

Accurate load balancing accelerates Lagrangian simulation of water ages on distributed, multi-GPU platforms

CHEN YANG¹, Reed M. Maxwell¹, and Richard Valent²

¹Princeton University

²National Center for Atmospheric Research

November 22, 2022

Abstract

Water age is a fundamental descriptor of source, storage, and mixing of water parcels in a watershed. The Lagrangian, particle tracking, approach is a powerful tool for physically-based modeling of water age distributions, but its application has been hampered since it is computationally demanding. In this study, we present a parallel approach for particle tracking simulations. This approach uses multi-GPU with MPI parallelism based on domain decomposition. An inherent challenge of distributed parallelization of Lagrangian approaches is the disparity in computational work or load imbalance (LIB) among different processing elements (PEs). Here, load balancing (LB) schemes were proposed to dynamically balance the distribution of particles across PEs during runtime. In the followed hillslope simulations, LIB was observed in all LB-disabled runs, e.g., with a load ratio of 423.62% by using 2-GPU in LW_Shrub case. LB schemes then accurately balanced the load distribution and improved the parallel scaling. Additionally, the parallel approach showed excellent overall speedup: a 60-fold improvement using 4-GPU relative to the serial run. A regional scale application further demonstrated the LB performance. The parallel time used by 8-GPU without LB was 31.33% reduced after LB was activated. When increasing 8-GPU with LB to 16-GPU with LB, it showed parallel scalability by reducing the parallel time of ~50%. This work shows how massively parallel computing can be applied to particle tracking in water age simulations. It also demonstrates the practical importance of load balancing in this context, which enables the large-scale simulations with an increased complexity of flow paths.

1 **Accurate load balancing accelerates Lagrangian simulation of water**
2 **ages on distributed, multi-GPU platforms**

3
4 Chen Yang^{1*}, Reed M. Maxwell^{1,2*}, Richard Valent³

5
6 ¹Department of Civil and Environmental Engineering, Princeton University,
7 Princeton, NJ 08544, USA

8 ²High Meadows Environmental Institute, Princeton University, Princeton, NJ
9 08544, USA

10 ³Computational and Information Systems Laboratory, National Center for
11 Atmospheric Research, Boulder, CO 80305, USA

12
13
14 Corresponding author:

15 Chen Yang cy15@princeton.edu

16 Reed Maxwell reedmaxwell@princeton.edu

17
18
19
20
21 Key points:

- 22 • Massively parallel computing of Lagrangian water-age simulations is
23 realized
 - 24 • Mechanisms of load imbalance are identified and LB schemes are proposed
 - 25 • Parallel performance of the approach is demonstrated at the regional scale
- 26
27
28
29
30
31
32
33
34
35
36

37 **Abstract**

38 Water age is a fundamental descriptor of source, storage, and mixing of water parcels in a
39 watershed. The Lagrangian, particle tracking, approach is a powerful tool for physically-based
40 modeling of water age distributions, but its application has been hampered since it is
41 computationally demanding. In this study, we present a parallel approach for particle tracking
42 simulations. This approach uses multi-GPU with MPI parallelism based on domain
43 decomposition. An inherent challenge of distributed parallelization of Lagrangian approaches is
44 the disparity in computational work or load imbalance (LIB) among different processing
45 elements (PEs). Here, load balancing (LB) schemes were proposed to dynamically balance the
46 distribution of particles across PEs during runtime. In the followed hillslope simulations, LIB
47 was observed in all LB-disabled runs, e.g., with a load ratio of 423.62% by using 2-GPU in
48 LW_Shrub case. LB schemes then accurately balanced the load distribution and improved the
49 parallel scaling. Additionally, the parallel approach showed excellent overall speedup: a 60-fold
50 improvement using 4-GPU relative to the serial run. A regional scale application further
51 demonstrated the LB performance. The parallel time used by 8-GPU without LB was 31.33%
52 reduced after LB was activated. When increasing 8-GPU with LB to 16-GPU with LB, it showed
53 parallel scalability by reducing the parallel time of ~50%. This work shows how massively
54 parallel computing can be applied to particle tracking in water age simulations. It also
55 demonstrates the practical importance of load balancing in this context, which enables the large-
56 scale simulations with an increased complexity of flow paths.

57

58

59

60 **Keywords**

61 Water age, Particle tracking, Multi-GPU with MPI, Domain decomposition, Load balancing

62 **1. Introduction**

63 Water age is an important metric that can unravel the journey that water-parcels and
64 pollutants take while traveling through a watershed (Botter, Bertuzzo, & Rinaldo, 2011).
65 Methods used to quantify water ages mainly include tracer data, lumped analytical models, and
66 distributed numerical models (Nicholas B. Engdahl & Maxwell, 2014). Tracer data are direct
67 observations of system behavior, however they have technical challenges such as tracer selection
68 and data interpretation (Sprenger et al., 2019). Moreover, the limited sampling cannot provide a
69 full view of the spatiotemporal variations of water ages (Nicholas B. Engdahl, McCallum, &
70 Massoudieh, 2016; McCallum, Engdahl, Ginn, & Cook, 2014). Analytical solutions can also
71 provide an understanding of the real system, although the limitations of this approach due to
72 simplifications have been well acknowledged by the community (Basu, Jindal, Schilling, Wolter,
73 & Takle, 2012). For distributed numerical models, the application of Eulerian framework is
74 hampered by the complications to solve the high-dimensional governing equation of water age
75 (Gomez & Wilson, 2013). Therefore, Lagrangian approach based on integrated hydrologic
76 modeling has become a promising tool to simulate transient age distributions (Nicholas B.
77 Engdahl & Maxwell, 2014; Jing et al., 2019; Wilusz, Harman, Ball, Maxwell, & Buda, 2019).

78 The Lagrangian approach is computationally demanding which limits widespread application
79 (Sprenger et al., 2019; Wilusz et al., 2019). Current studies of water ages using particle tracking
80 are limited to either catchments at small scales (Wilusz et al., 2019; J. Yang, Heidbüchel,
81 Musolff, Reinstorf, & Fleckenstein, 2018) or larger scales with a limited number of particles or
82 for steady-state conditions (Jing et al., 2020; Maxwell et al., 2016). Recently, as global water
83 security and climate change become increasing concerns, a growing number of studies have
84 proposed simulating water age at larger scales over long time-periods at high resolution
85 (Maxwell et al., 2016; McGuire et al., 2005; Starn, Kauffman, Carlson, Reddy, & Fienen, 2021).
86 Accomplishing these goals will significantly increase the computational burden of particle
87 tracking; massively parallel computing represents a promising solution.

88 Currently, there are few studies documenting parallelization of particle tracking approaches
89 for water age (Jing et al., 2020; Wilusz et al., 2019; J. Yang et al., 2018). Maxwell, Condon,
90 Danesh-Yazdi, and Bearup (2019) developed EcoSLIM, a particle tracking code, using OpenMP
91 (Open Multi-Processing) on CPU (Central Processing Unit). Yang and coauthors (C. Yang et al.,
92 2021), added CPU-based MPI (Message Passing Interface) and multi-GPU (Graphics Processing

93 Unit) with OpenMP into EcoSLIM. Ji, Luo, and Wang (2019) sped up MODPATH (Pollock,
94 2016) through multi-GPU with OpenMP/MPI based on domain decomposition (DDC). However,
95 their parallelized codes are limited to steady state simulations, and MODPATH is unable to
96 simulate evapotranspiration (ET) age and source water composition. N. B. Engdahl, Schmidt,
97 and Benson (2019) proposed two schemes for speeding up particle tracking simulations with
98 mass transfer to represent chemical reactions. One of the schemes also implemented MPI
99 parallelism through DDC. In DDC-based MPI parallelism, a typical problem encountered is load
100 imbalance among processing elements (MPI processes and/or GPUs), which can present a
101 challenge for good parallel efficiency. However, load imbalance has not been quantified in water
102 age simulations and its effects on parallel performance are unclear.

