

Author's Accepted Manuscript

Estimation of theoretical maximum speedup ratio for parallel computing of grid-based distributed hydrological models

Junzhi Liu, A-Xing Zhu, Cheng-Zhi Qin



www.elsevier.com/locate/cageo

PII: S0098-3004(13)00156-8
DOI: <http://dx.doi.org/10.1016/j.cageo.2013.04.030>
Reference: CAGEO3194

To appear in: *Computers & Geosciences*

Received date: 4 March 2013
Revised date: 26 April 2013
Accepted date: 30 April 2013

Cite this article as: Junzhi Liu, A-Xing Zhu, Cheng-Zhi Qin, Estimation of theoretical maximum speedup ratio for parallel computing of grid-based distributed hydrological models, *Computers & Geosciences*, <http://dx.doi.org/10.1016/j.cageo.2013.04.030>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 **Estimation of theoretical maximum speedup ratio for parallel**
2 **computing of grid-based distributed hydrological models**

3 Junzhi Liu^{a,b}, A-Xing Zhu^{a,c*}, Cheng-Zhi Qin^a

4 ^a State Key Lab of Resources and Environmental Information System, Institute of
5 Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences,
6 Beijing, 100101, China.

7 ^b University of Chinese Academy of Sciences, Beijing, 100049, China

8 ^c Department of Geography, University of Wisconsin-Madison, Madison, WI 53706,
9 USA

10 *Corresponding author. Email: axing@lreis.ac.cn; Fax: 86-10-64889630; Telephone:
11 86-13810787680; Address: 11A, Datun Road, Chaoyang district, Beijing, 100101,
12 China.

13

14 **Abstract:** Theoretical Maximum Speedup Ratio (TMSR) can be used as a goal for
15 improving parallel computing methods for distributed hydrological models. Different
16 types of distributed hydrological models need different TMSR estimation methods
17 because of the different computing characteristics of models. Existing TMSR
18 estimation methods, such as those for sub-basin based distributed hydrological models,
19 are inappropriate for grid-based distributed hydrological models. In this paper, we
20 proposed a TMSR estimation method suitable for grid-based distributed hydrological
21 models. With this method, TMSRs for hillslope processes and channel routing
22 processes are calculated separately and then combined to obtain the overall TMSR. A

1 branch-and-bound algorithm and a critical path heuristic algorithm are used to
2 estimate TMSRs for parallel computing of hillslope processes and channel routing
3 processes, respectively. The overall TMSR is calculated according to the proportions
4 of computing these two types of processes. A preliminary application showed that the
5 more the number of sub-basins, the larger the TMSRs and that the compact
6 watersheds had larger TMSRs than the long narrow watersheds.

7

8 **Keywords:** Theoretical maximum speedup ratio; Parallel computing; Grid-based
9 distributed hydrological model; Critical path heuristic algorithm

10

11 **1. Introduction**

12 Theoretical Maximum Speedup Ratio (TMSR) in parallel computing of
13 distributed hydrological models is defined as the potential maximum speedup ratio
14 value that can be obtained under a given number of processors (Wang et al., 2012). It
15 can be used as a goal for improving parallel computing methods for distributed
16 hydrological models. Different types of distributed hydrological models need different
17 TMSR estimation methods because of the different computing characteristics of
18 models.

19 Currently, there has been little published work on estimating TMSR for parallel
20 computing of grid-based distributed hydrological models, which are an important type
21 of hydrological model (Borah and Bera, 2004). Existing related researches mainly
22 focus on TMSR estimation for sub-basin based distributed hydrological models

1 (Apostolopoulos and Georgakakos, 1997; Li et al., 2011; Wang et al., 2012). For
2 example, Apostolopoulos and Georgakakos (1997) proposed a method for estimating
3 TMSR by representing sub-basins in a watershed using Directed Acyclic Graph
4 (DAG). Li et al. (2011) and Wang et al. (2012) estimated TMSR for parallel
5 computing of sub-basin based distributed hydrological models by analyzing the basin
6 width function.

