

UDC: 519.168

Multi-agent local voting protocol for online DAG scheduling

N. A. Zhitnukhin^{1,a}, A. Yu. Zhadan^{1,b}, I. V. Kondratov^{1,c},
A. L. Allahverdyan^{1,d}, O. N. Granichin^{1,e}, O. L. Petrosian^{1,f},
A. V. Romanovskii^{2,g}, V. S. Kharin^{2,h}

¹Federal State Budgetary Educational Institution of Higher Education
“Saint Petersburg State University” (SPbU),
7/9 Universitetskaya Embankment, Saint Petersburg, 199034, Russia
²Huawei Technologies Co., Ltd.,
17/2 Krylatskaya st., Moscow, 121614, Russia

E-mail: ^a kola1197mail@gmail.com, ^b anastasiya_markel@mail.ru, ^c kondratov.ivan.vladimirovich@gmail.com,
^d aleka_alexander@mail.ru, ^e o.granichin@spbu.ru, ^f petrosian.ovanes@yandex.ru,
^g aleksei.romanovskii@huawei.com, ^h vitaliy.kharin@huawei.com

Received 22.10.2024.

Accepted for publication 25.11.2024.

Scheduling computational workflows represented by directed acyclic graphs (DAGs) is crucial in many areas of computer science, such as cloud/edge tasks with distributed workloads and data mining. The complexity of online DAG scheduling is compounded by the large number of computational nodes, data transfer delays, heterogeneity (by type and processing power) of executors, precedence constraints imposed by DAG, and the nonuniform arrival of tasks. This paper introduces the Multi-Agent Local Voting Protocol (MLVP), a novel approach focused on dynamic load balancing for DAG scheduling in heterogeneous computing environments, where executors are represented as agents. The MLVP employs a local voting protocol to achieve effective load distribution by formulating the problem as a differentiated consensus achievement. The algorithm calculates an aggregated DAG metric for each executor-node pair based on node dependencies, node availability, and executor performance. The balance of these metrics as a weighted sum is optimized using a genetic algorithm to assign tasks probabilistically, achieving efficient workload distribution via information sharing and reaching consensus among the executors across the system and thus improving makespan. The effectiveness of the MLVP is demonstrated through comparisons with the state-of-the-art DAG scheduling algorithm and popular heuristics such as DONF, FIFO, Min-Min, and Max-Min. Numerical simulations show that MLVP achieves makespan improvements of up to 70% on specific graph topologies and an average makespan reduction of 23.99% over DONF (state-of-the-art DAG scheduling heuristic) across randomly generated diverse set of DAGs. Notably, the algorithm’s scalability is evidenced by enhanced performance with increasing numbers of executors and graph nodes.

Keywords: multi-agent, local voting protocol, scheduling, directed acyclic graph

Citation: *Computer Research and Modeling*, 2025, vol. 17, no. 1, pp. 29–44.

The authors acknowledge Saint Petersburg State University for a research project 95438429.

© 2025 Nikolay A. Zhitnukhin, Anastasia Yu. Zhadan, Ivan V. Kondratov, Alexander L. Allahverdyan, Oleg N. Granichin, Ovanes L. Petrosian, Aleksei V. Romanovskii, Vitaliy S. Kharin

This work is licensed under the Creative Commons Attribution-NoDerivs 3.0 Unported License.
To view a copy of this license, visit <http://creativecommons.org/licenses/by-nd/3.0/>
or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

УДК: 519.168

Многоагентный протокол локального голосования для онлайн-планирования DAG

Н. А. Житнухин^{1,a}, А. Ю. Жадан^{1,b}, И. В. Кондратов^{1,c},
А. Л. Аллахвердян^{1,d}, О. Н. Граничин^{1,e}, О. Л. Петросян^{1,f},
А. В. Романовский^{2,g}, В. С. Харин^{2,h}

¹Федеральное государственное бюджетное образовательное учреждение высшего образования
«Санкт-Петербургский государственный университет»,

Россия, 199034, г. Санкт-Петербург, Университетская набережная, д. 7/9

²Общество с ограниченной ответственностью «Техкомпания Хуавэй»,
Россия, 121614, г. Москва, ул. Крылатская, д. 17, корп. 2

E-mail: ^a kola1197mail@gmail.com, ^b anastasiya_markel@mail.ru, ^c kondratov.ivan.vladimirovich@gmail.com,
^d aleka_alexander@mail.ru, ^e o.granichin@spbu.ru, ^f petrosian.ovanes@yandex.ru,
^g aleksei.romanovskii@huawei.com, ^h vitaliy.kharin@huawei.com

Получено 22.10.2024.

Принято к публикации 25.11.2024.

Планирование вычислительных рабочих процессов, представленных направленными ациклическими графами (DAG), имеет ключевое значение для многих областей информатики, таких как облачные/edge задачи с распределенной рабочей нагрузкой и анализ данных. Сложность онлайн-планирования DAG усугубляется большим количеством вычислительных узлов, задержками передачи данных, неоднородностью (по типу и вычислительной мощности) исполнителей, ограничениями предшествования, накладываемыми DAG, и неравномерностью поступления задач. В данной статье представлен мультиагентный протокол локального голосования (MLVP) — новый подход, ориентированный на динамическое распределение нагрузки при планировании DAG в гетерогенных вычислительных средах, где исполнители представлены в виде агентов. MLVP использует протокол локального голосования для достижения эффективного распределения нагрузки, формулируя проблему как дифференцированное достижение консенсуса. Алгоритм вычисляет агрегированную метрику DAG для каждой пары исполнитель – узел на основе зависимостей между узлами, доступности узлов и производительности исполнителей. Баланс этих метрик как взвешенная сумма оптимизируется с помощью генетического алгоритма для вероятностного распределения задач, что позволяет добиться эффективного распределения рабочей нагрузки за счет обмена информацией и достижения консенсуса между исполнителями всей системы и, таким образом, улучшить время выполнения. Эффективность MLVP демонстрируется путем сравнения с современным алгоритмом планирования DAG и популярными эвристиками, такими как DONF, FIFO, Min-Min и Max-Min. Численное моделирование показывает, что MLVP достигает улучшения makespan до 70 % на определенных топологиях графов и среднего сокращения makespan на 23,99 % по сравнению с DONF (современная эвристика планирования DAG) на случайно сгенерированном разнообразном наборе DAG. Примечательно, что масштабируемость алгоритма подтверждается ростом производительности при увеличении числа исполнителей и узлов графа.