103 Load balancing (LB) schemes have been used to overcome these parallelization challenges
104 discussed above. Other disciplines using particle tracking, such as the molecular dynamics (MD)
105 and the smoothed particle hydrodynamics (SPH) (Boulmier, Raynaud, Abdennadher, & Chopard,
106 2019; Egorova, Dyachkov, Parshikov, & Zhakhovsky, 2019; Eibl & Rde, 2019; Fattebert,
107 Richards, & Glosli, 2012; Furuichi & Nishiura, 2017; Kunaseth et al., 2013) have presented LB
108 schemes that greatly improved parallel simulation performance. However, LB has not been
109 applied to hydrologic modeling based on particle tracking; even in the studies using MPI through
110 DDC mentioned above. Additionally, when applying the particle tracking at larger scales in
111 long-term simulations, spatiotemporal variations of water age drivers in the real-world
112 applications can increase the heterogeneity of flow paths. This heterogeneity causes uneven
113 particle distributions and velocities within the domain, presenting new challenges to efficiency.
114 This will further complicate the distribution of particles across different subdomains and thus the
115 processing elements (which we are using here as a generic term for compute resources such as
116 CPU cores or a GPU). Furthermore, it becomes challenging to implement LB when considering
117 the increasing complexity of code structure of particle tracking due to the growing capabilities
118 such as simulating ET age and source-water composition at transient state in EcoSLIM.

119 In this study, we present new LB approaches implemented in the EcoSLIM code which is a
120 particle tracking code simulating water age (ET, outflow, and groundwater) and source water
121 mixing (initial subsurface water, rainfall, and snow). EcoSLIM is a grid-based approach which is
122 different from the mesh-free particle tracking in other disciplines. EcoSLIM works seamlessly
123 with ParFlow.CLM (Kollet & Maxwell, 2008a), which is an integrated hydrologic model

124 simulating the coupled land-surface and subsurface water- and energy-processes at transient state.
125 ParFlow.CLM provides the temporally variant hydrodynamics and spatially variable subsurface-
126 properties for EcoSLIM as input files, such as saturation, precipitation minus ET, three-
127 dimensional velocity fields, and IDs and porosities of the subsurface units.

128 Objectives of this study are twofold. Firstly, we use MPI to manage multi-GPU instead of
129 OpenMP in our previous work (C. Yang et al., 2021). This changes the target of the
130 decomposition from computational load to modeling domain. MPI parallelism removes the
131 barrier of a limited number of GPUs on a single computational node by created when using
132 OpenMP and potentially extends the particle tracking applications to massively parallel
133 computing. Secondly, three schemes with increasing physical representation are proposed to
134 dynamically balance the load among different MPI-processes/GPUs during runtime, which is
135 crucial for parallel efficiency. In following sections, the EcoSLIM code, the implementation of
136 multi-GPU with MPI, and the load balancing schemes are introduced in section 2. The setup of
137 test cases and platforms are followed in section 3. In section 4, validation of the new code is
138 verified and parallel performances of the code with/without the LB schemes are illustrated.
139 Specifically, application of the code for a 40-year simulation at regional scale in the North China
140 Plain was shown. Finally, contributions and implications of this work to hydrologic modeling
141 using Lagrangian approach are concluded in section 5.

142 **2. Methodology**

143 **2.1. EcoSLIM code**

144 EcoSLIM is originally implemented in Fortran and further accelerated by GPUs using CUDA
145 (Compute Unified Device Architecture) Fortran (C. Yang et al., 2021). Its original structure is
146 briefly introduced here to understand the following implementations of multi-GPU with MPI and
147 LB schemes. For more details, please refer to our previous work (Maxwell et al., 2019; C. Yang
148 et al., 2021). In each timestep of a transient simulation, key steps are as follows:

- 149 (1) New particles are added into grid-cells where precipitation minus ET (PME) is positive.
- 150 (2) Mass balance of precipitation and ET is calculated based on PME.
- 151 (3) Advancing each active particle by a do-loop. Each particle either moves forward with an
152 increase of age or exits the modeling domain through outflow or ET.

153 (4) After the particle loop, inactive particles that have left the modeling domain via outflow
154 or ET are sorted out of the array of particles, the number of active particles is updated and
155 space at the end of this array is made available for new particles.

156 (5) The time loop moves to next timestep.

157 In the particle loop, other attributes of each particle are also updated, such as the
158 saturated/unsaturated travel time, the saturated/unsaturated travel length, and the travel
159 time/length in some specified subsurface units. Statistics of outflow and ET of the whole
160 modeling domain, such as mass, mass weighted age, and source water composition, are included.
161 Additionally, gridding information is recorded, such as mass, cell-averaged age, and source
162 water composition for each grid-cell. Particle loop is the main computational load in EcoSLIM
163 which occupies more than 99% of the total simulation time when serially executing the code (C.
164 Yang et al., 2021). Therefore, our previous work (C. Yang et al., 2021) and also this study
165 focused on parallelizing the particle loop.

166 **2.2. Multi-GPU with MPI**

167 Instead of the load decomposition in our previous work (C. Yang et al., 2021), MPI
168 parallelism is applied here based on DDC. The modeling domain is split into P and Q parts in x
169 and y directions respectively, which generates a computing topology using $P \times Q$ MPI processes.
170 The quantity of GPUs utilized in a simulation equals that of MPI processes. We use the method
171 in Ruetsch and Fatica (2014) to assign a unique GPU to each MPI rank. GPUs, MPI ranks, and
172 subdomains are all numbered from 0 to $P \times Q - 1$. Each GPU is responsible for a subdomain of
173 the same number. Source of particles on each GPU is from the assigned subdomain while
174 transport of these particles is in the whole modeling domain. For the particle loop, we adopt the
175 same GPU kernel in our previous work (C. Yang et al., 2021) with a few modifications for
176 tracking particles in specified subsurface units. After the mapping to subdomains, each MPI-
177 process/GPU works almost independently with a limited MPI collective communications. They
178 are the calculation of mass balance mentioned in item (2) in section 2.1 and the ET/outflow
179 statistics of the whole modeling domain also mentioned in section 2.1, which are performed
180 before and after the execution of the kernel respectively.

181 The transport of particles in the whole domain, instead of in the subdomain where they were
182 added, avoids the MPI communications of exchanging particles between subdomains when
183 particles move out of a subdomain. The disadvantage of such a design is the redundant copies of

184 the global information (e.g., velocities, porosities, saturations, and PME) for all MPI processes.
185 This is CPU-memory expensive. However, for present clusters which are commonly equipped
186 with 2/4/8 GPUs per node, the 2/4/8 copies of necessary information on one node are not a
187 bottleneck for regional modeling with scales of Tran, Zhang, Cohard, Condon, and Maxwell
188 (2020) and C. Yang et al. (2020). The multi-GPU with MPI was also conducted in our previous
189 work based on load decomposition (C. Yang et al., 2021). However, in each timestep, the
190 overhead of distributing load from rank 0 to others by MPI communication is over the speedup
191 by extending single GPU to multi-GPU. Though the DDC in this study avoids such a problem,
192 the load imbalance mentioned in section 1 becomes a new issue. Therefore, three schemes with
193 increasing physical representation are proposed to balance the load among GPUs/MPI-processes
194 during runtime in next section.

195 **2.3. Schemes of load balancing**

196 **2.3.1 Direct transfer (S1)**

197 The movement of particles in current EcoSLIM are independent, so the loads on GPUs can
198 be redistributed by directly transferring particles between MPI processes with a user specified
199 frequency. It is implemented by communications on CPU using the MPI functions of *MPI_Send*
200 and *MPI_Recv* after the sort of particles. At a given time, the number of active particles on each
201 MPI process is gathered. Thus, the old numbers of the starting and ending particles (np_{lo} and
202 np_{ro}) on each process ranked in a global queue are obtained. Then the starting and ending
203 numbers are updated to new ones (np_{ln} and np_{rn}) based on an even division of the global
204 queue. Then the transfer is accomplished on each MPI process by the following four steps: (1)
205 sending particles to the upstream process if np_{lo} is smaller than np_{ln} , (2) receiving particles
206 from the downstream process if np_{ro} is smaller than np_{rn} , (3) sending particles to the
207 downstream process if np_{ro} is larger than np_{rn} , and (4) receiving particles from the upstream
208 process if np_{lo} is larger than np_{ln} . The number of active particles is updated after each transfer.
209 The overhead of this scheme is determined by the quantity of particles transferred and the
210 bandwidth of the computing platform.

211 **2.3.2 Cyclic mapping (S2)**

212 In this scheme, the DDC is static which is determined after initialization of the simulation. In
213 a simulation using n GPUs/MPI-processes, the mapping between subdomains and GPUs (both
214 numbered from 0 to $n-1$) is continuously shifted with a user specified frequency. For instance, at

215 a given time t_1 , the mapping is shifted from ‘the m th subdomain \rightarrow the m th GPU’ to ‘the $(m+1)$ th
 216 subdomain \rightarrow the m th GPU’ where m ranges from 0 to $n-1$. If $m+1$ is larger than $n-1$, the
 217 subdomain numbered with the remainder of $m+1$ and n will be mapped to the m th GPU. Table 1
 218 showed the cyclic mapping between subdomains and GPUs in a simulation using four MPI
 219 processes. Using this scheme, each GPU traverses the loads of all subdomains periodically, and
 220 thus the load distribution is dynamically balanced among GPUs. The overhead of shifting the
 221 mapping is almost negligible.

