performance gains due to their architectural advantages, and
735 heterogeneous supercomputing architectures centered on "CPU+GPU" now dominate the landscape
736 of advanced high-performance computing systems worldwide. Taking the 66th TOP500 list released
737 in November 2025 as an example, 9 out of the top 10 supercomputers adopt heterogeneous
738 architectures. Hence, heterogeneous computing represents the primary direction for the future
739 development of supercomputing. That said, the advancement of numerical models relies
740 fundamentally on the support of supercomputing capabilities. The team awarded the Gordon Bell
741 Prize for Climate Modeling in November 2025 leveraged the powerful computational capacity of
742 NVIDIA GH200 GPUs to accomplish the world's first 1.25 km ultra-high-resolution Earth system
743 simulation (Klocke et al., 2025). To achieve kilometer-scale, ultra-high-resolution simulations of
744 atmospheric pollution, it is both essential and imperative to adapt air quality models to
745 heterogeneous supercomputing environments and enable efficient parallel computing.

746 Moreover, results from offline computational performance tests in this study and coupled
747 performance tests in Cao et al. (2024) consistently show that the acceleration effect of GPUs
748 becomes more pronounced as the computational scale increases, which sufficiently demonstrates
749 the potential of GPUs in large-scale, high-resolution application scenarios. The simulation case
750 configured in this study, with a resolution of only 15 km over the whole of China, may not yet fully
751 reflect the parallel computing advantages of GPUs. In the future, higher-resolution simulation cases
752 will be designed to better highlight the performance advantages of GPU acceleration.

753 Finally, in future work, we will employ further optimization techniques to enhance the
754 computational performance of the advection module and other computationally intensive modules
755 in the EPICCC-Model on domestic heterogeneous clusters, including but not limited to:

756 (1) Firstly, priority should be given to optimizing CPU-GPU data transfer efficiency by reducing



757 communication overhead through strategies such as unified memory architecture, asynchronous
758 communication, mixed-precision methods, and minimizing non-essential variable I/O in air
759 quality forecasting.

760 (2) Second, while GPU-accelerated modules including the gas-phase chemistry module (Cao et al.,
761 2025) and advection module have been individually developed, their systematic integration into
762 EPICCC-Model requires architectural refinement to increase GPU code coverage. We will
763 analyze the computational characteristics of the code in other modules of the model. For code
764 segments involving iterative computations, we will first decouple the iterative computations by
765 creating intermediate variables, thereby eliminating dependencies between successive
766 calculation steps. Subsequently, the computational code of other modules after iterative
767 decoupling will be rewritten in the form of Kernel functions to increase the proportion of code
768 executed on the GPU.

769 (3) Finally, in current heterogeneous architecture supercomputing systems, the number of CPU
770 processes within a computing node typically exceeds the number of GPUs. Employing the
771 current matching scheme of one CPU process to one GPU accelerator results in the waste of
772 remaining CPU computing resources. In the future, on one hand, while avoiding data
773 transmission competition between the CPU and GPU, we will consider designing a more
774 sophisticated mechanism for matching multiple CPU processes with a single GPU. On the other
775 hand, drawing on Cao et al. (2024), we plan to introduce an OpenMP shared-memory parallel
776 scheme into the EPICCC-Model. Through multi-level hybrid parallelism, while porting
777 computationally intensive modules to the GPU for parallel computing, other modules running
778 on the CPU will utilize OpenMP multithreading parallelism, thereby fully leveraging CPU
779 computing resources.

780 Consequently, achieving efficient parallel computation of air quality models on GPU
781 accelerators necessitates collaboration between researchers well-versed in the intrinsic
782 physicochemical mechanisms of the models and engineers with expertise in GPU hardware and
783 software. Herein, we appeal for the participation of engineers specializing in GPU software and
784 hardware development to jointly advance the progress of air quality modeling.