Ключевые слова: многоагентные системы, протокол локального голосования, построение расписаний, направленный ациклический граф

Договор на проведение научно-исследовательских работ TC20210202013.

© 2025 Николай Алексеевич Житнухин, Анастасия Юрьевна Жадан, Иван Владимирович Кондратов, Александр Львович Аллахвердян, Олег Николаевич Граничин, Ованес Леонович Петросян, Алексей Валентинович Романовский, Виталий Сергеевич Харин

Статья доступна по лицензии Creative Commons Attribution-NoDerivs 3.0 Unported License.
Чтобы получить текст лицензии, посетите веб-сайт <http://creativecommons.org/licenses/by-nd/3.0/>
или отправьте письмо в Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

Introduction

In heterogeneous computing systems, task scheduling and resource allocation present significant challenges. These systems integrate various types of processors, such as CPUs, GPUs, and FPGAs, to execute tasks efficiently and optimize resource utilization. Workflows in these systems are often modeled as directed acyclic graphs (DAGs), which clearly represent task dependencies, where nodes denote individual tasks and edges reflect precedence and data transfer costs between tasks.

Numerous studies examine online and dynamic scheduling, often using the terms interchangeably. Typically, papers specify the scheduling model and context, which depend on executor parameters, especially in fog computing environments [Alizadeh et al., 2020]. Task arrival attributes, whether random-based for non-DAG models [Alizadeh et al., 2020] or involving random arrival times for known DAGs [Quintin, Wagner, 2012], also define this. When DAGs aren't fully visible, schedulers have limited visibility and cannot use heuristics like the Critical-Path method like CPOP and its successors like HEFT [Ghose, Dey, 2022; Grandl et al., 2014].

Significant research on online and dynamic scheduling for DAG-modeled applications has been contributed by major companies. Microsoft's research team has focused on resource management [Mao et al., 2016] and cluster scheduling [Grandl et al., 2014]. Cloud providers, including Microsoft [Jalaparti et al., 2015; Yan et al., 2016], Amazon [Durillo, Prodan, 2013], Google [Soualhia, Khomh, Tahar, 2015], and IBM [Feitelson, 1994], have published findings on process optimization and scheduling. Distributed application frameworks like Spark [Duan, Wang, Wu, 2020] and Flink [Li et al., 2020] utilize DAGs to depict applications and employ various online scheduling methodologies.

This paper introduces the Multi-Agent Local Voting Protocol (MLVP), a novel approach to DAG scheduling in heterogeneous computing environments. The MLVP leverages a multi-agent framework where executors, acting as agents, use a local voting protocol to dynamically balance loads across the system and in turn reducing makespan. The protocol calculates an aggregated DAG metric for each task-executor pair, optimized via a genetic algorithm, to ensure efficient task assignment. This method addresses the limitations of traditional scheduling algorithms by accommodating the dynamic and online nature of real-world computing tasks.

The remainder of this paper is organized as follows. Section "Related works" provides an overview of the online scheduling problem and dynamic DAG scheduling approaches. Section "Online DAG Scheduling Problem" describes distributed heterogeneous computing system, DAG model and performance criteria. Section "DAG Scheduling Algorithms" details the proposed MLVP-based scheduling algorithm and its architecture. Section "Numerical simulation" evaluates the MLVP-based scheduling algorithm performance under different workloads.

Related works

In online scheduling in general, review by [Alizadeh et al., 2020] identifies popular scheduling algorithms: MCT, Min-Min, and Max-Min used in fog computing for task allocation based on task size and completion times. Review by [Khallouli, Huang, 2022] assesses scheduling methods in YARN, Borg, and Kubernetes, and advocates for "smarter" approaches, like reinforcement learning and decision trees to refine multi-objective optimization in cloud resource allocation, addressing the complexities of heterogeneous jobs and resources.

Online DAG scheduling often deals with incomplete knowledge of tasks and communications. To address this, [Quintin, Wagner, 2012] introduces WSCOM, a variation of the work-stealing algorithm, for scheduling file execution tasks. [Grandl et al., 2014] tackles online task scheduling in data-parallel clusters that require diverse resources such as CPU, memory, disk, and network. They propose the Tetris scheduling method, which strives to optimally assign tasks based on comprehensive resource needs.

The consensus approach is increasingly popular for addressing various practical issues, such as cooperative control in multi-vehicle networks [Granichin et al., 2012; Ren, Beard, Atkins, 2007], distributed control of robotic networks [Bullo, Cortés, Martínez, 2009], the flocking problem [Yu, Chen, Cao, 2010; Virágh et al., 2013], and optimizing sensor networks [Kar, Moura, 2010]. It holds promise for load balancing in computer, production, transport, logistics, and service networks by framing the problem as a consensus issue among network nodes [Amelina et al., 2015]. The study in [Vergados et al., 2018] examines optimal task redistribution in stochastic networks with variable priorities by determining an ideal step-size for a consensus-type protocol.

Online DAG Scheduling Problem

This paper addresses the challenge of scheduling a single DAG within a distributed heterogeneous system, specifically when executors have visibility limited to only the initial layers of the working front. The scheduling model is composed of four distinct parts.

Directed acyclic graph $G = (C, E)$, where $C = 1, 2, 3 \dots, c$ is a set of nodes and E is the set of edges, where:

- Edge $(v, s) \in E$, $v, s \in C$ denotes the precedence constraint such that node j must wait until task v finishes its execution.
- Cost of communication $b_{v,s}$ between nodes v and s , $(v, s) \in E$ should be taken into account if node i and s assigned to different executors, otherwise there is no cost of communication. If node s has several parents (v, \dots, k) , $k \in O$ that were performed on executors other than the executor assigned to node v , then the cost of communication is taken into account from each parent $b_{v,s} + \dots + b_{k,s}$.
- Each node v has a type that denotes on what type of executor this node should be executed.
- Set of immediate predecessors of node i in a DAG is expressed as $\text{pred}(v)$. A node without any predecessor is called an entry node. Multiple entry nodes may exist in a DAG.
- Set of immediate successors of node v is expressed as $\text{succ}(v)$. A node without any successor is called an exit node v_{exit} . There may be multiple exit nodes in a DAG.