222 **Table 1. Cyclic mapping in a simulation using four GPUs/MPI-processes**

t_0 (Initial)	Subdomain number	0	1	2	3
	GPU number	0	1	2	3
t_1	Subdomain number	1	2	3	0
	GPU number	0	1	2	3
t_2	Subdomain number	2	3	0	1
	GPU number	0	1	2	3

223
 224 **2.3.3 Dynamic DDC (S3)**
 225 In the initialization of a simulation, particles are evenly distributed in space, so we
 226 decompose the modeling domain into subdomains of an equal size. PME is spatiotemporally
 227 variable, so the number of particles added into each subdomain is different at a given timestep
 228 and such a difference varies with time. More importantly, particles added into different
 229 subdomains have different exits from the whole modeling domain. Both source and exit are
 230 responsible for heterogeneity of the flow paths. A subdomain of more source particles and longer
 231 flow paths imposes heavier load on its corresponding GPU. As a result, after the initial even-
 232 decomposition, we dynamically update the decomposition during runtime based on flow paths of
 233 particles. At a given timestep, the initial location of an active particle is identified and the
 234 corresponding grid-cell in the top layer get one score. After traversing all the active particles, we
 235 get the accumulated scores for each grid-cell in the top layer. This was implemented through
 236 atomic operations on a two-dimensional matrix in the GPU kernel mentioned in section 2.2.
 237 Then we conduct DDC based on this weight matrix. The frequency of such a dynamic DDC can
 238 be specified by users.

239 The orthogonal recursive bisection (ORB) method is used for DDC, which is popular in MD
 240 and SPH (Egorova et al., 2019; Fattebert et al., 2012). The domain or each subdomain is divided
 241 into two in a direction at one time. By switching the direction, the whole domain is recursively

242 divided into the scheduled number of subdomains ($P \times Q$). DDC is only implemented in x and y
243 directions in this study. It starts in a direction which will be divided into more pieces. For
244 example, if Q is larger than P , DDC will start in y direction, otherwise, it starts in x direction.
245 Two algorithms are provided in the code to determine the dividing line. One calculates the
246 accumulated particles of columns (rows) in x (y) direction. Once the number of particles of n -1
247 columns (rows) is less than half the total particles in this subdomain while that of n columns
248 (rows) is more than half the total particles in this subdomain, the dividing line is found as column
249 (row) n . The other algorithm to find the dividing line with higher efficiency is the typical
250 dichotomizing search.

251 **3. Test setup**

252 **3.1. Hillslope model**

253 Tests were conducted based on a hillslope model (Maxwell et al., 2019; C. Yang et al., 2021).
254 The modeling domain has the length of 100-, 1-, and 9.4-m in x , y , and z directions, respectively.
255 It was divided into 20 columns, 5 rows, and 20 layers with constant resolutions in x and y
256 directions. In vertical direction, the layer-thickness was variable: 0.5 m for the bottom 18 layers
257 while 0.3- and 0.1-m for the top 2 layers. Soil has the homogeneous properties: saturated
258 hydraulic conductivity of 0.05 m/h, Manning's N of $10^{-6} \text{ m}^{1/3}\text{h}^{-1}$, porosity of 0.2, and van
259 Genuchten parameters with α of 1.0 m^{-1} and exponent n of 2.0. Two real meteorological-forcings
260 were used to drive ParFlow.CLM, representing a high elevation, snow dominated mountain
261 headwaters (ER) and a semiarid, rain-dominated plains system (LW). Two homogeneous land-
262 cover types were used which are the Shrub plant functional type (Shrub) and the Evergreen
263 Needleleaf plant functional type (Trees). Thus, four cases were tested with a combination of the
264 meteorological forcings and the land-cover types, which were named as: ER_Shrub, ER_Trees,
265 LW_Shrub, and LW_Trees. For simulations of both ParFlow.CLM and EcoSLIM, no flux
266 boundaries were adopted except the land surface which was open for precipitation, outflow and
267 ET. For each case, ParFlow.CLM simulation of 5 years was conducted using hourly timestep.
268 One-year forcing data were repeatedly used in the whole simulation. Dynamic equilibrium of the
269 flow field was approached at the end of simulation. Hence, the transient flow field of the last
270 year in ParFlow.CLM simulation was repeatedly used in EcoSLIM simulation of 20 years with
271 hourly timestep. At the end of each simulation, EcoSLIM system achieved the dynamic
272 equilibrium.

273 By injecting 2 particles into the modeling domain per precipitation event, the average
274 particle-numbers in the last year of the simulations for four cases are 0.39-, 0.83-, 1.30-, and
275 1.03-million, respectively. Such quantities of particles are comparable to those in most of the
276 previous studies (Danesh-Yazdi, Klaus, Condon, & Maxwell, 2018; Nicholas B. Engdahl &
277 Maxwell, 2015; Jing et al., 2019; Jing et al., 2020; Kollet & Maxwell, 2008b; Maxwell et al.,
278 2016; Weill, Lesparre, Jeannot, & Delay, 2019; Wilusz et al., 2019). It has a maximum of 6.8
279 million in Weill et al. (2019) to the best of our knowledge. In fact, the particle-number in most
280 previous studies is the total injected particles while that during runtime is a fewer quantity.
281 However, the number in this study is the active particles in the modeling domain and the
282 particles out of the domain through ET and outflow are not included. Thus, the particle-
283 quantities are much more than those in previous studies.

284 **3.2. Test platform**

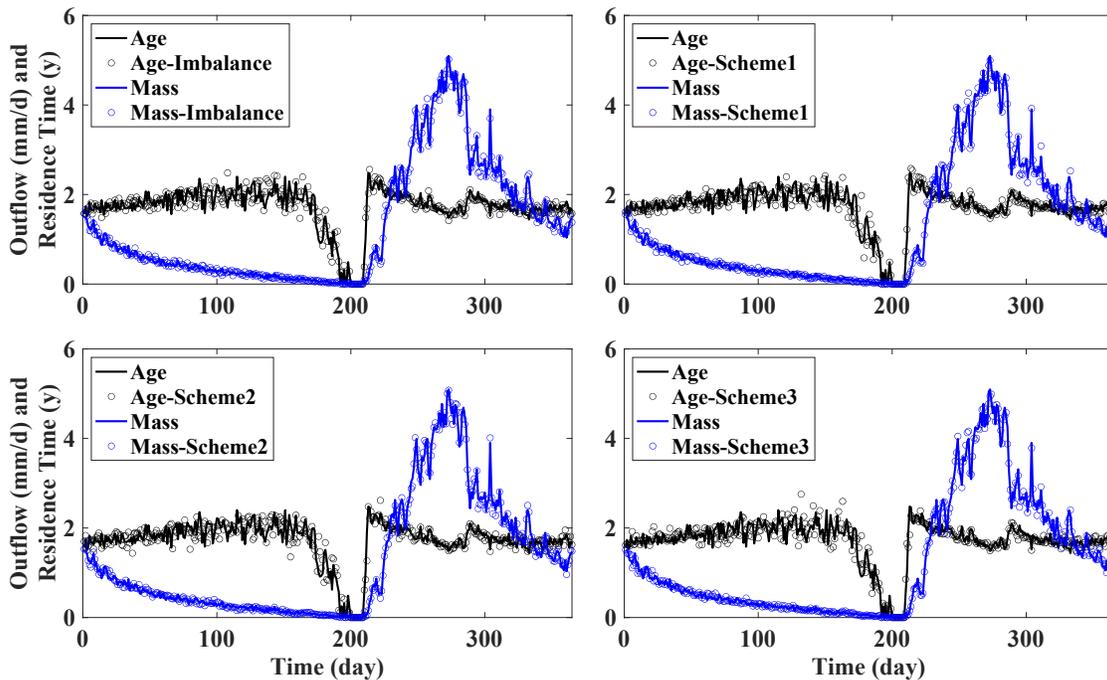
285 Tests were conducted on the Casper cluster in the Computational and Information Systems
286 Laboratory at the National Center for Atmospheric Research. The computational node used for
287 the following simulations is equipped with 2.3-GHz Intel[®] Xeon[®] Gold 6140 processors and
288 NVIDIA Tesla V100 32GB SXM2 GPUs with NVLink. The compiler is NVIDIA HPC SDK of
289 version 20.11, the MPI is implemented using Open MPI of version 4.0.5, the GPU driver version
290 is 450.51.06, and the CUDA version is 11.0.3. Tests were also repeated on a personal
291 workstation (WS). The WS is equipped with 2.00-GHz Intel[®] Xeon[®] E5-2683 v3 processors
292 together with four GPUs of 12 GB GeForce GTX 1080 Ti. Other necessary setups of the WS
293 environment are NVIDIA HPC SDK 20.11, Open MPI 3.1.5, GPU driver 440.118.02, and
294 CUDA 10.2.