785



786 **Code and data availability.** The source codes of EPICC-Model V1.6.0 are available online via
787 ZENODO (<https://zenodo.org/records/20303367>, EPICC-Model Working Group, 2026). The GPU
788 acceleration code of HADVPPM scheme in the EPICC-Model, datasets, and the offline test code
789 related to this paper are available online via ZENODO (<https://zenodo.org/records/20605319>, Cao
790 et al., 2026).

791

792

793 **Author contributions.** KC and QW refactored the existing code, visualization, and prepared the
794 materials. QW, XT, JZ, and ZW planned and organized the project. KC, QW, JinL, XC, HC, WW,
795 and LK optimized the GPU-based codes. HC, WW, HW, and JieL prepared the data and conducted
796 the simulation. KC, QW, TX, XC, WW, JieL, JZ, and ZW carried out formal analysis of the model
797 results. KC, QW, TX, JinL, XC, HC, WW, LK, and ZW took part in the discussion.

798

799

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

801

802

803 **Acknowledgements.** The Strategic Priority Research Program of the Chinese Academy of Sciences
804 (grant no. XDB0760401), the National Key Research and Development Program of China (grant
805 no. 2023YFC3705705), the State Key Laboratory of Atmospheric Environment and Extreme
806 Meteorology (grant no. 2024QN08), the National Natural Science Foundation of China (grant nos.
807 42507147 and 42377105), and the Key Research and Development Program of Henan Province
808 (grant no. 241111212300) funded this work. The authors would like to thank for the technical
809 support of the National large Scientific and Technological Infrastructure “Earth System Numerical
810 Simulation Facility” (<https://estr.cn/31134.02.EL>).

811

812

813 **Financial support.** This research has been supported by the Strategic Priority Research Program of
814 the Chinese Academy of Sciences (grant no. XDB0760401), the National Key Research and



815 Development Program of China (grant no. 2023YFC3705705), the State Key Laboratory of
816 Atmospheric Environment and Extreme Meteorology (grant no. 2024QN08), and the National
817 Natural Science Foundation of China (grant nos. 42507147 and 42377105).

818 **Reference**

819 Adani, M., D'Isidoro, M., Mircea, M., Guarnieri, G., Vitali, L., D'Elia, I., Ciancarella, L., Gualtieri,
820 M., Briganti, G., Cappelletti, A., Piersanti, A., Stracquadiano, M., Righini, G., Russo, F.,
821 Cremona, G., Villani, M. G., and Zanini, G.: Evaluation of air quality forecasting system
822 FORAIR-IT over Europe and Italy at high resolution for year 2017, *Atmos. Pollut. Res.*, 13,
823 10.1016/j.apr.2022.101456, 2022.

824 Alvanos, M. and Christoudias, T.: GPU-accelerated atmospheric chemical kinetics in the
825 ECHAM/MESSy (EMAC) Earth system model (version 2.52), *Geosci. Model Dev.*, 10, 3679-
826 3693, 10.5194/gmd-10-3679-2017, 2017.

827 AMD: ROCm Documentation Release 5.7.1, Advanced Micro Devices Inc.,
828 <https://rocm.docs.amd.com/en/docs-5.7.1/> (last access: 26 May 2025), 2023.

829 Appel, K. W., Bash, J. O., Fahey, K. M., Foley, K. M., Gilliam, R. C., Hogrefe, C., Hutzell, W. T.,
830 Kang, D., Mathur, R., Murphy, B. N., Napelenok, S. L., Nolte, C. G., Pleim, J. E., Pouliot, G.
831 A., Pye, H. O. T., Ran, L., Roselle, S. J., Sarwar, G., Schwede, D. B., Sidi, F. I., Spero, T. L.,
832 and Wong, D. C.: The Community Multiscale Air Quality (CMAQ) model versions 5.3 and
833 5.3.1: system updates and evaluation, *Geosci. Model Dev.*, 14, 2867-2897, 10.5194/gmd-14-
834 2867-2021, 2021.

835 Atkinson, R. W., Fuller, G. W., Anderson, H. R., Harrison, R. M., and Armstrong, B.: Urban ambient
836 particle metrics and health: a time-series analysis, *Epidemiology (Cambridge, Mass.)*, 21, 501-
837 511, 10.1097/EDE.0b013e3181debc88, 2010.