7 These TMSR estimation methods for sub-basin based distributed hydrological
8 models all assume that the amounts of computation for each sub-basin are the same
9 (Wang et al., 2012), which is not true for grid-based distributed hydrological models
10 as the numbers of grid cells in each sub-basin vary widely. So existing TMSR
11 estimation methods for sub-basin based distributed hydrological models are
12 inappropriate for grid-based distributed hydrological models. The objective of this
13 paper is to propose a new TMSR estimation method suitable for grid-based distributed
14 hydrological models considering their computing characteristics.

15 **2. Basic idea**

16 This paper also uses sub-basin as the basic unit for parallel computing because
17 communication among sub-basins is low and sub-basin based parallel computing has
18 been proved to be effective (Vivoni et al., 2011; Wang et al., 2011). According to the
19 computation characteristics processes in sub-basins can be divided into two types:
20 hillslope processes (such as infiltration and evaporation) and channel routing
21 processes (such as channel flow and sediment routing). For hillslope processes
22 calculations are independent among sub-basins while for channel routing processes

1 calculations for one sub-basin depends on the calculation results of its upstream
2 sub-basins. Thus, different strategies should be adopted for parallel computing of
3 these two types of processes. Different TMSR estimation methods should also be used
4 accordingly. Our idea is first to estimate the TMSRs for these two types of processes,
5 respectively and then combine them to obtain an overall TMSR.

6 **3. TMSR estimation method**

7 **3.1 Assumptions**

8 The assumptions for the TMSR estimation method in this paper are as follows:

- 9 (1) Groundwater follows the direction of overland flow and there is only
10 interaction between upstream and downstream sub-basins (Li et al., 2011).
- 11 (2) Communication overhead among sub-basins is low, so it can be ignored for
12 simplicity (Li et al., 2011; Wang et al., 2012).

13 **3.2 Estimation process**

14 **Step 1: Estimating TMSR for parallel computing of hillslope processes**

15 For parallel computing of hillslope processes, due to the fact that there are no
16 dependencies among sub-basins, the goal of task scheduling is to assign n
17 independent tasks (one task for each sub-basin) to m processors with the objective of
18 minimizing execution time. The execution time of hillslope processes for a sub-basin
19 can be represented by the number of hillslope cells in this sub-basin because generally
20 the more the number of cells, the more the execution time there will be. Thus, the
21 serial computing time can be represented by the total number of hillslope cells in the
22 watershed.

1 Estimating TMSR in this situation is a classical Multiprocessor Scheduling
2 Problem. In this paper, a branch-and-bound algorithm for exact solution of the
3 problem (Dell'Amico and Martello, 1995) is used to get the minimum execution time
4 of parallelizing hillslope processes under m processors. Then TMSR for hillslope
5 processes is calculated by dividing the serial computing time by the minimum parallel
6 computing time.

7 **Step 2: Estimating TMSR for parallel computing of channel routing processes**

8 According to assumption (1) the dependence among sub-basins for calculations of
9 channel routing processes can be represented by a tree-structured DAG
10 (Apostolopoulos and Georgakakos, 1997; Li et al., 2010) (Fig. 1). Each node in the
11 DAG represents the calculation task of channel routing processes for a sub-basin. The
12 goal of task scheduling for channel routing processes is to distribute tasks in the DAG
13 among a given number of processors to achieve minimum execution time. The
14 execution time of channel routing processes for a sub-basin can be represented by the
15 number of channel cells in this sub-basin and the serial computing time can be
16 represented by the total number of channel cells in the watershed.

17 (Fig. 1 is about here)

18 The scheduling problem for this case belongs to the static task scheduling
19 problems and the critical path heuristic algorithm has been proved to be effective on
20 solving such problems (Shirazi et al., 1990). This paper uses the critical path heuristic
21 algorithm to obtain the minimum execution time of channel routing processes
22 parallelized on a given number of processors. To illustrate this algorithm, some terms

1 are defined as below.

2 (1) The “accumulation time” of a processor P_i , $AT(P_i)$, is the total time needed for
3 a processor to finish all the tasks assigned to it.