295 **4. Results and discussion**

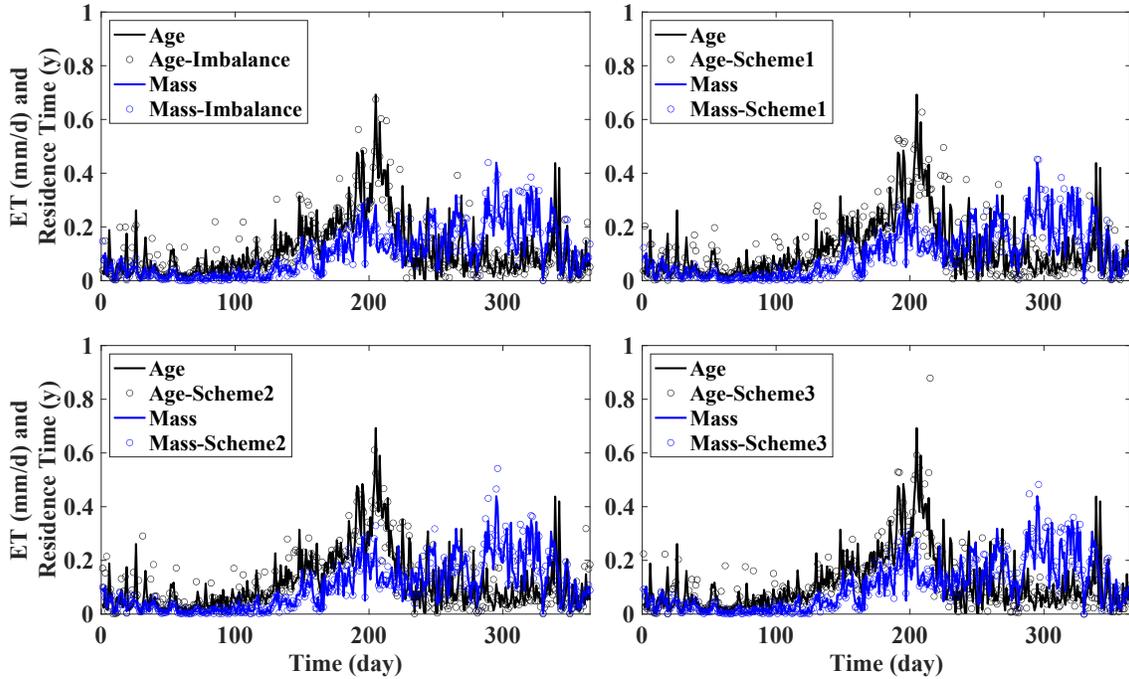
296 **4.1. Code-to-code verification**

297 To verify the availability of the new code, simulation results using the new code were
298 compared to those of the original OpenMP version (Maxwell et al., 2019). Comparisons were
299 performed for all four cases while that of the ER_Shrub case was shown in Figures 1 (outflow)
300 and 2 (ET). Results of other test cases had performances as good as that of ER_Shrub. Tests in
301 this study were conducted using one, two, and four GPUs successively while the results using
302 four GPUs were illustrated in Figures 1 and 2. Subplots in Figures 1 and 2 were for results
303 without LB and with each LB scheme. The water-age and -mass for both outflow and ET

304 simulated by the new code well fitted those generated by the original code. The deviations
 305 between them were attributed to the generation of pseudo-random numbers (PRNs). Though the
 306 ensembles of the PRNs were statistically the same for each run, the PRN for a specific particle
 307 probably changed due to the invoking sequence of the generation-function which was dependent
 308 on the parallelism, i.e., the OpenMP or the multi-GPU with MPI. For the same parallelism, if
 309 different numbers of CPU-threads/GPUs were used, there were also such deviations during our
 310 tests. The fitness of outflow was better than that of ET because ET in EcoSLIM were directly
 311 dependent on PRNs.

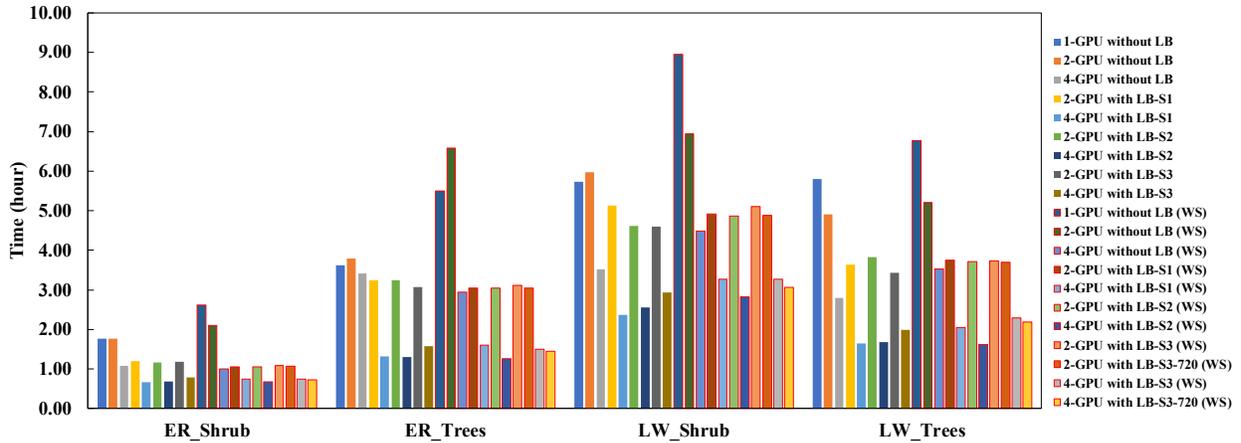


312
 313 **Figure 1. Comparisons for age and mass of outflow based on ER_Shruh case between the original**
 314 **EcoSLIM code and that parallelized in this study.**
 315



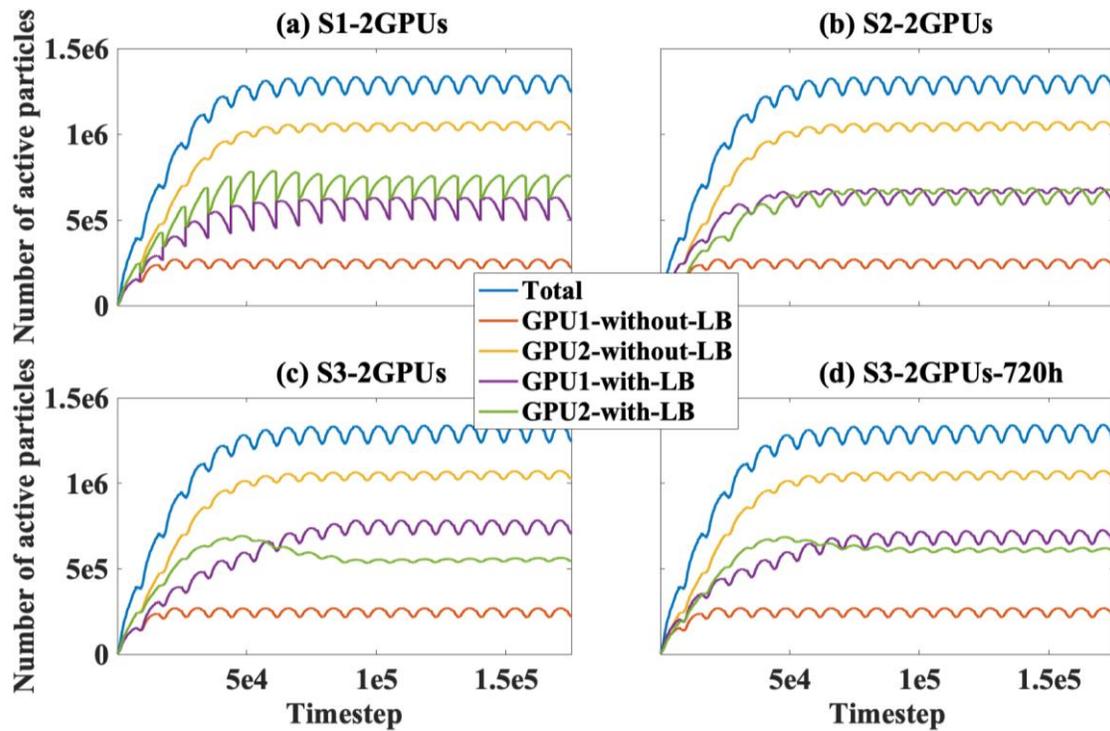
316
 317 **Figure 2. Comparisons for age and mass of ET based on ER_shrub case between the original**
 318 **EcoSLIM code and that parallelized in this study.**