838 Byun, D. W. and Dennis, R.: Design artifacts in eulerian air quality models: Evaluation of the effects
839 of layer thickness and vertical profile correction on surface ozone concentrations, *Atmos.*
840 *Environ.*, 29, 105-126, [https://doi.org/10.1016/1352-2310\(94\)00225-A](https://doi.org/10.1016/1352-2310(94)00225-A), 1995.

841 Cao, K., Wu, Q., and Tang, Xiao.: The dataset of the manuscript "Enhancing the advection module
842 performance in the EPICC-Model V1.6.0 via GPU-HADVPPM4HIP V1.0 coupling and GPU-



- 843 optimized strategies", Zenodo [data set], <https://doi.org/10.5281/zenodo.20605319>, 2026.
- 844 Cao, K., Wu, Q., Wang, L., Wang, N., Cheng, H., Tang, X., Li, D., and Wang, L.: GPU-HADVPPM
845 V1.0: a high-efficiency parallel GPU design of the piecewise parabolic method (PPM) for
846 horizontal advection in an air quality model (CAMx V6.10), *Geosci. Model Dev.*, 16, 4367-
847 4383, [10.5194/gmd-16-4367-2023](https://doi.org/10.5194/gmd-16-4367-2023), 2023.
- 848 Cao, K., Wu, Q., Wang, L., Guo, H., Wang, N., Cheng, H., Tang, X., Li, D., Liu, L., Li, D., Wu, H.,
849 and Wang, L.: GPU-HADVPPM4HIP V1.0: using the heterogeneous-compute interface for
850 portability (HIP) to speed up the piecewise parabolic method in the CAMx (v6.10) air quality
851 model on China's domestic GPU-like accelerator, *Geosci. Model Dev.*, 17, 6887-6901,
852 [10.5194/gmd-17-6887-2024](https://doi.org/10.5194/gmd-17-6887-2024), 2024.
- 853 Cao, K., Tang, X., Chen, H., Ma, J., Wu, Q., Wang, W., Chen, X., Li J., Wang, Z.: Porting and
854 Parallel Optimization of the Gas-phase Chemistry Module of the Air Quality Model EPICCC-
855 Model on China's Domestic Accelerators, *Frontiers of Data&Computing*,
856 [10.11871/jfdc.issn.2096-742X.2025.05.010](https://doi.org/10.11871/jfdc.issn.2096-742X.2025.05.010), 2025 (in Chinese).
- 857 Chai, Z., Zhang, H., Zhang, M., Tang, X., Zheng, W., Zhu, J., Zhou, G., Cao, J., and Zeng, Q.:
858 China's EarthLab—Forefront of Earth System Simulation Research, *Adv. Atmos. Sci.*, 38,
859 1611-1620, [10.1007/s00376-021-1175-y](https://doi.org/10.1007/s00376-021-1175-y), 2021.
- 860 Chang, J. S., Brost, R. A., Isaksen, I. S. A., Madronich, S., Middleton, P., Stockwell, W. R., and
861 Walcek, C. J.: A three-dimensional Eulerian acid deposition model: Physical concepts and
862 formulation, *J. Geophys. Res.-Atmos.*, 92, 14681-14700,
863 <https://doi.org/10.1029/JD092iD12p14681>, 1987.
- 864 Colella, P. and Woodward, P. R.: The Piecewise Parabolic Method (PPM) for gas-dynamical
865 simulations, *J. Comput. Phys.*, 54, 174-201, [https://doi.org/10.1016/0021-9991\(84\)90143-8](https://doi.org/10.1016/0021-9991(84)90143-8),
866 1984.
- 867 Crippa, M., Guizzardi, D., Pagani, F., Schiavina, M., Melchiorri, M., Pisoni, E., Graziosi, F.,
868 Muntean, M., Maes, J., Dijkstra, L., Van Damme, M., Clarisse, L., and Coheur, P.: Insights into
869 the spatial distribution of global, national, and subnational greenhouse gas emissions in the
870 Emissions Database for Global Atmospheric Research (EDGAR v8.0), *Earth Syst. Sci. Data*,
871 16, 2811-2830, [10.5194/essd-16-2811-2024](https://doi.org/10.5194/essd-16-2811-2024), 2024.