Distributed heterogeneous computing system is represented by the graph $D = (N, L)$, where:

- Graph has $N = 1, 2, 3, \dots, n$ executors that work in parallel, L bidirectional links exist between any two executors.
- At every moment of time t each executor performs his task or chooses a new one from the DAG ready to perform nodes $H^t = 1, 2, 3, \dots, h^t$, $H^t \subset O$.
- Each node i , $i \in H^t$ of the DAG has $R^i = 1, 2, 3, \dots, r^i$ child nodes and each node j^m , $j^m \in R^i$ could have not only node i as a mother node.
- Executor m , $m \in N$ chooses the new node to perform in the state $x^m(t)$. The state of executor m is the evaluation of the probability of choosing a node P_i^m for each DAG node i , at each moment of time.
- τ_i^m is the execution time for executor m to process the node i .
- All executors are divided in equal shares according to the types of DAG nodes.

- Each executor m has a type that denotes what type of tasks can be defined for this executor. In order to assign a node i to an executor m , it is necessary that the type of node i and the type of executor m match.
- Executor m receives a penalty $w_{v,m}$ for each parent node of v , that is not performed by node m .

Online sliding informational window: at each given moment the scheduler observes only ready-to-execute nodes and their immediate successors within the informational window with length 2. This limits the amount of information that can be utilized for scheduling and makes problem online, i. e., scheduler cannot utilize information about nature of task arrival. So, on the one hand, in comparison with the static scheduling there is significantly less information, which makes it impossible to utilize classic DAG scheduling algorithms like CPOP or HEFT, that requires full information about DAG. On the other hand, limited information about DAG topology and parameters of not yet observed tasks combines limitations described in various online DAG scheduling papers (see “Introduction”) and makes the problem more difficult. If compared with purely online scheduling, there is more information available such as dependencies and parameters of tasks inside the sliding informational window.

Performance criterion for scheduling: before presenting the final scheduling objective function, we first define the Makespan, Earliest Start Time (EST), Earliest Finish Time (EFT) attributes.

- Makespan is the finish time of the last node in the scheduled DAG. It is defined by $\text{makespan} = \max\{AFT(v_{\text{exit}})\}$ where $AFT(v_{\text{exit}})$ is the actual finish time (AFT) of exit node v_{exit} . In the case where there are multiple exit nodes, the makespan is the maximum AFT of all exit nodes.
- $EST(v, m)$ denotes the earliest start time of node v on executor m and it is defined as

$$EST(v, m) = \max \left\{ T_{\text{Ava}}(m), \max_{\{v \in \text{pred}(v)\}} \{AFT(v) + b_{v,s}\} \right\}.$$

$T_{\text{Ava}}(m)$ is the earliest ready time of executor m .

- $EFT(v, m)$ denotes the earliest finish time of node v on executor m and is defined as

$$EFT(v, m) = EST(v, m) + w_{v,m}.$$

The objective function of the DAG scheduling is to determine the assignment policies of an application’s node to heterogeneous executors so that the makespan is minimized.

Graph generation

We utilize the open-source project DAGGEN to generate random DAGs [DAGGEN, 2022]. It is a popular tool for generating DAGs to evaluate the performance of scheduling heuristics and has been used in studies on DONF [Lin et al., 2019], CPOP [Sih, Lee, 1993], HEFT [Topcuoglu, Hariri, Wu, 1999], and PETS [Ilavarasan, Thambidurai, 2007]. DAGGEN creates random, synthetic task graphs for simulation purposes, with the process varying according to the configured parameters:

1. Generate the tasks according to -n, -fat, -regularity.
2. Generate the dependencies according to -density, -jump, -ccr.
3. Add transfer costs, these costs derive from the size of the data handled by the initiator of the transfer.

Table 1 presents a description of the graph generation parameters.

Table 1. DAG generation parameters

Parameter	Description
n	Number of computation nodes in the DAG (i. e., application “tasks”). Values: 3000, 6000, 12 000, 24 000.
fat	Width of the DAG, the maximum number of tasks that can be executed concurrently. A small value will lead to a thin DAG (e. g., chain) with low task parallelism, while a large value induces a fat DAG (e. g., fork-join) with high parallelism. Values: 0.2, 0.5.
$density$	Determines the number of dependencies between tasks of two consecutive DAG levels. Values: 0.1, 0.4, 0.8.
$regularity$	Regularity of the distribution of tasks between the different levels of the DAG. Values: 0.2, 0.8.
$jump$	Maximum number of levels spanned by inter-task communications. This allows generating DAGs with execution paths of different lengths. Values: 2, 4.
ccr	The ratio of the sum of edge weights to the sum of node weights. Values: 0.2, 0.8.

DAG Scheduling Algorithms

The Multi-Agent Local Voting Protocol (MLVP) algorithm is designed to optimize task scheduling for directed acyclic graphs (DAGs) in heterogeneous computing environments. The algorithm achieves efficient task distribution by calculating an aggregated metric for each executor-DAG node pair and assigning probability of a given node execution on this executor. Aggregated metric is derived from three key factors:

1. Incoming Connections: The ratio of the total incoming connections from unfinished DAG nodes to their child nodes.
2. Available Nodes: The number of DAG nodes currently available for execution compared to nodes of the same type.
3. Executor Performance: The performance capability of the executor in processing the task.

A genetic algorithm is employed to optimize the coefficients of these factors in a linear combination, forming the aggregated DAG metric. This metric serves as a basis for determining the probability of assigning a DAG node to a specific executor. By optimizing these coefficients, the genetic algorithm ensures that the metric accurately reflects the importance of each factor, leading to better task prioritization and scheduling efficiency.

During the Local Voting Procedure, the genetic algorithm further refines the assignment probabilities by dynamically adjusting them to balance the workload across executors. This iterative optimization process helps the system adapt to varying task and executor conditions, thereby reducing makespan and enhancing overall performance.

Consensus plays a crucial role in the MLVP algorithm by coordinating task distribution among executors. The local voting protocol (LVP) facilitates consensus by enabling executors to share their current task assignment probabilities and states with neighboring executors. The protocol iteratively adjusts these probabilities to ensure a balanced workload distribution, which is achieved when the difference in task assignment states between any executor and its neighbors falls below a specified threshold (ϵ), which in this case was 0.05. The algorithm’s pseudocode is provided as Algorithm 1.