4 (2) A node is said to be “mature” if it is ready to be assigned to a processor, that
5 is, this node has no precedent nodes or all its precedent nodes have been
6 already completed.

7 (3) The “exit path” of a node means the path from this node to the exit node
8 (representing the outlet sub-basin) and the length of exit path is the sum of
9 execution time of every node in this path..

10 The longest exit path (the critical path, e.g. G-D-C-A in Fig.1) determines the
11 minimum possible execution time for parallel computing. So the basic principle of
12 this algorithm is to assign the mature node with the largest exit path length at current
13 time. Load balance is achieved by assigning this mature node to the processor with
14 the minimum accumulation time. The steps of the algorithm are given as follows:

15 (1) Calculate the length of exit path for each node in the DAG. This can be
16 accomplished by a recursive algorithm:

17 a) The length of exit path for the exit node is equal to the execution time of
18 this node;

19 b) The length of exit path for other nodes is equal to the execution time of
20 this node plus the exit path length of its downstream node.

21 (2) Assign tasks following the steps below:

22 a) Choose the processor with the minimum accumulation time, P_{min} ;

- 1 b) Find a mature node with the largest exit path length;
- 2 c) If this mature node can be found, assign this node to P_{min} ;
- 3 d) If there is no mature node available, assign a dummy node to P_{min} . The
- 4 length of the dummy node is equal to the difference between $AT(P_i)$ and
- 5 $AT(P_{min})$. Here $AT(P_i)$ is the smallest accumulation time that is strictly
- 6 larger than $AT(P_{min})$.
- 7 e) Repeat a)-d) until all tasks are completed.

8 When the minimum parallel computing time is obtained using the above

9 algorithm TMSR for channel routing processes is calculated by dividing the serial

10 computing time by the minimum parallel computing time.

11 **Step 3: Estimating the overall TMSR**

12 The overall TMSR can be calculated by combing the TMSRs of hillslope

13 processes and of channel routing processes:

$$14 \quad R_{all} = 1 / (p_{hs} / R_{hs} + p_{ch} / R_{ch}) \quad (1)$$

15 where R_{all} is the overall TMSR; R_{hs} and R_{ch} are the TMSRs for hillslope processes and

16 channel routing processes, respectively; p_{hs} and p_{ch} are the proportions of computing

17 hillslope processes and channel routing processes in the total amount of computations

18 respectively. $p_{hs} + p_{ch} = 1$. The values of p_{hs} and p_{ch} are determined by the simulation

19 methods used for hillslope and channel processes and these values can be determined

20 by profiling model simulations using a performance profiler (for example, Intel[®]

21 VTune[™]).

1 curves reached plateaus after the numbers of processes exceeded certain thresholds.
2 These thresholds represent the maximum numbers of processors that parallel
3 computing of distributed hydrological models can make use of. Parallelizing channel
4 routing processes has much smaller TMSRs than parallelizing hillslope processes, so
5 reducing p_{ch} improves overall TMSRs dramatically.

6 (Fig. 3 is about here)

7 Fig. 4 shows the estimated TMSRs under enough number of processors for
8 different numbers of sub-basins in the six watersheds assuming p_{ch} is 10%. It is
9 obvious that the more the number of sub-basins, the larger the TMSRs were. For the
10 hillslope processes, the relationships between TMSR and number of sub-basins were
11 quantitatively similar for watersheds with different areas and shapes. This reflected
12 the similarities among all watersheds to some extent. Meanwhile for the channel
13 routing processes and the whole model, the relationships between TMSR and number
14 of sub-basins were affected by the shapes of watersheds. Generally, TMSRs are larger
15 for compact watersheds (e.g. Fenkeng2) than for long narrow watersheds (e.g.
16 Qingshuihe1 and Fenkeng1).

17 (Fig. 4 is about here)

18 **5 Conclusion**

19 This paper proposed a TMSR estimation method for parallel computing of
20 grid-based distributed hydrological models. In this method, TMSRs for hillslope
21 processes and channel routing processes are calculated separately. Then they are
22 combined to obtain overall TMSR. A preliminary application showed that the more

1 the number of sub-basins, the larger the TMSRs and that the compact watersheds had
2 larger TMSRs than the long narrow watersheds.