- 872 Crippa, M., Guizzardi, D., Butler, T., Keating, T., Wu, R., Kaminski, J., Kuenen, J., Kurokawa, J.,
873 Chatani, S., Morikawa, T., Pouliot, G., Racine, J., Moran, M. D., Klimont, Z., Manseau, P. M.,
874 Mashayekhi, R., Henderson, B. H., Smith, S. J., Suchyta, H., Muntean, M., Solazzo, E., Banja,
875 M., Schaaf, E., Pagani, F., Woo, J. H., Kim, J., Monforti-Ferrario, F., Pisoni, E., Zhang, J.,
876 Niemi, D., Sassi, M., Ansari, T., and Foley, K.: The HTAP_v3 emission mosaic: merging
877 regional and global monthly emissions (2000–2018) to support air quality modelling and
878 policies, *Earth Syst. Sci. Data*, 15, 2667-2694, 10.5194/essd-15-2667-2023, 2023.
- 879 Damian, V., Sandu, A., Damian, M., Potra, F., and Carmichael, G. R.: The kinetic preprocessor KPP-
880 a software environment for solving chemical kinetics, *Comput. Chem. Eng.*, 26, 1567-1579,
881 [https://doi.org/10.1016/S0098-1354\(02\)00128-X](https://doi.org/10.1016/S0098-1354(02)00128-X), 2002.
- 882 Elbern, H.: Parallelization and load balancing of a comprehensive atmospheric chemistry transport
883 model, *Atmos. Environ.*, 31, 3561-3574, [https://doi.org/10.1016/S1352-2310\(97\)00157-](https://doi.org/10.1016/S1352-2310(97)00157-X)
884 [X](https://doi.org/10.1016/S1352-2310(97)00157-X), 1997.
- 885 Emery, C., Baker, K., Wilson, G., and Yarwood, G.: Comprehensive Air Quality Model with
886 Extensions: Formulation and Evaluation for Ozone and Particulate Matter over the US,
887 *Atmosphere*, 15, 1158. <https://doi.org/10.3390/atmos15101158>, 2024.
- 888 EPICC-Model Working Group.: Description and evaluation of the Emission and atmospheric
889 Processes Integrated and Coupled Community (EPICC) Model version 1.0. *Adv. Atmos. Sci.*,
890 <http://www.iapjournals.ac.cn/aas/en/article/doi/10.1007/s00376-025-4384-y>, 2025.
- 891 EPICC-Model Working Group. : Emission and atmospheric Processes Integrated and Coupled
892 Community Model (EPICC-Model). Zenodo. <https://doi.org/10.5281/zenodo.20303367>, 2026.
- 893 Gao, Z. and Zhou, X.: A review of the CAMx, CMAQ, WRF-Chem and NAQPMS models:
894 Application, evaluation and uncertainty factors, *Environ. Pollut.*, 343,
895 10.1016/j.envpol.2023.123183, 2024.
- 896 Georgiou, G. K., Christoudias, T., Proestos, Y., Kushta, J., Hadjinicolaou, P., and Lelieveld, J.: Air
897 quality modelling in the summer over the eastern Mediterranean using WRF-Chem: chemistry
898 and aerosol mechanism intercomparison, *Atmos. Chem. Phys.*, 18, 1555-1571, 10.5194/acp-
899 18-1555-2018, 2018.
- 900 Georgiou, G. K., Christoudias, T., Proestos, Y., Kushta, J., Pikridas, M., Sciare, J., Savvides, C., and