319
 320 **4.2. Parallel performance**

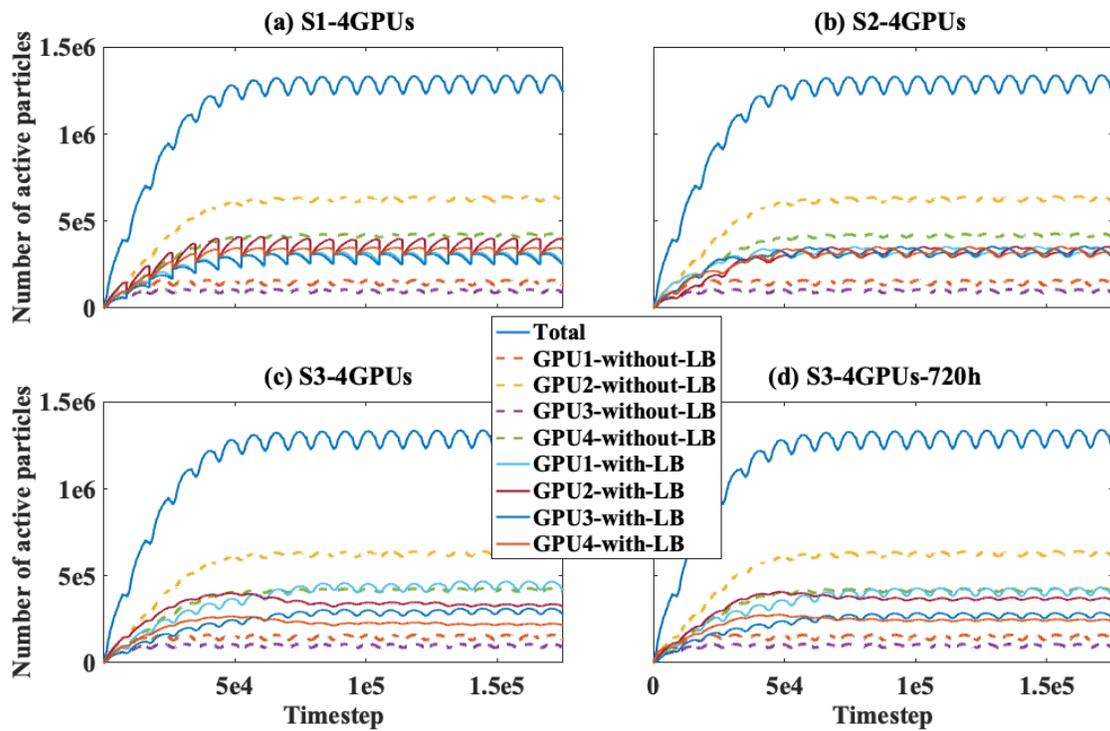


321
 322 **Figure 3. Wall-clock time consumption of each test. n -GPU represents the number of GPUs used in**
 323 **simulations. Sn represents different LB schemes. WS indicates tests conducted on workstation**
 324 **while others on Casper. 720 indicates S3 worked every 720-hour while others worked every 8760-**
 325 **hour.**

326
 327



328
 329 **Figure 4.** Load distribution for LW_Shrub by using 2-GPU. S_n represents LB schemes. 720h
 330 represents that S3 worked every 720-hour while others were 8760-hour.
 331



332
 333 **Figure 5.** Load distribution for LW_Shrub by using 4-GPU. S_n represents LB schemes. 720h
 334 represents that S3 worked every 720-hour while others were 8760-hour.
 335
 336

337 Figure 3 shows wall-clock time consumption of each test on both Casper and WS. Load
338 distributions for tests using 2- and 4-GPU were shown in Figures 4 and 5, respectively. Tests
339 using 2-GPU ($P = 2$ and $Q = 1$) were performed by evenly dividing the domain in x direction
340 while those using 4-GPU ($P = 2$ and $Q = 2$) had an additional division in y direction. All LB
341 schemes worked every 8760-hour except that in Figures 4d and 5d were 720-hour for S3.
342 Though only LW_Shruh was taken as an example in Figures 4 and 5, other cases had similar
343 performances. Time used for four cases by one CPU-thread based on the original code were
344 tested on WS, which were 36.41-, 76.34-, 121.56-, and 92.40-hour for ER_Shruh, ER_Trees,
345 LW_Shruh, and LW_Trees, respectively. Speedup of each case was then calculated and listed in
346 Table 2.

347 **Table 2. Speedup of each test relative to the serial run using one CPU-thread**

Platform	Case name	Speedup								
		Without LB			Scheme 1		Scheme 2		Scheme 3	
		1-GPU	2-GPU	4-GPU	2-GPU	4-GPU	2-GPU	4-GPU	2-GPU	4-GPU
Casper	ER_Shruh	20.5403	20.5693	33.5134	30.5237	54.8938	31.0469	52.7465	30.6935	46.3058
	ER_Trees	21.1316	20.1687	22.3125	23.5297	57.5619	23.5857	58.7581	24.8849	48.6518
	LW_Shruh	21.1858	20.3505	34.5421	23.7138	51.5341	26.3410	47.6590	26.4183	41.5489
	LW_Trees	15.9127	18.8541	32.9764	25.3874	56.0544	24.1910	54.8824	26.9184	46.2856
WS	ER_Shruh	13.9637	17.2900	36.0809	34.4456	48.8535	34.4522	54.4261	33.4696	49.0443
	ER_Trees	13.8955	11.5971	26.0117	25.0581	47.8438	25.1381	60.4948	24.5480	51.0523
	LW_Shruh	13.5910	17.5129	27.0710	24.7589	37.1392	24.9531	43.0452	23.8150	37.1698
	LW_Trees	13.6440	17.7556	26.1808	24.6492	45.1304	24.8835	57.2562	24.7449	40.4952

348 In tests without LB, time used by 2-GPU was even more than that by 1-GPU for ER_Trees
349 and LW_Shruh on Casper (Figure 3). This performance degradation has two reasons. Firstly,
350 severe load imbalance can be observed in Figure 4 when using 2-GPU with a particle-number
351 ratio of 423.62%. The larger one (yellow lines in Figure 4) which determined the parallel
352 efficiency almost achieved the total load. It confirmed that load imbalance also exists in
353 Lagrangian hydrologic modeling and decreases the parallel performance. Secondly, collective
354 MPI communications mentioned in section 2.2 introduced overhead when increasing the GPU
355 number from one to two. When using 4-GPU without LB, the maximum load largely decreased
356 (Figure 5) due to the additional decomposition in y direction. The hillslope model is quasi-three-
357 dimensional with a x slope of 0.1 and a y slope of 0. Hence, the movement of particles in y
358 direction can be neglected which formed the parallel flow paths along x direction. Along x
359 direction, particles added upstream had much longer flow paths than those added downstream.
360 As a result, the division in y direction was much more effective than that in x direction to

361 improve the parallel performance. However, load imbalance was still significant (Figure 5) and
362 parallel scalability was not shown by increasing 2-GPU to 4-GPU (Figure 3).

363 When S1 was activated, time used by 2- and 4-GPU were dramatically decreased relative to
364 those without LB. With S1, time used by 4-GPU was even less than half of that used by 2-GPU
365 on Casper (Figure 3). The overhead of particle transfer was small for all four cases, which was
366 less than 3 seconds in each 20-year simulation. S2 had performance as good as S1 (Figure 3). It
367 was mentioned in section 2.3.2 that the overhead of S2 was small enough to be neglected. For S3,
368 more time was used by 4-GPU when compared to that of S1 and S2 (Figure 3). In Figures 4c and
369 5c, the load distribution was not well balanced relative to that of S1 and S2. The flow paths of
370 particles were transient, so the weight matrix at a moment cannot effectively balance the load for
371 a long-period of simulation (i.e., 8760-hour). When S3 was activated with a higher frequency of
372 every 720-hour, the improved load balance can be observed in Figures 4d and 5d and the further
373 speedup was indicated in Figure 3. However, difference of the loads between GPUs 1-2 and 3-4
374 was still observed in Figure 5d. This is due to the model dimension which is five grid-cells in y
375 direction. Hence it cannot be evenly divided by ORB introduced in section 2.3.3.