Evaluation of probability for nodes

At each moment of time t the executors choose $\frac{|H^t|}{0.8n}$ nodes for estimation at time moment t randomly. Coefficient 0.8 is selected to limit the number of ready-to-execute tasks on extremely large (24 000 nodes) and wide (with large “fat” parameter) to reduce scheduling complexity. To synchronize

Algorithm 1. Multi-Agent Local Voting Protocol (MLVP)

```

1: Input: Set of executors  $E$ , Set of DAG nodes  $N$ 
2: Output: Assignment probabilities for each DAG node to each executor
3: Initialize assignmentProbabilities to an empty data structure
4: for each executor  $e \in E$  do
5:   for each node  $n \in N$  do
6:     if  $n$  is not finished and  $n$  is ready for execution then
7:       incomingConnections  $\leftarrow$  Sum of weights of all incoming edges to  $n$ 
8:       availableNodes  $\leftarrow$  Count of nodes ready for execution of the same type as  $n$ 
9:       executorPerformance  $\leftarrow$  Performance metric of executor  $e$ 
10:      aggregatedMetric  $\leftarrow$  LinearCombination(incomingConnections, availableNodes,
executorPerformance)
11:      assignmentProbabilities[ $e$ ][ $n$ ]  $\leftarrow$  CalculateProbability(aggregatedMetric)
12:    end if
13:  end for
14: end for
15: Find coefficients for LinearCombination function based on historical data
16: for each node  $n \in N$  do
17:   if  $n$  is ready for execution then
18:     for each executor  $e \in E$  do
19:       Tune assignmentProbabilities[ $e$ ][ $n$ ] using genetic algorithm to balance workload
20:     end for
21:   end if
22: end for
23: return assignmentProbabilities

```

the data of all nodes the local voting protocol (LVP) is used. The LVP is necessary to reach a consensus on the state between the executors.

For each node i at time moment t the following values are calculated by the algorithm:

1. Ratio of the sum of all incoming connections from not finished DAG node into child's nodes of node i to the number of child nodes of node i :

$$W_t = \sum_{j^i \in R^i} \frac{inp_t(j^i)}{R^i},$$

where R^i is the number of child nodes of i node, $inp_t(j^i)$ is the number of **not finished** mother nodes of j_i node.

2. Sum of the ratio of DAG nodes number currently available for execution to the number of the same as j^i types:

$$Q_i = \sum_{j^i \in R^i} \frac{|H^t|}{size_{j^i}(|H^t|)},$$

where $size_{j^i}(|H^t|)$ is the number of DAG nodes ready to execute with the same type as node j^i .

3. Performance of executor m is taken into account. For executor m the rate of performing the node i is compared with executor $e \in N$ which is the fastest executor for task i :

$$Z_i = \min_{e \in N} (exTime_e(i)) - exTime_m(i),$$

where $exTime$ is the time when the executor needs to finish the node i .

4. Sum of previous elements with coefficients:

$$result_i = wW_i + qQ_i + zZ_i.$$

Coefficients w, u, z were defined using the genetic algorithm [Carson, 2017] for 2 types of graph topologies. Random population of 100 elements is created, after testing the top 10 go to the next stage, we also take 40 crossed from the top 10, 39 crossed with additional mutations, 10 copies of the top 10 with mutations and 1 random DNA, with the target $makespan \rightarrow \min$ and stopping condition if makespan improvement with previous is less than 0.1 %.

5. The probability of choosing the executor m to perform the node i is

$$P_i^m = \frac{result_i}{\sum_{j \in nodesToExecute} result_j}.$$

Local Voting Protocol

The state of executor m is the evaluation of the probability of choosing a node P_i^m for each DAG node i , at each moment of time. The dynamics of the state of the node i has the form:

$$x_{t+1}^m = x_t^m + f_t^m + u_t^m$$

where f_t^m is the number of new nodes of DAG that were opened for execution.

At time instant t executor m sends to the $\frac{n}{4}$ nodes its own load x_m^t . We assume that to form the control strategy each executor m has observations about its other agents state with noise and delay:

$$y_t^{m,c} = x_{t-d_t^{m,c}}^j + w_t^{m,c}, \quad c \in N,$$

where $d_t^{m,c}$ is delay of transfer data, and $w_t^{m,c}$ is a noise.

The message exchange is undergone only once per system cycle. Applying the local voting algorithm, we obtain a parameter characterizing the state of executor m relative to its neighbors:

$$u_t^m = \gamma \sum_{c \in N} (y_t^{m,c} - x_t^{m,m})$$

where $\gamma > 0$ is a step-size, which represents the sensitivity of the algorithm to the difference between agents states. The consensus is achieved when the state of the last node (the node to which the message was intended) differs from its neighbors by no more than a specified one ϵ -consensus.

Algorithms for comparison

In this section we consider a detailed description of the graph scheduling algorithm that we use to compare with the proposed algorithm. The Degree-of-Node-First algorithm is selected in accordance with the review article for graph scheduling in heterogeneous environments [Suman, Kumar, 2018; Murad et al., 2022; Li et al., 2010; Kanani, Maniyar, Mohammad, 2015; Hayatunnufus, Riasetiawan, Ashari, 2020]. Selection was based on the following considerations: select state-of-the-art DAG scheduling (DONF) and online scheduling heuristics (Optimized Min-Min and Max-Min), as well as widely used FIFO. Various Machine Learning and Reinforcement Learning-based algorithms have been considered in the preparation of this paper, but most of the modern works either assume that the whole graph is known in advance (as opposed to the online sliding informational window in our case), and/or the algorithm is not designed for heterogeneous types of executors [Zhang et al., 2022; Jeon et al., 2023; Koslovski, Pereira, Albuquerque, 2024; Roeder, Pimentel, Grelck, 2023]. While classical and

graph heuristics are quite easy to transfer to such a setting, ML/RL based approaches usually cannot be transferred, and require retraining/modification. The proper transfer of such models from one task to another can be considered as a separate research task, which is beyond the scope of this article. Another significant problem is the dimensionality of the graphs on which comparisons are made in articles using ML/RL approaches. Usually papers, with even remotely comparable settings [Zhang et al., 2022; Jeon et al., 2023; Koslovski, Pereira, Albuquerque, 2024; Roeder, Pimentel, Grelck, 2023] usually consider DAGs with no more than 2000 nodes, and even less, while our approach is designed for large-scale DAGs with more than 3000 nodes.

- **Degree of Node First (DONF) algorithm** [Lin et al., 2019]. It is feasible to maintain a higher parallelism during the scheduling process in order to make full use of heterogeneous system resources. Thus, the chosen node should have the property of enlarging parallelism as much as possible. The Degree-of-node scheduling procedure can be shortly described as nodes with larger out-degree should be scheduled earlier. The weighted out-degree (WOD) of node v_i is defined by:

$$\text{WOD}(v_i) = \sum_{v \in \text{succ}(s)} \frac{1}{\text{ID}(s)}, \quad (1)$$

where $\text{ID}(s)$ is the in-degree of node s . Each executor in its turn selects best node by DONF algorithm until all nodes are executed. The DONF algorithm has a computational complexity $O(v \times o \times p)$ where v is number of nodes in DAG, p is the number of executors, and o is the maximum out-degree of any node in DAG.