- 901 Lelieveld, J.: Evaluation of WRF-Chem model (v3.9.1.1) real-time air quality forecasts over
902 the Eastern Mediterranean, *Geosci. Model Dev.*, 15, 4129-4146, 10.5194/gmd-15-4129-2022,
903 2022.
- 904 Guevara, M., Jorba, O., Tena, C., Denier van der Gon, H., Kuenen, J., Elguindi, N., Darras, S.,
905 Granier, C., and Pérez García-Pando, C.: Copernicus Atmosphere Monitoring Service
906 TEMPORal profiles (CAMS-TEMPO): global and European emission temporal profile maps
907 for atmospheric chemistry modelling, *Earth Syst. Sci. Data*, 13, 367-404, 10.5194/essd-13-
908 367-2021, 2021.
- 909 Gupta, M. and Mohan, M.: Validation of WRF/Chem model and sensitivity of chemical mechanisms
910 to ozone simulation over megacity Delhi, *Atmos. Environ.*, 122, 220-229,
911 10.1016/j.atmosenv.2015.09.039, 2015.
- 912 Hong, S.-Y., Noh, Y., and Dudhia, J.: A New Vertical Diffusion Package with an Explicit Treatment
913 of Entrainment Processes, *Mon. Weather Rev.*, 134, 2318-2341,
914 <https://doi.org/10.1175/MWR3199.1>, 2006.
- 915 Intel Software: Quick Reference Guide to Optimization with Intel C++ and Fortran Compilers v19,
916 Intel, [https://www.intel.cn/content/www/cn/zh/developer/articles/technical/intel-fortran-and-
917 c-compiler-documentation.html](https://www.intel.cn/content/www/cn/zh/developer/articles/technical/intel-fortran-and-c-compiler-documentation.html). (last access: 26 May 2025), 2020.
- 918 Kim, K.-H., Kabir, E., and Kabir, S.: A review on the human health impact of airborne particulate
919 matter, *Environ. Int.*, 74, 136-143, <https://doi.org/10.1016/j.envint.2014.10.005>, 2015.
- 920 Klocke, D., Frauen, C., Engels, J. F., Alexeev, D., Redler, R., Schnur, R., Haak, H., Kornblueh, L.,
921 Brüggemann, N., Chegini, F., Römmner, M., Hoffmann, L., Griessbach, S., Bode, M., Coles, J.,
922 Gila, M., Sawyer, W., Calotoiu, A., Budanaz, Y., Mazumder, P., Copik, M., Weber, B., Herten,
923 A., Bockelmann, H., Hoefler, T., Hohenegger, C., and Stevens, B.: Computing the Full Earth
924 System at 1km Resolution, *Proceedings of the International Conference for High Performance
925 Computing, Networking, Storage and Analysis*, 10.1145/3712285.3771789, 2025.
- 926 Kong, L., Tang, X., Wang, Z., Zhu, J., Li, J., Wu, H., Wu, Q., Chen, H., Zhu, L., Wang, W., Liu, B.,
927 Wang, Q., Chen, D., Pan, Y., Li, J., Wu, L., and Carmichael, G. R.: Changes in air pollutant
928 emissions in China during two clean-air action periods derived from the newly developed
929 Inversed Emission Inventory for Chinese Air Quality (CAQIEI), *Earth Syst. Sci. Data*, 16,