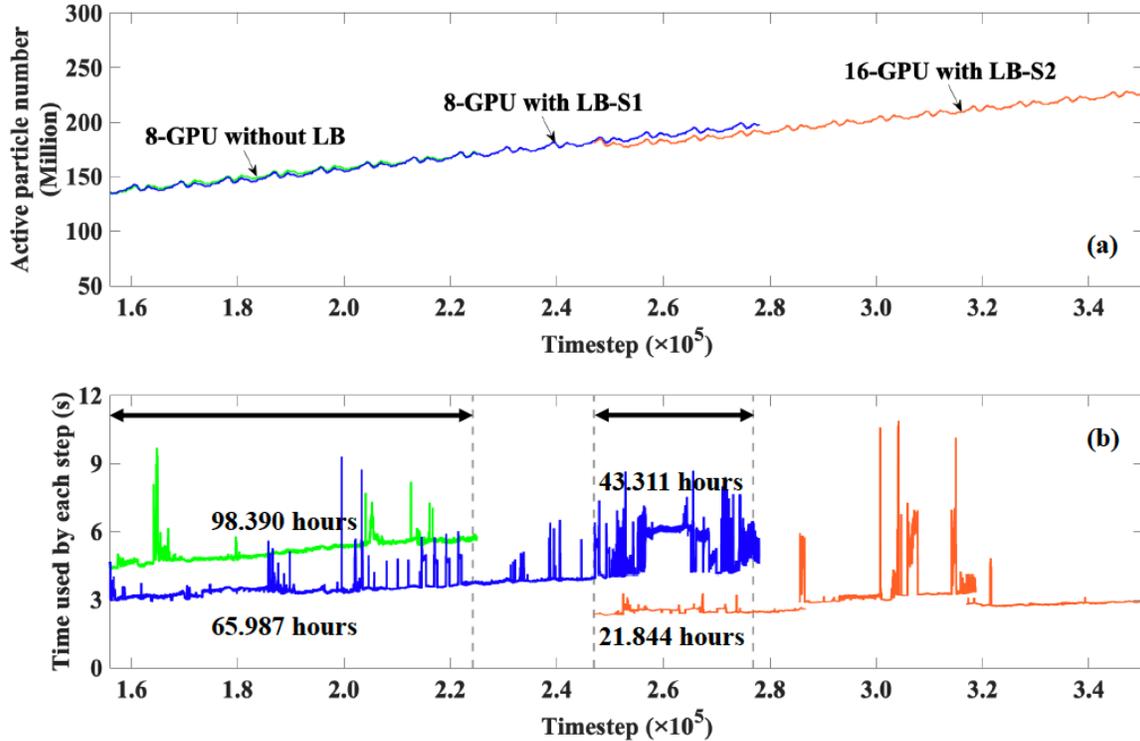
376 S3, a physically-based scheme, not only aims to balance the load but also to understand the
377 load imbalance in Lagrangian hydrological modeling. To build S3, we also tried DDC based on
378 the source of particles (i.e., PME). The score/weight of a grid-cell is determined by the
379 accumulated particle-number added into it in a period. However, it didn't show good
380 performance. Current implementation actually integrates the effects of both the quantity of
381 source particles and the flow-path lengths. This trial and error indicates that the flow-path
382 lengths instead of the quantity of added particles dominate the load distribution. This has
383 important implications to efficiently build other physically-based LB schemes. Generally, with
384 LB, the new code showed excellent parallel performance in the tests on both Casper and WS
385 (Table 2). The speedup by one GPU is ~13-fold on WS with 1080 Ti while ~21-fold on Casper
386 with Tesla V100. The speedup by 4-GPU is over 50-fold on both Casper and WS and has a
387 maximum over 60-fold. With LB schemes, the code showed parallel scalability from 2-GPU to
388 4-GPU.

389 **4.3. Application in the North China Plain**

390 We applied the new parallel code on a North China Plain (NCP) domain to demonstrate its
391 capacity for large-scale simulations. To the best of our knowledge, there have been no previous

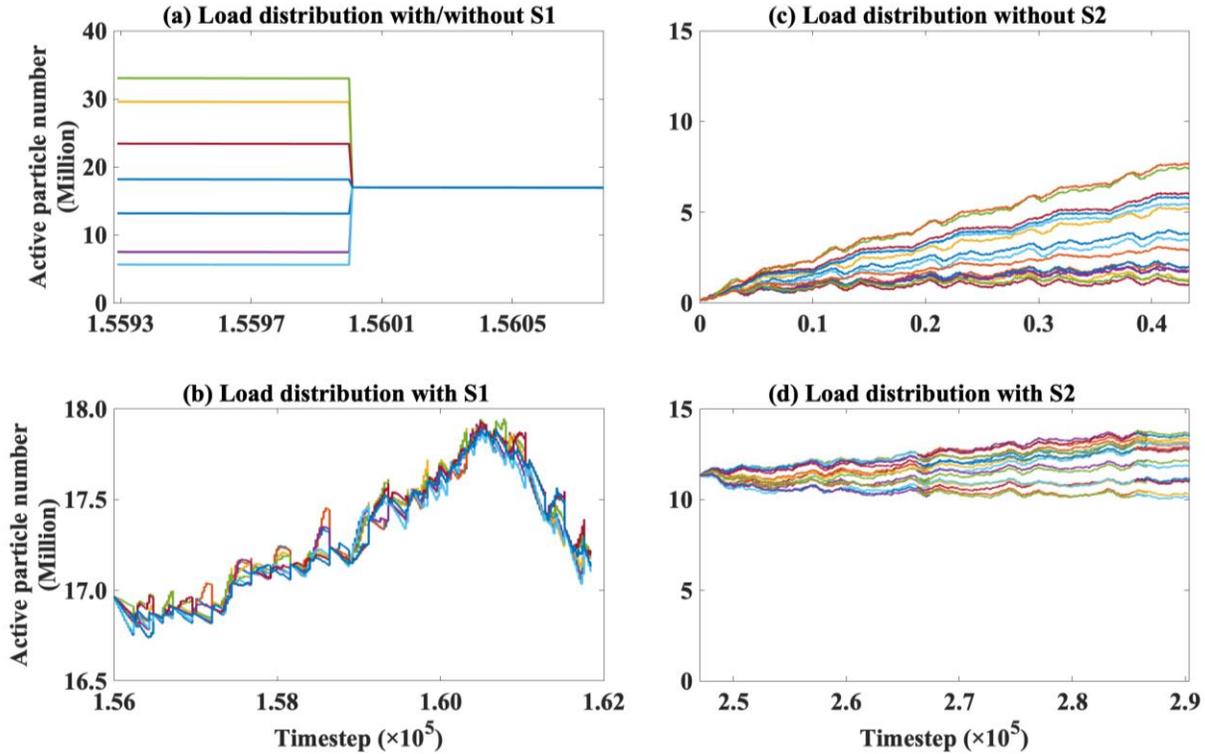
392 studies based on particle tracking at such a regional scale for water ages of ET, groundwater
393 (GW), and outflow in a unified framework. The NCP ParFlow.CLM model was adopted from C.
394 Yang et al. (2020) with a few modifications. The model has 509 and 921 grid-cells in x and y
395 directions respectively while it is discretized into five layers in vertical direction. The horizontal
396 resolution is 1 km while the layer thickness from bottom to top is 100-, 1-, 0.6-, 0.3-, and 0.1-m.
397 Thus, the NCP model has a dimension of $509 \text{ km} \times 921 \text{ km} \times 102 \text{ m}$ in total. We conducted an
398 EcoSLIM simulation of 40 years on Casper with hourly timestep, in which the hourly outputs of
399 one-year simulation from ParFlow.CLM were repeatedly used.