- **Optimized Min-Min task scheduling algorithm** [Murad et al., 2022]. Min-Min does the opposite by choosing and assigning small tasks in resources that can perform the task with minimum execution time. The objective of the Min-Min algorithm is to ensure, firstly, that all tasks with a minimum completion time are completed. The Min-Min has the following steps:
 1. For all nodes of DAG and all executors the time of performance is calculated.
 2. Each executor in its turn selects and executes the node with the lowest execution time.

The Optimized Min-Min algorithm has a computational complexity of $O(v^2)$, where v is the number of nodes in DAG.

- **Max-Min Task scheduling Algorithm** [Kanani, Maniyar, Mohammad, 2015]. The Max-Min is quite similar with a previous algorithm, the main difference is in the node selection criterion – in this case the executor should select the biggest node by weight. The Max-Min has the following steps:
 1. For all nodes of DAG performance weight is calculated.
 2. Each executor in its turn selects and executes the node with the biggest performance weight.

The Min-Max algorithm has a computational complexity of $O(v^2)$, where v is the number of nodes in DAG.

- **FIFO** [Hayatunnufus, Riasetiawan, Ashari, 2020]. This is the default Scheduler. The tasks are placed in a queue and the tasks are performed in their submission order. In this method, once the job is scheduled, no intervention is allowed. The FIFO algorithm has the following steps:
 1. Each node receives its place in the queue when ready for execution.
 2. Each executor in its turn selects and executes the first node in the queue.

The FIFO algorithm has a computational complexity of $O(v)$, where v is the number of nodes in DAG.

MLVP differs mainly in the stage of executor selection, where executors need to pass information between each other in order to facilitate consensus and information exchange occurs only one once per scheduling cycle. With MVLP complexity $O(k \times o \times p^2)$, where k is the maximum number of ready-to-execute nodes in DAG, o is the maximum out-degree of any node in DAG, and p is the number of executors. We consider k in complexity estimation primarily because each executor selects tasks (nodes) to be scheduled after consensus is reached simultaneously, i. e. MLVP processes a batch of nodes at the same time. The worst case scenario can be considered in two aspects: if DAG is akin to a circle, i. e., all nodes are ready to be executed from the start, that will yield complexity of $O(v \times o \times p^2)$, and if the maximum out-degree of any node is as large as $v - 1$ then the complexity will become $O(v^2 \times p^2)$, which is worse than the DONF worst-case scenario of $O(v^2 \times p)$. Since coefficients selected by the genetic algorithm are fixed in the online scheduling phase, we consider only online stage complexity. Given the size of graphs and the diversity of their topologies, finding suitable coefficients took roughly 2 days (i. e., 48 hours) of computations combined for all DAG types. Note that this figure is not relevant with respect to the length of time it takes to construct the schedule in an online setting. Also, genetic algorithms in and of itself are not particularly important: makespan in the given task is a nondifferentiable function, and any optimization approach that is suitable for this should be sufficient.

We do not conduct evaluation in terms of concrete utilization of memory, computational power and execution time since this is not common in the DAG scheduling problem in general (DONF [Lin et al., 2019], CPOP/HEFT [Ghose, Dey, 2022; Grandl et al., 2014], PEFT [Arabnejad, Barbosa, 2014]) and specifically in online setting (Optimized Min-Min [Murad et al., 2022], Min-Max [Kanani, Maniyar, Mohammad, 2015], FIFO [Hayatunnufus, Riasetiawan, Ashari, 2020]). Resource utilization is significantly influenced by the specific implementation of the scheduling environment and specific algorithms. Such information would not be representative of the algorithms compared in a general sense and is not provided in the literature listed earlier, and thus is beyond the scope of the comparison.

Numerical simulation

In this section, the design of an experiment is explained and the performance of the proposed MLVP-based algorithm is evaluated. The hardware used for computations:

- Server 1: CPU AMD 9 5950x 32 threads; GPU RTX 3080 ti; RAM 128 GB.
- Server 2: CPU AMD Ryzen Threadripper 3990X 128 threads; GPU RTX 3080 ti; RAM 192 GB.

Simulation environment

The proposed Local Voting based algorithm is tested to measure its generalization and robustness depending on different system configurations. The working environment is determined by different heterogeneous configurations from three to twenty four executors, 4 configurations in total. The types of executors are evenly distributed (equally for executors of types 1, 2, and 3), and their powers were set randomly from a range of 10 to 300 GFlops. Each workspace is intended for testing on a fixed DAG dimension. The details of the configuration test cases are shown in Table 2.

In the simulation, the executors inside a computer node are “fully connected”: links between different executors do not obstruct each other. Links between computer nodes also have similar properties, however, later data transmission starts only when all the former ones finish.

For testing and training algorithms, a real-time graph traversal C++ simulation environment was developed (see Fig. 1). The logic of the simulation is as follows: after starting the calculation, executors,

Table 2. Description of test cases for Large-Scale DAGs

Workspace ID	Total number of executors	Executors of each type	Dimension DAGs
1	3	1	3000
2	6	3	6000
3	12	4	12 000
4	24	8	24 000