- 930 4351-4387, 10.5194/essd-16-4351-2024, 2024.
- 931 Li, J., Wang, Z., Zhuang, G., Luo, G., Sun, Y., and Wang, Q.: Mixing of Asian mineral dust with
932 anthropogenic pollutants over East Asia: a model case study of a super-duststorm in March
933 2010, *Atmos. Chem. Phys.*, 12, 7591-7607, 10.5194/acp-12-7591-2012, 2012.
- 934 Li, J., Chen, X., Wang, Z., Du, H., Yang, W., Sun, Y., Hu, B., Li, J., Wang, W., Wang, T., Fu, P., and
935 Huang, H.: Radiative and heterogeneous chemical effects of aerosols on ozone and inorganic
936 aerosols over East Asia, *Sci. Total Environ.*, 622-623, 1327-1342,
937 <https://doi.org/10.1016/j.scitotenv.2017.12.041>, 2018.
- 938 Li, M., Liu, H., Geng, G., Hong, C., Liu, F., Song, Y., Tong, D., Zheng, B., Cui, H., Man, H., Zhang,
939 Q., and He, K.: Anthropogenic emission inventories in China: a review, *Natl. Sci. Rev.*, 4, 834-
940 866, 10.1093/nsr/nwx150, 2017.
- 941 Linford, J. C., Michalakes, J., Vachharajani, M., and Sandu, A.: Automatic Generation of Multicore
942 Chemical Kernels, *IEEE T. Parall. Distr.*, 22, 119-131, 10.1109/tpds.2010.106, 2011.
- 943 Liu, J., Han, Y., Tang, X., Zhu, J., and Zhu, T.: Estimating adult mortality attributable to PM_{2.5}
944 exposure in China with assimilated PM_{2.5} concentrations based on a ground monitoring
945 network, *Sci. Total. Environ.*, 568, 1253-1262, <https://doi.org/10.1016/j.scitotenv.2016.05.165>,
946 2016.
- 947 Milton, L. A. and White, A. R.: The potential impact of bushfire smoke on brain health, *Neurochem.*
948 *Int.*, 139, 104796, 10.1016/j.neuint.2020.104796, 2020.
- 949 NVIDIA: CUDA C++ Programming Guide Version 10.2, NVIDIA Corporation,
950 https://docs.nvidia.com/cuda/archive/10.2/pdf/CUDA_C_Programming_Guide.pdf (last
951 access: 26 May 2025), 2020.
- 952 Nenes, A., Pandis, S. N., and Pilinis, C.: ISORROPIA: A New Thermodynamic Equilibrium Model
953 for Multiphase Multicomponent Inorganic Aerosols, *Aquat. Geochem.*, 4, 123-152,
954 10.1023/A:1009604003981, 1998.
- 955 Podrascanin, Z.: Setting-up a Real-Time Air Quality Forecasting system for Serbia: a WRF-Chem
956 feasibility study with different horizontal resolutions and emission inventories, *Environ. Sci.*
957 *Pollut. R.*, 26, 17066-17079, 10.1007/s11356-019-05140-y, 2019.
- 958 Quevedo, D., Do, K., Delic, G., Rodríguez-Borbón, J., Wong, B. M., and Ivey, C. E.: GPU



- 959 Implementation of a Gas-Phase Chemistry Solver in the CMAQ Chemical Transport Model,
960 ACS ES&T Air, 2, 226-235, 10.1021/acsestair.4c00181, 2025.
- 961 Sandu, A., Verwer, J. G., Van Loon, M., Carmichael, G. R., Potra, F. A., Dabdub, D., and Seinfeld,
962 J. H.: Benchmarking stiffode solvers for atmospheric chemistry problems-I. implicit vs explicit,
963 Atmos. Environ., 31, 3151-3166, [https://doi.org/10.1016/S1352-2310\(97\)00059-9](https://doi.org/10.1016/S1352-2310(97)00059-9), 1997.
- 964 Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Barker, D.M., Duda, M. G., Huang, X. Y.,
965 Wang, W., and Powers, J. G.: A Description of the Advanced Research WRF Version3
966 (No.NCAR/TN-475CSTR), University Corporation for Atmospheric Research, NCAR,
967 <https://doi.org/10.5065/D68S4MVH>, 2008.
- 968 Sun, J., Fu, J. S., Drake, J. B., Zhu, Q., Haidar, A., Gates, M., Tomov, S., and Dongarra, J.:
969 Computational Benefit of GPU Optimization for the Atmospheric Chemistry Modeling, J. Adv.
970 Model. Earth. Sy., 10, 1952-1969, <https://doi.org/10.1029/2018MS001276>, 2018.
- 971 Tang, Y., Campbell, P. C., Lee, P., Saylor, R., Yang, F., Baker, B., Tong, D., Stein, A., Huang, J.,
972 Huang, H.-C., Pan, L., McQueen, J., Stajner, I., Tirado-Delgado, J., Jung, Y., Yang, M.,
973 Bourgeois, I., Peischl, J., Ryerson, T., Blake, D., Schwarz, J., Jimenez, J.-L., Crawford, J.,
974 Diskin, G., Moore, R., Hair, J., Huey, G., Rollins, A., Dibb, J., and Zhang, X.: Evaluation of
975 the NAQFC driven by the NOAA Global Forecast System (version 16): comparison with the
976 WRF-CMAQ during the summer 2019 FIREX-AQ campaign, Geosci. Model Dev., 15, 7977-
977 7999, 10.5194/gmd-15-7977-2022, 2022.
- 978 Thompson, T. M. and Selin, N. E.: Influence of air quality model resolution on uncertainty
979 associated with health impacts, Atmos. Chem. Phys., 12, 9753-9762, 10.5194/acp-12-9753-
980 2012, 2012.
- 981 Top500: Supercomputing Top500 list, TOP500 international organization,
982 <https://www.top500.org/lists/top500/2024/11/> (last access: 26 May 2025), 2024.
- 983 Váňa, F., Düben, P., Lang, S., Palmer, T., Leutbecher, M., Salmond, D., and Carver, G.: Single
984 Precision in Weather Forecasting Models: An Evaluation with the IFS, Mon. Weather Rev.,
985 145, 495–502, <https://doi.org/10.1175/mwr-d-16-0228.1>, 2017.
- 986 Walcek, C. J. and Aleksic, N. M.: A simple but accurate mass conservative, peak-preserving, mixing
987 ratio bounded advection algorithm with FORTRAN code, Atmos. Environ., 32, 3863-3880,