400 The EcoSLIM simulation was started using 8-GPU ($P = 2$ and $Q = 4$) without LB
401 (abbreviated as R1 hereafter). From the 155,928th hour, S1 was activated every 240-hour (R2)
402 while R1 was continued for the following 7.8 years. At the 247,032th hour, the load of R2 was
403 evenly divided into 16 portions and a new run (R3) was started using 16-GPU ($P = 4$ and $Q = 4$)
404 with S2 activated every 240-hour. The overlap between R2 and R3 is 3.4 years. We also tried 16-
405 GPU without LB for the first five years of the simulation (R4). Figure 6 showed the active-
406 particle-number and the wall-clock time consumption of each timestep during the latter 22 years
407 of the simulation (the 155,929th to the 350,400th hour). R1, R2, and R3 were indicated by green,
408 blue, and red in Figure 6 respectively. The active-particle-number is around 200 million during
409 this simulation time-interval (Figure 6a). The discrepancy of the particle number between
410 different runs in the overlaps are due to the generation of random numbers discussed in section
411 4.1.



412
 413 **Figure 6. Computational load (a) and wall-clock time consumption (b) for the EcoSLIM simulation**
 414 **in the North China Plain. The time interval is from the 155929th to the 350400th hour in the 40-**
 415 **year simulation.**

416 For the overlap between R1 and R2, parallel time of the particle loop was 98.390- and
 417 65.987-hour for R1 and R2 respectively (Figure 6b). The overhead of S1 for transferring data
 418 was 1.573-hour. S1 (overhead included) decreased 31.33% of the time used by R1, which
 419 demonstrated the high efficiency of S1. Figures 7a and 7b showed the well-balanced load by S1.
 420 For the overlap between R2 and R3, the parallel time was 43.311- and 21.844-hour for R2 and
 421 R3 respectively (Figure 6b). Though it showed 50% decrease of the parallel time, the obvious
 422 jitters of the time in R2 has to be considered. Based on their baselines, the time used by R3 was a
 423 little longer than half the time used by R2. This should be due to the better load balancing effect
 424 of S1 than that of S2, which was shown in Figures 7b and 7d. However, when comparing the
 425 load distribution between R3 and R4 for a time interval of the same length (4.94-year), the load
 426 balancing effect of S2 was significant. The difference between the maximum- and minimum-
 427 load at the end of the comparing time-interval in R4 was 6.66 million (Figure 7d) while that in
 428 R3 was 3.52 million (Figure 7c), which was 47.21% decrease of the load variance. Additionally,
 429 based on the increasing trend in Figures 7c and 7d, the load variance in R3 with S2 gradually
 430 achieved a steady state while that in R4 continued increasing.



431
 432 **Figure 7. Load distributions in the application in the North China Plain. Load distribution on 8-**
 433 **GPU before and after using S1 in R2 (a), on 8-GPU with S1 in R2 (b), on 16-GPU without LB in R4**
 434 **(c), and on 16-GPU with S2 in R3. (b) was magnified from (a).**

435
 436 **5. Conclusions**

437 Water age can reveal the source, storage, and mixing of water parcels in a watershed. Though
 438 data- and model-driven methods have significantly advanced our understanding of water ages,
 439 the quantification of water ages is still technically challenging. Lagrangian particle tracking is an
 440 invaluable tool for physically-based transient modeling of water ages, but it is computationally
 441 expensive. When considering climate change and global water security, it is essential to conduct
 442 simulations of water ages at large scale with high resolution, which makes the implementation of
 443 massively parallel computing in particle tracking for this purpose pressing. Though parallel
 444 computing is widely implemented for Eulerian hydrological modeling, applications to
 445 Lagrangian based simulations are developing. This is likely due to the inherent difficulties such
 446 as load imbalance across computational resources which will become more challenging when
 447 modeling a real hydrologic system with high spatiotemporal variability.

448 In this study, multi-GPU with MPI parallelism based on domain decomposition (DDC) was
 449 implemented in the Lagrangian, particle tracking code EcoSLIM, to accelerate simulations of

450 water age and source-water mixing. Three load balancing (LB) schemes with increasing physical
451 representation (i.e., direct transfer, cyclic mapping, and dynamic DDC) were built to
452 dynamically balance the quantity of particles across GPUs during runtime. With LB, the code
453 showed excellent parallel performance in the hillslope simulations on two different platforms,
454 e.g., a maximum of 60-fold speedup on 4-GPUs and the parallel scalability from 2-GPU to 4-
455 GPU that is almost ideal. A 40-year simulation conducted in the North China Plain further
456 demonstrated the high parallel efficiency of LB for a large-scale application. Using 8-GPU with
457 LB, it reduced 31.33% of the parallel time using 8-GPU without LB. When increasing 8-GPU
458 with LB to 16-GPU with LB, ~50% reduction of the parallel time demonstrated the parallel
459 scalability.

460 More importantly, results confirmed the load imbalance in Lagrangian hydrologic modeling.
461 In LW_Shrub case using 2-GPU, the particle-number ratio achieved 423.62%, which severely
462 degraded the parallel performance without LB. For LB schemes, physically-based dynamic DDC
463 performed as well as other schemes in hillslope simulations. Trial and error of building this
464 scheme identified that the distribution of flow-path lengths in the domain instead of the quantity
465 of particles added into the domain dominates the load distribution. This illustrated both the
466 mechanisms of load imbalance and the directions to build efficient physically-based LB schemes
467 in this context. This study realized the massively parallel computing of particle tracking in water
468 age simulations which is lacking in hydrologic modeling. It also demonstrated that LB have
469 practical importance enabling its applications at large scales with increased heterogeneity of flow
470 paths. The LB schemes can be borrowed to other hydrologic models using Lagrangian approach
471 and the parallelized EcoSLIM is a promising tool to accelerate the scientific progress of water
472 age studies.

473

474 **Acknowledgements**

475 This work was supported by the U.S. Department of Energy Office of Science, Offices of
476 Advanced Scientific Computing Research and Biological and Environmental Sciences IDEAS
477 project and Watershed Function Scientific Focus Area under Award Number DE-AC02-
478 05CH11231. The authors acknowledge high-performance computing support from Cheyenne
479 (doi:10.5065/D6RX99HX) provided by NCAR's Computational and Information Systems
480 Laboratory, sponsored by the National Science Foundation. Fortran/CUDA-Fortran code

481 (<https://github.com/aureliayang/EcoSLIM/tree/multi-GPU>) which can reproduce the results is
482 located on Github. Once the manuscript is accepted, this repository will be archived in Zenodo to
483 get a DOI for citation purpose.

484

485 **References**

- 486 Basu, N. B., Jindal, P., Schilling, K. E., Wolter, C. F., & Takle, E. S. (2012). Evaluation of
487 analytical and numerical approaches for the estimation of groundwater travel time
488 distribution. *Journal of Hydrology*, 475, 65-73.
489 doi:<https://doi.org/10.1016/j.jhydrol.2012.08.052>
- 490 Botter, G., Bertuzzo, E., & Rinaldo, A. (2011). Catchment residence and travel time distributions:
491 The master equation. *Geophysical Research Letters*, 38(11). doi:10.1029/2011gl047666
- 492 Boulmier, A., Raynaud, F., Abdennadher, N., & Chopard, B. (2019, 23-26 Sept. 2019). *On the*
493 *Benefits of Anticipating Load Imbalance for Performance Optimization of Parallel*
494 *Applications*. Paper presented at the 2019 IEEE International Conference on Cluster
495 Computing (CLUSTER).
- 496 Danesh-Yazdi, M., Klaus, J., Condon, L. E., & Maxwell, R. M. (2018). Bridging the gap
497 between numerical solutions of travel time distributions and analytical storage selection
498 functions. *Hydrological Processes*, 32(8), 1063-1076. Retrieved from <Go to
499 ISI>://WOS:000430466700006
- 500 Egorova, M. S., Dyachkov, S. A., Parshikov, A. N., & Zhakhovsky, V. V. (2019). Parallel SPH
501 modeling using dynamic domain decomposition and load balancing displacement of
502 Voronoi subdomains. *Computer Physics Communications*, 234, 112-125.
503 doi:<https://doi.org/10.1016/j.cpc.2018.07.019>
- 504 Eibl, S., & Rde, U. (2019). A systematic comparison of runtime load balancing algorithms for
505 massively parallel rigid particle dynamics. *Computer Physics Communications*, 244, 76-
506 85. doi:<https://doi.org/10.1016/j.cpc.2019.06.020>
- 507 Engdahl, N. B., & Maxwell, R. M. (2014). Approximating groundwater age distributions using
508 simple streamtube models and multiple tracers. *Advances in Water Resources*, 66, 19-31.
509 doi:<https://doi.org/10.1016/j.advwatres.2014.02.001>
- 510 Engdahl, N. B., & Maxwell, R. M. (2015). Quantifying changes in age distributions and the
511 hydrologic balance of a high-mountain watershed from climate induced variations in
512 recharge. *Journal of Hydrology*, 522, 152-162.
513 doi:<https://doi.org/10.1016/j.jhydrol.2014.12.032>
- 514 Engdahl, N. B., McCallum, J. L., & Massoudieh, A. (2016). Transient age distributions in
515 subsurface hydrologic systems. *Journal of Hydrology*, 543, 88-100.
516 doi:<https://doi.org/10.1016/j.jhydrol.2016.04.066>
- 517 Engdahl, N. B., Schmidt, M. J., & Benson, D. A. (2019). Accelerating and Parallelizing
518 Lagrangian Simulations of Mixing-Limited Reactive Transport. *Water Resources*
519 *Research*, 55(4), 3556-3566. doi:10.1029/2018wr024361
- 520 Fattebert, J. L., Richards, D. F., & Glosli, J. N. (2012). Dynamic load balancing algorithm for
521 molecular dynamics based on Voronoi cells domain decompositions. *Computer Physics*
522 *Communications*, 183(12), 2608-2615. doi:<https://doi.org/10.1016/j.cpc.2012.07.013>