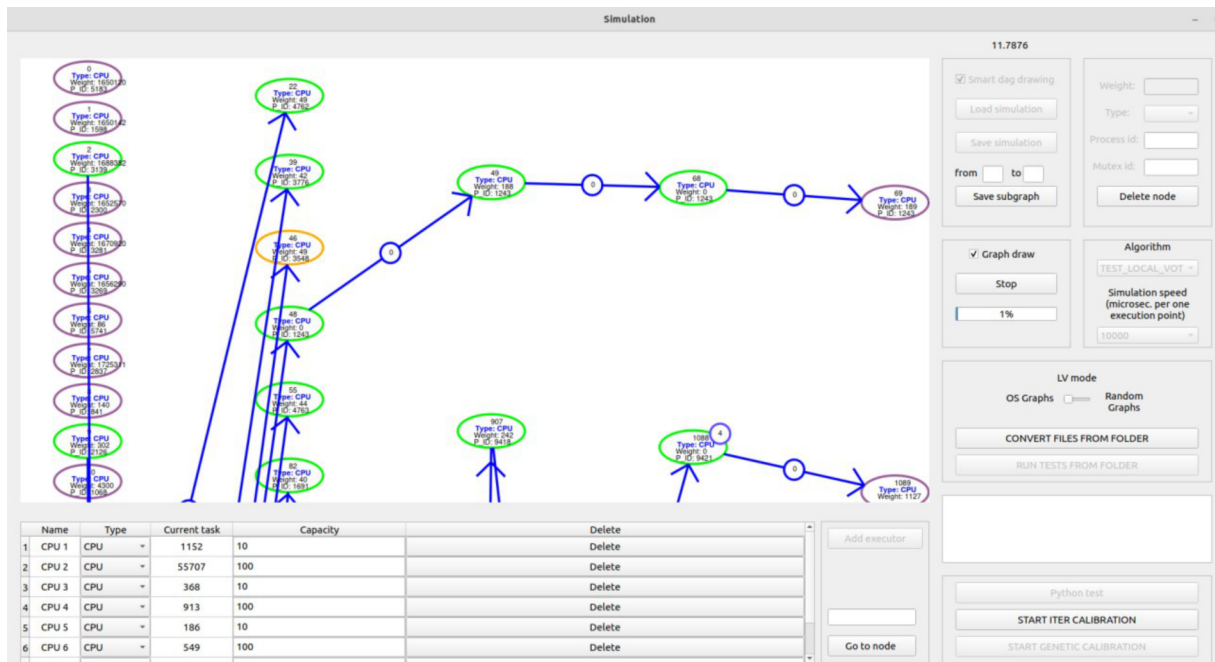


Figure 1. An example of a C++ simulation environment

working in parallel and exchanging service data, choose a node for execution, and then work on it for the time period, depending on node weight and executor power. After execution of a node the operation is repeated, and so on until the moment when all nodes in the DAG are executed.

For training and evaluation of the proposed MLVP-based algorithm, a batch of random DAG are used (described in Section “Graph generation”). At fixed dimensions there exist 48 topologies. For training, 3 graphs are generated for each topology for training (totally 144 for each dimension 3000, 6000, 12 000, 24 000) and 10 graphs for evaluation (totally 480 for each dimension 3000, 6000, 12 000, 24 000). A total of 1920 DAGs were tested. An example of a DAG shown in Fig. 2. A comparative analysis is presented in the section “Results”.

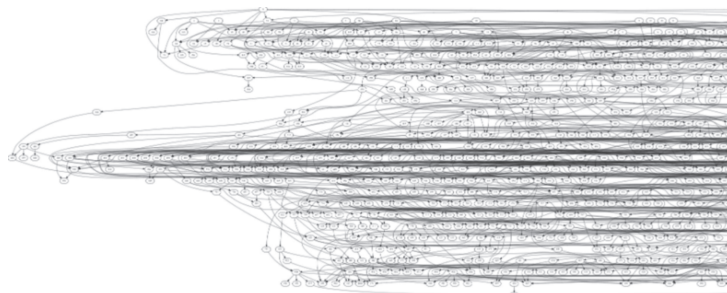


Figure 2. Examples of large-scale DAG

Results

We present the comparison results between the MLVP-based algorithm and the state-of-the-art DAG scheduling algorithms DONF, FIFO, Min-Min, Min-Max for different workspaces.

The comparison was carried out in accordance with the following metric for each tested graph:

$$p = \frac{(makespan_{MLVP} - makespan_{other}) \cdot 100 \%}{makespan_{other}}, \quad (2)$$

where $makespan_{other}$ is the finish time of the last node in the scheduled DAG using one of the algorithms from the list above, $makespan_{MLVP}$ is the finish time of the last node in the scheduled DAG using the MLVP-based algorithm. The obtained values are summarized and averaged for each workspace and DAG dimension DAG. The proposed MLVP-based algorithm is tested on large-scale DAGs and is compared to the state-of-the-art algorithms, to determine its generalization and robustness depending on dimensions DAGs and different system configurations.