- 988 [https://doi.org/10.1016/S1352-2310\(98\)00099-5](https://doi.org/10.1016/S1352-2310(98)00099-5), 1998.
- 989 Wang, H., Chen, H., Wu, Q., Lin, J., Chen, X., Xie, X., Wang, R., Tang, X., and Wang, Z.:
990 GNAQPMS v1.1: accelerating the Global Nested Air Quality Prediction Modeling System
991 (GNAQPMS) on Intel Xeon Phi processors, *Geosci. Model Dev.*, 10, 2891-2904,
992 [10.5194/gmd-10-2891-2017](https://doi.org/10.5194/gmd-10-2891-2017), 2017.
- 993 Wang, H., Lin, J., Wu, Q., Chen, H., Tang, X., Wang, Z., Chen, X., Cheng, H., and Wang, L.: MP
994 CBM-Z V1.0: design for a new Carbon Bond Mechanism Z (CBM-Z) gas-phase chemical
995 mechanism architecture for next-generation processors, *Geosci. Model Dev.*, 12, 749-764,
996 [10.5194/gmd-12-749-2019](https://doi.org/10.5194/gmd-12-749-2019), 2019.
- 997 Wang, P., Jiang, J., Lin, P., Ding, M., Wei, J., Zhang, F., Zhao, L., Li, Y., Yu, Z., Zheng, W., Yu, Y.,
998 Chi, X., and Liu, H.: The GPU version of LASG/IAP Climate System Ocean Model version 3
999 (LICOM3) under the heterogeneous-compute interface for portability (HIP) framework and its
1000 large-scale application, *Geosci. Model Dev.*, 14, 2781–2799, [https://doi.org/10.5194/gmd-14-](https://doi.org/10.5194/gmd-14-2781-2021)
1001 [2781-2021](https://doi.org/10.5194/gmd-14-2781-2021), 2021.
- 1002 Wang, W., Chen, H., Wang, Z., Li, J., Chen, X., Yu, F., Fan, X., Zhao, S., Hu, B., Wang, W., Tang,
1003 X., Wang, Z., Ge, B., and Wu, J.: Development and evaluation of photolysis and gas-phase
1004 reaction scheme in EPICCC-model: Impacts on tropospheric ozone simulation, *Atmos. Environ.*,
1005 359, <https://doi.org/10.1016/j.atmosenv.2025.121373>, 2025.
- 1006 Wesely, M. L.: Parameterization of surface resistances to gaseous dry deposition in regional-scale
1007 numerical models, *Atmos. Environ.*(1967), 23, 1293-1304, [https://doi.org/10.1016/0004-](https://doi.org/10.1016/0004-6981(89)90153-4)
1008 [6981\(89\)90153-4](https://doi.org/10.1016/0004-6981(89)90153-4), 1989.
- 1009 Wu, H., Tang, X., Wang, Z., Wu, L., Lu, M., Wei, L., and Zhu, J.: Probabilistic Automatic Outlier
1010 Detection for Surface Air Quality Measurements from the China National Environmental
1011 Monitoring Network, *Adv. Atmos. Sci.*, 35, 1522-1532, [10.1007/s00376-018-8067-9](https://doi.org/10.1007/s00376-018-8067-9), 2018.
- 1012 Wu, Q. Z., Xu, W. S., Shi, A. J., Li, Y. T., Zhao, X. J., Wang, Z. F., Li, J. X., and Wang, L. N.: Air
1013 quality forecast of PM₁₀ in Beijing with Community Multi-scale Air Quality Modeling (CMAQ)
1014 system: emission and improvement, *Geosci. Model Dev.*, 7, 2243-2259, [10.5194/gmd-7-2243-](https://doi.org/10.5194/gmd-7-2243-2014)
1015 [2014](https://doi.org/10.5194/gmd-7-2243-2014), 2014.
- 1016 Yarwood, G., Shi, Y., Beardsley, R.: Impact of cb6r5 mechanism changes on air pollutant modeling