523 Furuichi, M., & Nishiura, D. (2017). Iterative load-balancing method with multigrid level
524 relaxation for particle simulation with short-range interactions. *Computer Physics*
525 *Communications*, 219, 135-148. doi:<https://doi.org/10.1016/j.cpc.2017.05.015>

526 Gomez, J. D., & Wilson, J. L. (2013). Age distributions and dynamically changing hydrologic
527 systems: Exploring topography-driven flow. *Water Resources Research*, 49(3), 1503-
528 1522. doi:<https://doi.org/10.1002/wrcr.20127>

529 Ji, X. H., Luo, M. L., & Wang, X. S. (2019). Accelerating Streamline Tracking in Groundwater
530 Flow Modeling on GPUs. *Groundwater*. doi:10.1111/gwat.12959

531 Jing, M., Heße, F., Kumar, R., Kolditz, O., Kalbacher, T., & Attinger, S. (2019). Influence of
532 input and parameter uncertainty on the prediction of catchment-scale groundwater travel
533 time distributions. *Hydrol. Earth Syst. Sci.*, 23(1), 171-190. doi:10.5194/hess-23-171-
534 2019

535 Jing, M., Kumar, R., Heße, F., Thober, S., Rakovec, O., Samaniego, L., & Attinger, S. (2020).
536 Assessing the response of groundwater quantity and travel time distribution to 1.5, 2, and
537 3 °C global warming in a mesoscale central German basin. *Hydrol. Earth Syst.*
538 *Sci.*, 24(3), 1511-1526. doi:10.5194/hess-24-1511-2020

539 Kollet, S. J., & Maxwell, R. M. (2008a). Capturing the influence of groundwater dynamics on
540 land surface processes using an integrated, distributed watershed model. *Water Resources*
541 *Research*, 44(2). doi:Artn W0240210.1029/2007wr006004

542 Kollet, S. J., & Maxwell, R. M. (2008b). Demonstrating fractal scaling of baseflow residence
543 time distributions using a fully-coupled groundwater and land surface model.
544 *Geophysical Research Letters*, 35(7). Retrieved from <Go to
545 ISI>://WOS:000254716500001

546 Kunaseth, M., Richards, D. F., Glosli, J. N., Kalia, R. K., Nakano, A., & Vashishta, P. (2013).
547 Analysis of scalable data-privatization threading algorithms for hybrid MPI/OpenMP
548 parallelization of molecular dynamics. *The Journal of Supercomputing*, 66(1), 406-430.
549 doi:10.1007/s11227-013-0915-x

550 Maxwell, R. M., Condon, L. E., Danesh-Yazdi, M., & Bearup, L. A. (2019). Exploring source
551 water mixing and transient residence time distributions of outflow and evapotranspiration
552 with an integrated hydrologic model and Lagrangian particle tracking approach.
553 *Ecohydrology*, 12(1). Retrieved from <Go to ISI>://WOS:000454601400016

554 Maxwell, R. M., Condon, L. E., Kollet, S. J., Maher, K., Haggerty, R., & Forrester, M. M. (2016).
555 The imprint of climate and geology on the residence times of groundwater. *Geophysical*
556 *Research Letters*, 43(2), 701-708. doi:10.1002/2015gl066916

557 McCallum, J. L., Engdahl, N. B., Ginn, T. R., & Cook, P. G. (2014). Nonparametric estimation
558 of groundwater residence time distributions: What can environmental tracer data tell us
559 about groundwater residence time? *Water Resources Research*, 50(3), 2022-2038.
560 doi:10.1002/2013wr014974

561 McGuire, K. J., McDonnell, J. J., Weiler, M., Kendall, C., McGlynn, B. L., Welker, J. M., &
562 Seibert, J. (2005). The role of topography on catchment-scale water residence time.
563 *Water Resources Research*, 41(5). doi:10.1029/2004wr003657

564 Pollock, D. W. (2016). *User guide for MODPATH Version 7—A particle-tracking model for*
565 *MODFLOW* (2016-1086). Retrieved from Reston, VA:
566 <http://pubs.er.usgs.gov/publication/ofr20161086>

567 Ruetsch, G., & Fatica, M. (2014). *CUDA Fortran for scientists and engineers : best practices for*
568 *efficient CUDA Fortran programming*. Amsterdam Boston: Morgan Kaufmann, an
569 imprint of Elsevier.

570 Sprenger, M., Stumpp, C., Weiler, M., Aeschbach, W., Allen, S. T., Benettin, P., . . . Werner, C.
571 (2019). The Demographics of Water: A Review of Water Ages in the Critical Zone.
572 *Reviews of Geophysics*, 57(3), 800-834. doi:<https://doi.org/10.1029/2018RG000633>

573 Starn, J. J., Kauffman, L. J., Carlson, C. S., Reddy, J. E., & Fienen, M. N. (2021). Three-
574 Dimensional Distribution of Groundwater Residence Time Metrics in the Glaciated
575 United States Using Metamodels Trained on General Numerical Simulation Models.
576 *Water Resources Research*, 57(2), e2020WR027335.
577 doi:<https://doi.org/10.1029/2020WR027335>

578 Tran, H., Zhang, J., Cohard, J.-M., Condon, L. E., & Maxwell, R. M. (2020). Simulating
579 Groundwater-Streamflow Connections in the Upper Colorado River Basin. *Groundwater*,
580 58(3), 392-405. doi:10.1111/gwat.13000

581 Weill, S., Lesparre, N., Jeannot, B., & Delay, F. (2019). Variability of Water Transit Time
582 Distributions at the Strengbach Catchment (Vosges Mountains, France) Inferred Through
583 Integrated Hydrological Modeling and Particle Tracking Algorithms. *Water*, 11(12), 2637.
584 Retrieved from <https://www.mdpi.com/2073-4441/11/12/2637>

585 Wilusz, D. C., Harman, C. J., Ball, W. B., Maxwell, R. M., & Buda, A. R. (2019). Using particle
586 tracking to understand flow paths, age distributions, and the paradoxical origins of the
587 inverse storage effect in an experimental catchment. *Water Resources Research*, n/a(n/a),
588 e24397. doi:10.1029/2019wr025140

589 Yang, C., Li, H.-Y., Fang, Y., Cui, C., Wang, T., Zheng, C., . . . Yang, X. (2020). Effects of
590 Groundwater Pumping on Ground Surface Temperature: A Regional Modeling Study in
591 the North China Plain. *Journal of Geophysical Research: Atmospheres*, 125(9),
592 e2019JD031764. doi:10.1029/2019jd031764

593 Yang, C., Zhang, Y.-K., Liang, X., Olschanowsky, C., Yang, X., & Maxwell, R. (2021).
594 Accelerating the Lagrangian particle tracking of residence time distributions and source
595 water mixing towards large scales. *Computers & Geosciences*, 104760.
596 doi:<https://doi.org/10.1016/j.cageo.2021.104760>

597 Yang, J., Heibüchel, I., Musolff, A., Reinstorf, F., & Fleckenstein, J. H. (2018). Exploring the
598 Dynamics of Transit Times and Subsurface Mixing in a Small Agricultural Catchment.
599 *Water Resources Research*, 54(3), 2317-2335. doi:10.1002/2017wr021896

600