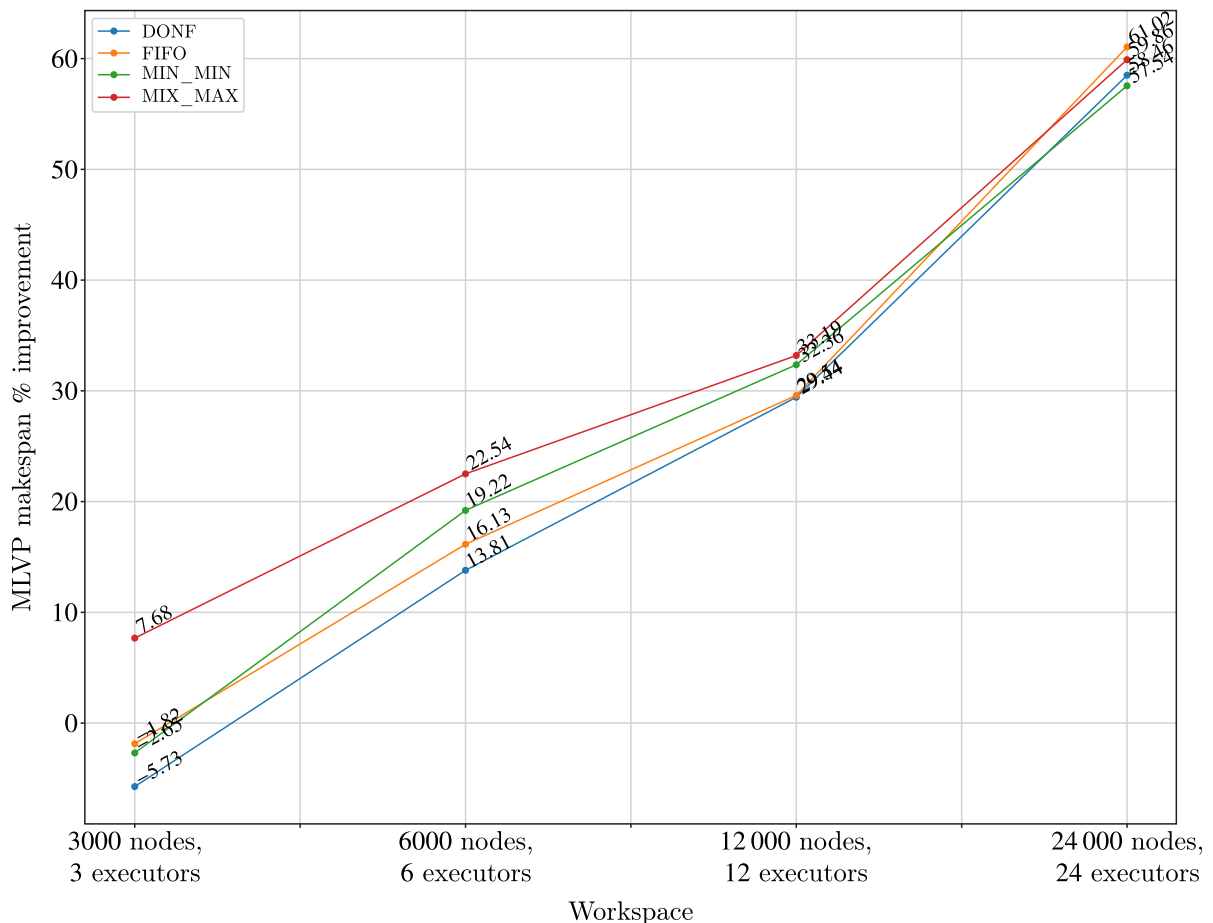


Figure 3. Percentage improvements MLVP-based approach compared to other algorithms

Figure 3 shows a comparison of the MLVP-based approach compared to DONF, FIFO, Min-Min, Min-Max algorithms for various test cases, where X-Axis is test cases and Y-Axis is the percentage of how much the MLVP-based exceeds/inferior compared with state-of-the-art algorithms. On average, the proposed approach outperforms the DONF by 23.99 %, FIFO by 26.21 %, Min-Min 26.61 %, Min-Max 30.81 % on all test cases.

Figure 4 presents box plots illustrating the percentage improvements in makespans of the MLVP algorithm relative to various scheduling algorithms—DONF, FIFO, Min-Min, and Min-Max—across different workspaces and DAG dimensions. Each workspace, differentiated by executor counts of 3, 6, 12, and 24, is analyzed separately to delineate the algorithm’s performance under varying computational resources.

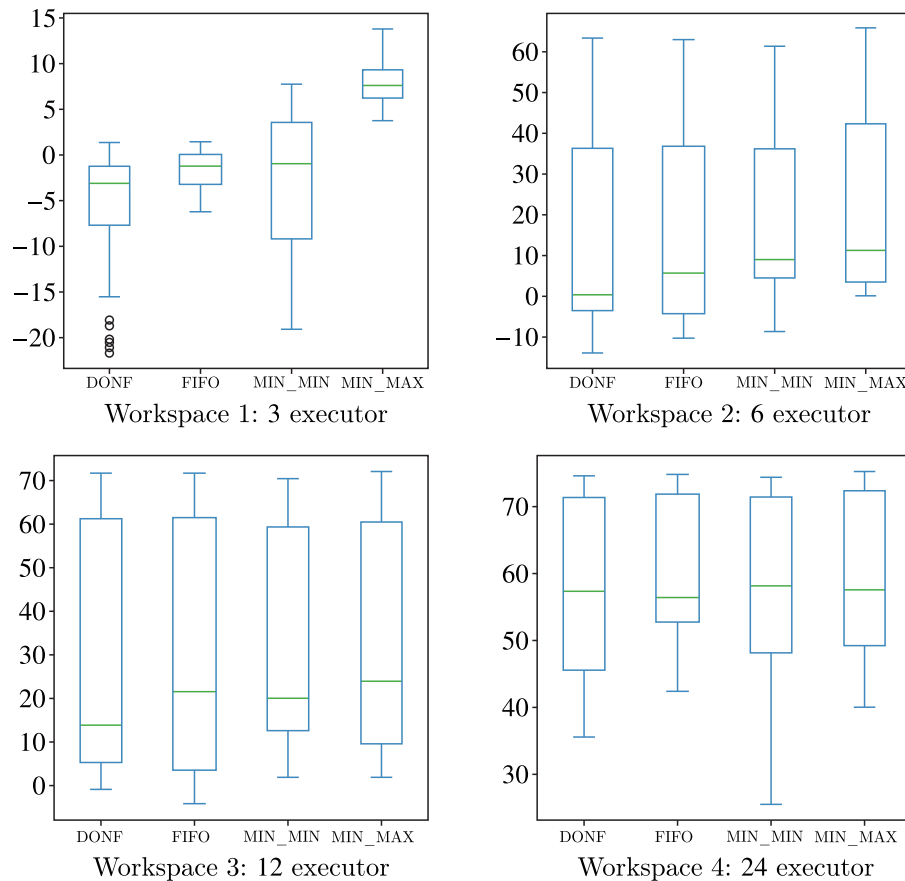


Figure 4. Box plot percentage improvements in makespans MLVP compared to DONF, FIFO, Min-Min,Min-Max for each workspace and each DAG dimension

In the first workspace with three executors, the MLVP algorithm exhibits substantial variation in performance relative to the reference algorithms, with a wider interquartile range. Conversely, in workspaces with six, twelve, and twenty-four executors, the MLVP algorithm demonstrates more consistent improvements, as indicated by the narrower interquartile ranges.

The observing the median values, it is evident that the MLVP algorithm tends to offer improvements in makespan across all workspaces, albeit with variations in the magnitude of improvement. Particularly, the algorithm shows consistent positive enhancements compared to the FIFO, Min-Min, and Min-Max algorithms across all executor counts. In the three-executor workspace, MLVP presents around a 10% improvement over DONF and FIFO and approximately a 5% improvement over Min-Min and Min-Max. In the six-executor workspace, MLVP shows nearly 20% improvement across all compared algorithms. In larger workspaces with twelve and twenty-four executors, MLVP consistently exhibits around a 60% improvement over the reference algorithms.

The MLVP algorithm demonstrates a promising capacity to improve makespan in various workspaces across different topologies and against multiple scheduling algorithms. It exhibits consistency and adaptability across different numbers of executions, underlining its potential to scale

effectively in heterogeneous computing environments (Fig. 5). Future work might focus on exploring the robustness and scalability of MLVP in more diverse and complex scheduling scenarios to further validate its applicability and performance.