3

4 **Acknowledgements**

5 This study was funded by the National High-Tech Research and Development Program of
6 China (No. 2011AA120305) and the National Natural Science Foundation of China (No.
7 41023010). This study was also partly funded by the Program of International S&T
8 Cooperation, MOST of China (No. 2010DFB24140). We thank Prof. Mauro Dell'Amico for
9 his kindly providing the programming library for solving the multiprocessor scheduling
10 problem.

11

12 **References**

- 13 Apostolopoulos, T.K., Georgakakos, K.P., 1997. Parallel computation for streamflow
14 prediction with distributed hydrologic models. *Journal of Hydrology* 197(1-4),
15 1-24.
- 16 Borah, D.K., Bera, M., 2004. Watershed-scale hydrologic and nonpoint-source
17 pollution models: Review of applications. *Transactions of the ASAE* 47(3),
18 789-803.
- 19 Dell'Amico, M., Martello, S., 1995. Optimal scheduling of tasks on identical parallel
20 processors. *ORSA Journal on Computing* 7(2), 191-200.
- 21 Li, T.J., Wang, G.Q., Chen, J., 2010. A modified binary tree codification of drainage
22 networks to support complex hydrological models. *Computers & Geosciences*

- 1 36(11), 1427-1435.
- 2 Li, T.J., Wang, G.Q., Chen, J., Wang, H., 2011. Dynamic parallelization of
3 hydrological model simulations. *Environmental Modelling & Software* 26(12),
4 1736-1746.
- 5 Schumm, S. A., 1956. Evolution of drainage systems and slopes in badlands at Perth
6 Amboy, New Jersey. *Geological Society of American Bulletin*, 67, 597-646.
- 7 Shirazi, B., Wang, M., Pathak, G., 1990. Analysis and evaluation of heuristic methods
8 for static task scheduling. *Journal of Parallel and Distributed Computing* 10(3),
9 222-232.
- 10 Vivoni, E.R., Mascaro, G., Mniszewski, S., Fasel, P., Springer, E.P., Ivanov, V.Y., Bras,
11 R.L., 2011. Real-world hydrologic assessment of a fully-distributed hydrological
12 model in a parallel computing environment. *Journal of Hydrology* 409(1-2),
13 483-496.
- 14 Wang, H., Wang, G.Q., Gao, J., Fu, X.D., 2010. Parallel characteristics of river basin
15 based on temporal-spatial-discrete approach. *Sciencepaper Online* 5(7), 494-498
16 [in Chinese].
- 17 Wang, H., Fu, X.D., Wang, G.Q., Li, T.J., Gao, J., 2011. A common parallel computing
18 framework for modeling hydrological processes of river basins. *Parallel*
19 *Computing* 37(6-7), 302-315.
- 20 Wang, H., Zhou, Y., Fu, X.D., Gao, J., Wang, G.Q., 2012. Maximum speedup ratio
21 curve (MSC) in parallel computing of the binary-tree-based drainage network.
22 *Computers & Geosciences* 38(1), 127-135.

1 **Figure List**

2 Fig. 1: Example of representing the dependence relationships among sub-basins (a)
3 using a tree-structured Directed Acyclic Graph (b). Node A is the exit node,
4 representing the outlet sub-basin. The first number in a DAG node represents the
5 execution time of corresponding channel routing processes and the second number in
6 a DAG node represents the exit path length correspondingly. Fig. 1(a) was modified
7 from Wang et al. (2010).

8 Fig.2: Location, DEM, and the drainage network of the study areas, including the
9 Qingshuihe watershed and its two sub-watersheds (Qingshuihe1 and Qingshuihe2),
10 the Fenkeng watershed and its two sub-watersheds (Fenkeng1 and Fenkeng2).