1017 in Texas. Final report prepared for the Texas Commission on Environmental Quality, Austin,
1018 TX. https://www.tceq.texas.gov/airquality/airmod/project/pj_report_pm.html, 2020.

1019 Yang, J., Qu, Y., Chen, Y., Zhang, J., Liu, X., Niu, H., and An, J.: Dominant physical and chemical
1020 processes impacting nitrate in Shandong of the North China Plain during winter haze events,
1021 *Sci. Total Environ.*, 912, 169065, <https://doi.org/10.1016/j.scitotenv.2023.169065>, 2024.

1022 Zaveri, R. A. and Peters, L. K.: A new lumped structure photochemical mechanism for large-scale
1023 applications, *J. Geophys. Res.-Atmos.*, 104, 30387-30415, 10.1029/1999jd900876, 1999.

1024 Zhang, J., Lian, C., Wang, W., Ge, M., Guo, Y., Ran, H., Zhang, Y., Zheng, F., Fan, X., Yan, C.,
1025 Daellenbach, K. R., Liu, Y., Kulmala, M., and An, J.: Amplified role of potential HONO
1026 sources in O₃ formation in North China Plain during autumn haze aggravating processes,
1027 *Atmos. Chem. Phys.*, 22, 3275-3302, 10.5194/acp-22-3275-2022, 2022.

1028 Zhang, L., Brook, J. R., and Vet, R.: A revised parameterization for gaseous dry deposition in air-
1029 quality models, *Atmos. Chem. Phys.*, 3, 2067-2082, 10.5194/acp-3-2067-2003, 2003.

1030 Zhu, T., Tang, M., Gao, M., Bi, X., Cao, J., Che, H., Chen, J., Ding, A., Fu, P., Gao, J., Gao, Y., Ge,
1031 M., Ge, X., Han, Z., He, H., Huang, R.-J., Huang, X., Liao, H., Liu, C., Liu, H., Liu, J., Liu, S.
1032 C., Lu, K., Ma, Q., Nie, W., Shao, M., Song, Y., Sun, Y., Tang, X., Wang, T., Wang, T., Wang,
1033 W., Wang, X., Wang, Z., Yin, Y., Zhang, Q., Zhang, W., Zhang, Y., Zhang, Y., Zhao, Y., Zheng,
1034 M., Zhu, B., and Zhu, J.: Recent Progress in Atmospheric Chemistry Research in China:
1035 Establishing a Theoretical Framework for the "Air Pollution Complex", *Adv. Atmos. Sci.*,
1036 40, 1339-1361, 10.1007/s00376-023-2379-0, 2023.

1037