11 Fig. 3: Estimated TMSR curves for parallel computing of both sub-processes and the
12 model in Qingshuihe watershed when the number of sub-basins is 101. “Hillslope
13 processes” and “Channel routing processes” in the legend means the TMSR curves for
14 parallel computing of hillslope processes and channel routing processes, respectively.
15 “Overall” means the overall TMSR curve of the model. “ p_{ch} ” means the proportion of
16 computing channel routing processes in the total amount of computations.

17 Fig. 4: Estimated TMSRs under enough number of processors for different numbers
18 of sub-basins in the six watersheds assuming p_{ch} (the proportion of computing channel
19 routing processes in the total amount of computations) is 10%.

20

21 **Table List**

22 Table 1: Area, shape and sub-basin division information of the six watersheds.

1

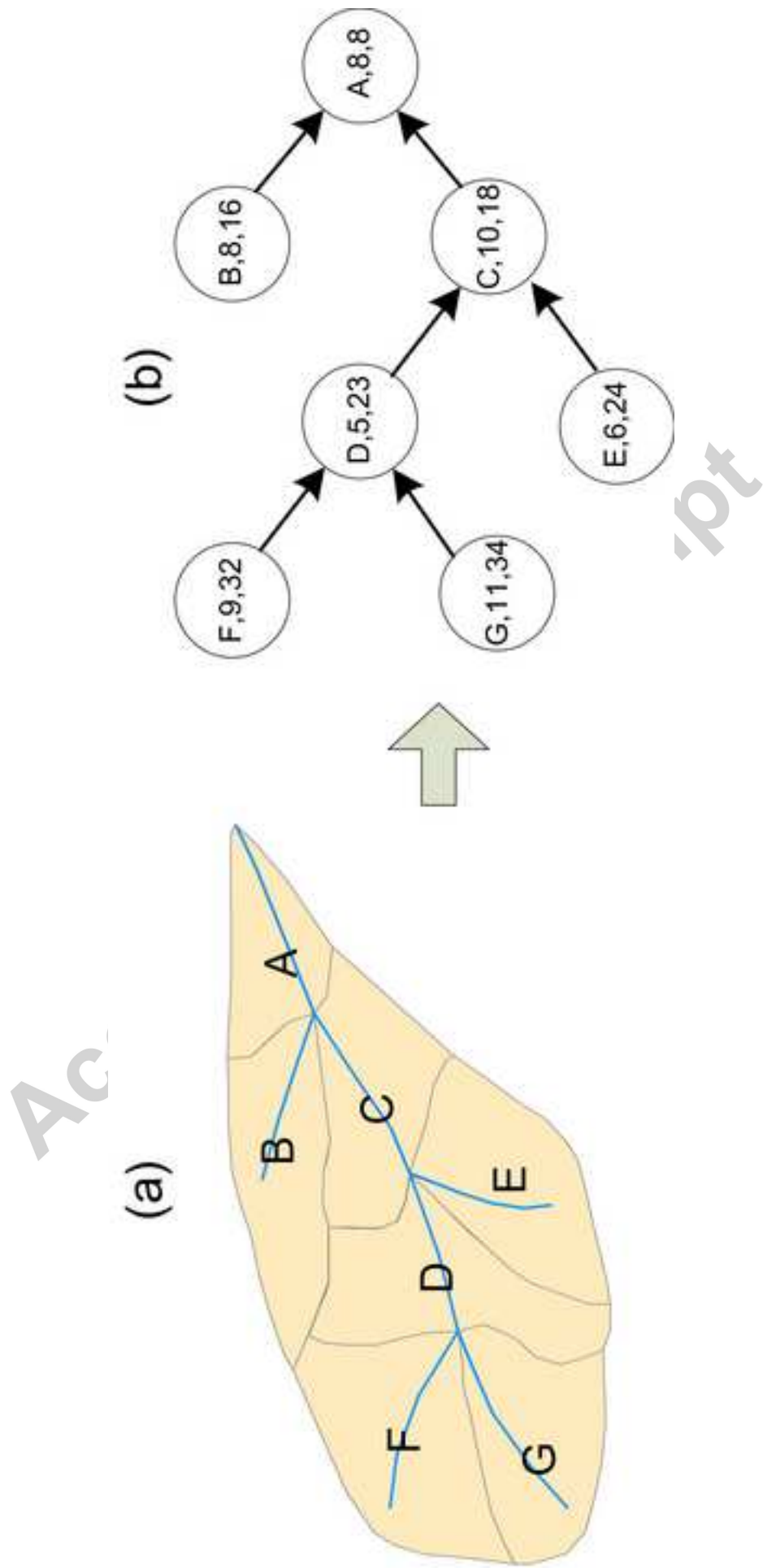
2 **Table 1**

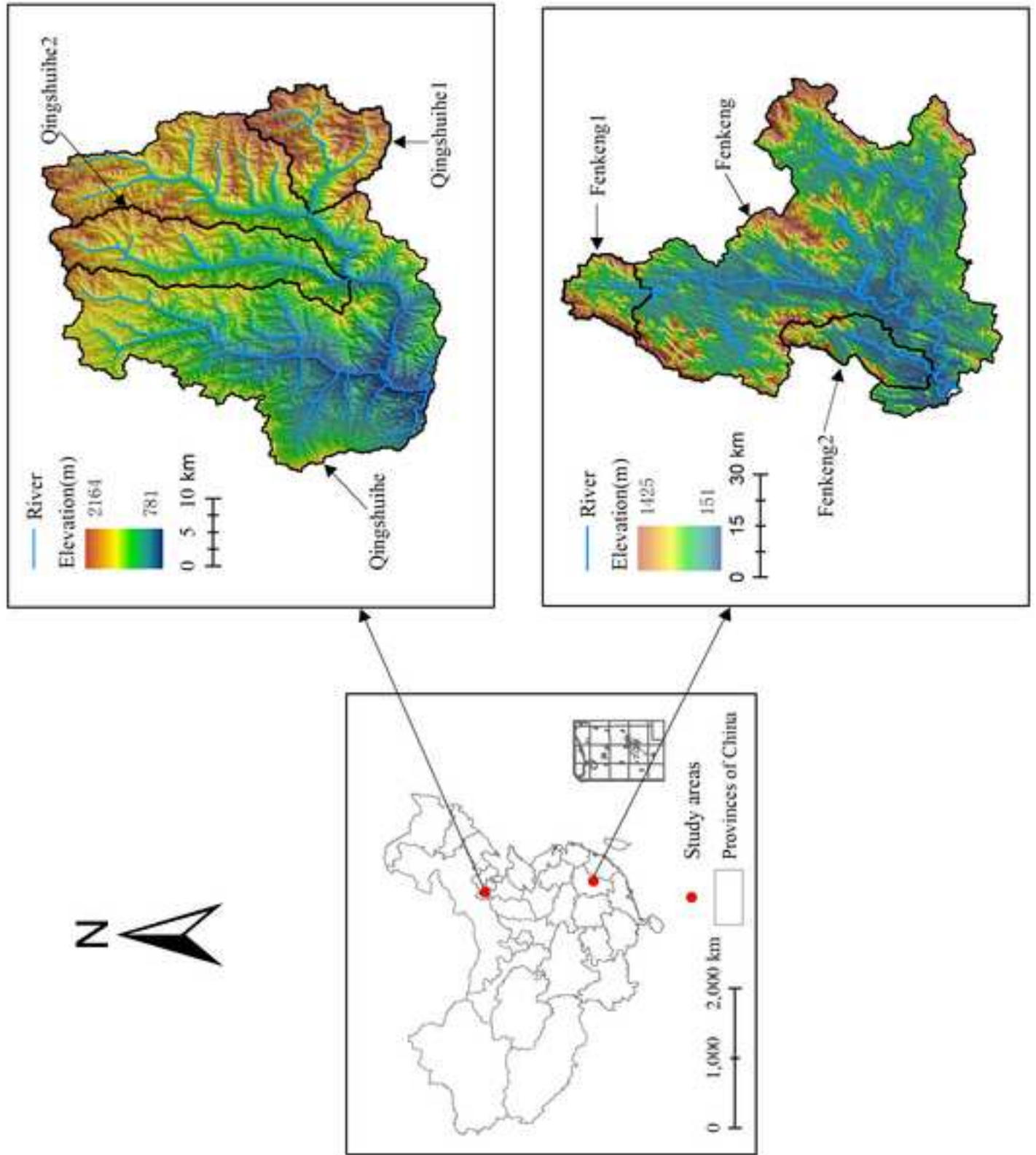
3 Area, shape and sub-basin division information of the six watersheds.

Name	Area (km ²)	Basin elongation ratio	List of sub-basin division numbers
Qingshuihe	2178	0.5	[55, 101, 203, 277, 491]
Qingshuihe1	231	0.56	[53, 134, 253, 346, 498]
Qingshuihe2	354	0.38	[45, 95, 227, 324, 507]
Fenkeng	6323	0.47	[57, 103, 193, 341, 501]
Fenkeng1	373	0.62	[45, 95, 235, 324, 487]
Fenkeng2	475	0.36	[45, 79, 193, 257, 476]

4

- 5 ➤ A new method for estimating theoretical maximum speedup ratio was proposed.
- 6 ➤ TMSRs for hillslope processes and channel processes were estimated
- 7 respectively.
- 8 ➤ Channel processes have much smaller TMSRs than hillslope processes.
- 9 ➤ For one watershed, the more the number of sub-basins, the larger the TMSRs.
- 10 ➤ The compact watersheds had larger TMSRs than the long narrow watersheds.





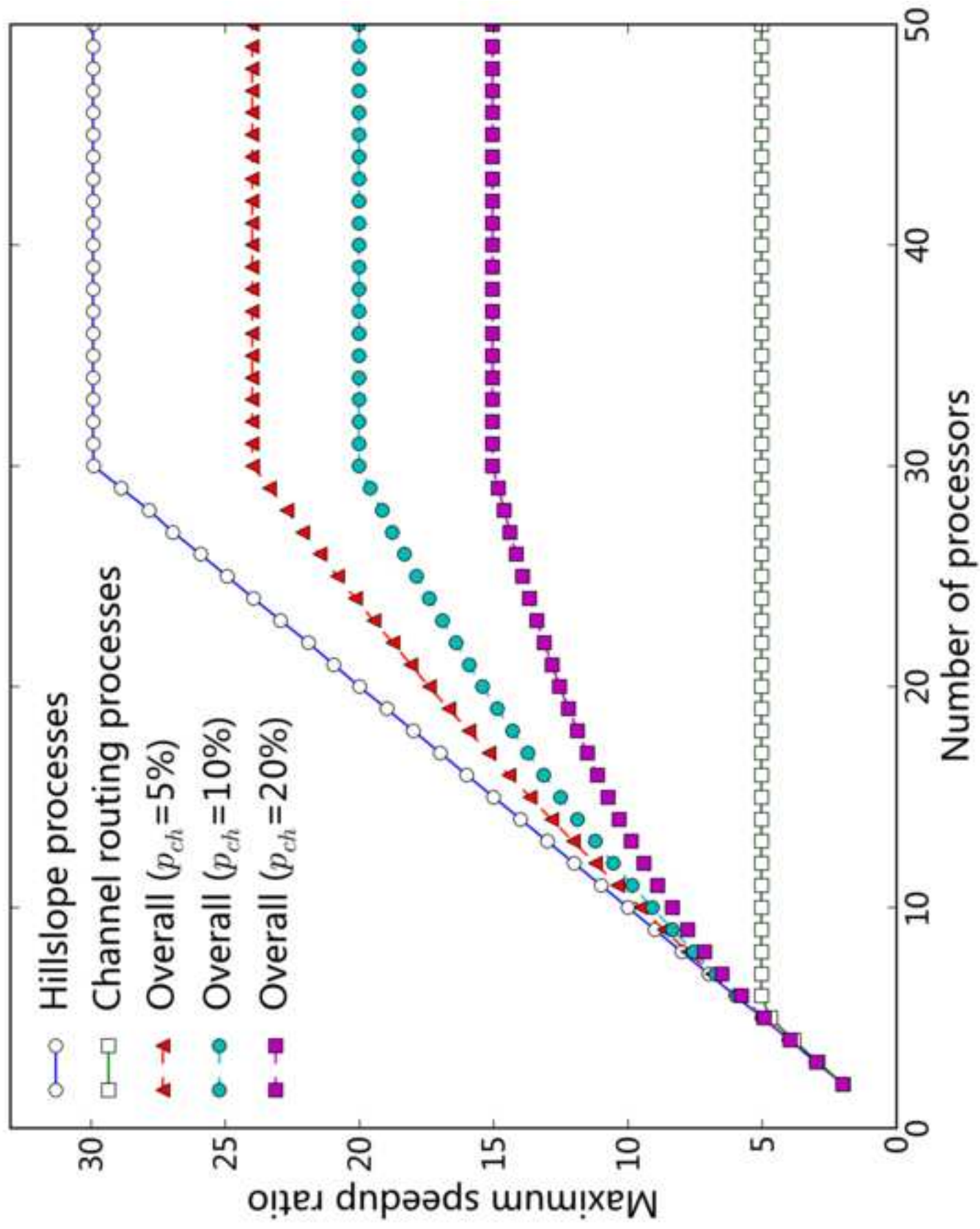


Fig. 3

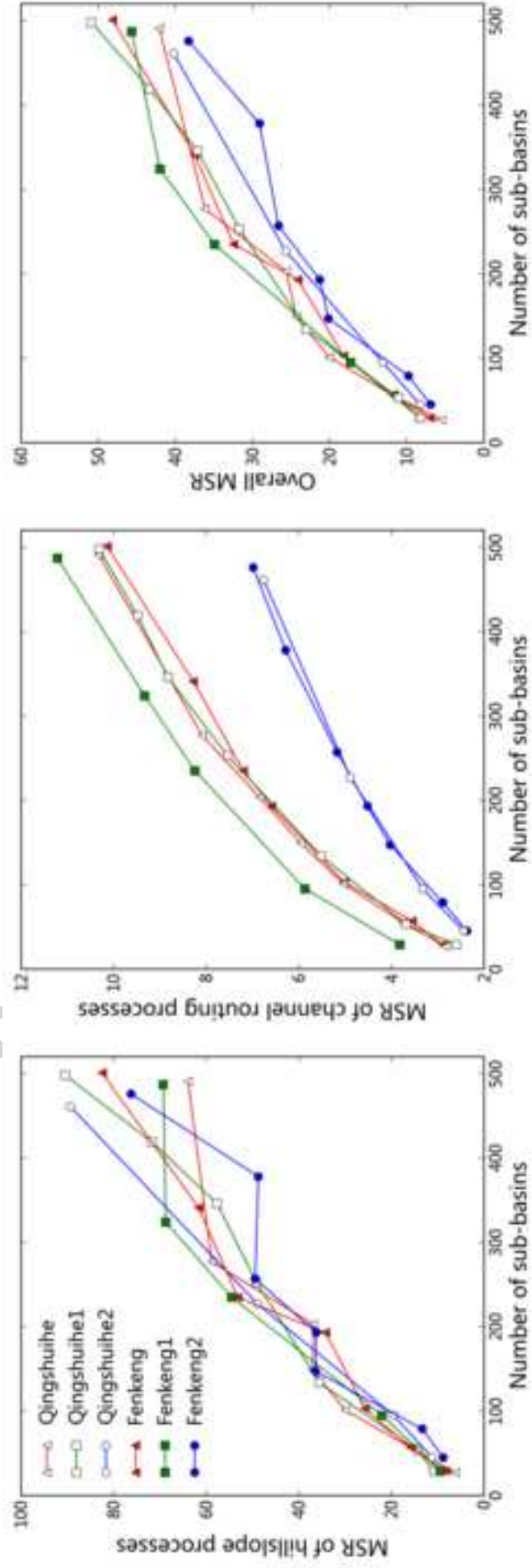


Fig. 4