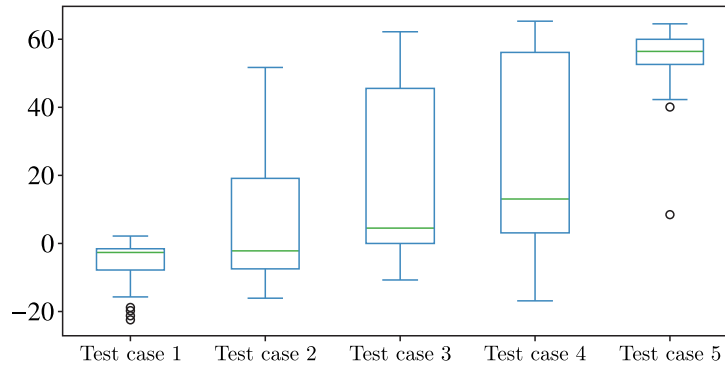


Figure 5. Box plot with makespan percentage difference comparisons between MLVP and DONF for each workspace

Conclusion

This paper presents the Multi-Agent Local Voting Protocol (MLVP), a novel approach to scheduling tasks in directed acyclic graphs (DAGs) within heterogeneous computing environments. The MLVP effectively addresses the challenges of online DAG scheduling by dynamically balancing the load across multiple executors. This is achieved through a local voting protocol that calculates an aggregated metric for each executor-DAG node pair, leveraging factors such as incoming connections, available nodes, and executor performance. The optimization of combination of these metrics is performed using a genetic algorithm, allowing for probabilistic task assignment that improves makespan.

The MLVP's innovative approach demonstrates significant improvements over traditional scheduling algorithms, with simulations showing up to a 70% reduction in makespan on specific graph topologies and an average performance gain of 23.99% over the Degree of Node First (DONF) algorithm. These results highlight the MLVP's scalability and effectiveness in various system configurations, making it a promising solution for real-world applications in cloud computing, data mining, and beyond.

Looking forward, future research could explore the extension of the MLVP to support multi-dimensional executor requirements and handling dynamic conditions such as fluctuating numbers of executors. Additionally, further investigation into optimizing resource reservation based on anticipated critical system tasks could enhance the algorithm's adaptability and efficiency.

References

- Alizadeh M.R., Khajehvand V., Rahmani A.M., Akbari E.* Task scheduling approaches in fog computing: A systematic review // *International Journal of Communication Systems*. — 2020. — Vol. 33, No. 16. — P. e4583.
- Amelina N., Fradkov A., Jiang Y., Vergados D.* Approximate consensus in stochastic networks with application to load balancing // *IEEE Transactions on Information Theory*. — 2015. — Vol. 61, No. 4. — P. 1739–1752.

- Arabnejad H., Barbosa J. G.* List scheduling algorithm for heterogeneous systems by an optimistic cost table // *IEEE Transactions on Parallel and Distributed Systems*. — 2014. — Vol. 25, No. 3. — P. 682–694.
- Bullo F., Cortés J., Martínez S.* Distributed control of robotic networks: A mathematical approach to motion coordination algorithms / *Princeton Series in Applied Mathematics*. — Princeton University Press, 2009. — Vol. 27. — 336 p.
- Carson J.* Genetic algorithms: Advances in research and applications // *Genetic Algorithms: Advances in Research and Applications*. — 2017. — P. 1–117.
- Duan Y., Wang N., Wu J.* Reducing makespans of DAG scheduling through interleaving overlapping resource utilization / *2020 IEEE 17th International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*. — 2020. — P. 392–400.
- Durillo J., Prodan R.* Multi-objective workflow scheduling in Amazon EC2 // *Cluster Computing*. — 2013. — Vol. 17. — P. 169–189.
- Feitelson D. G.* A survey of scheduling in multiprogrammed parallel systems // *IBM T.J. Watson Research Center research report*. — 1994.
- Ghose A., Dey S.* FGFS: Feature guided frontier scheduling for SIMT DAGs // *J. Supercomput.* — 2022. — Vol. 78, No. 9. — P. 11702–11743.
- Grandl R., Ananthanarayanan G., Kandula S., Rao S., Akella A.* Multi-resource packing for cluster schedulers // *SIGCOMM Comput. Commun. Rev.* — 2014. — Vol. 44, No. 4. — P. 455–466.
- Granichin O., Skobelev P., Lada A., Mayorov I., Tsarev A.* Comparing adaptive and non-adaptive models of cargo transportation in multi-agent system for real time truck scheduling // *International Conference on Evolutionary Computation Theory and Applications*. — 2012. — Vol. 2. — P. 282–285.
- Hayatunnufus H., Riasetiawan M., Ashari A.* Performance analysis of FIFO and Round Robin scheduling process algorithm in IoT operating system for collecting landslide data / *2020 International Conference on Data Science, Artificial Intelligence, and Business Analytics (DATABIA)*. — 2020. — P. 63–68.
- Ilavarasan E., Thambidurai P.* Low complexity performance effective task scheduling algorithm for heterogeneous computing environments // *Journal of Computer Science*. — 2007. — Vol. 3.
- Jalaparti V., Bodik P., Menache I., Rao S., Makarychev K., Caesar M.* Network-aware scheduling for data-parallel jobs: plan when you can // *SIGCOMM Comput. Commun. Rev.* — 2015. — Vol. 45, No. 4. — P. 407–420.
- Jeon W., Gagrani M., Bartan B., Zeng W. W., Teague H., Zappi P., Lott C.* Neural DAG scheduling via one-shot priority sampling / *The Eleventh International Conference on Learning Representations*. — 2023.
- Kanani B., Maniyar B., Mohammad S.* Review on Max-Min task scheduling algorithm for cloud computing // *SSRN Electronic Journal*. — 2015. — Vol. 2. — P. 781–784.
- Kar S., Moura J.* Distributed consensus algorithms in sensor networks: quantized data and random link failures // *IEEE Transactions on Signal Processing*. — 2010. — Vol. 58, No. 3. — P. 1383–1400.
- Khallouli W., Huang J.* Cluster resource scheduling in cloud computing: literature review and research challenges // *The Journal of Supercomputing*. — 2022. — Vol. 78.
- Koslovski G. P., Pereira K., Albuquerque P. R.* DAG-based workflows scheduling using Actor-Critic Deep Reinforcement Learning // *Future Generation Computer Systems*. — 2024. — Vol. 150. — P. 354–363.
- Li A., Yang X., Kandula S., Zhang M.* CloudCmp: Comparing public cloud providers / *Proceedings of the ACM SIGCOMM Internet Measurement Conference, IMC*. — 2010. — P. 1–14.

- Li Z., Yu J., Bian C., Pu Y., Wang Y., Zhang Y., Guo B.* Flink-ER: An elastic resource-scheduling strategy for processing fluctuating mobile stream data on flink // *Mobile Information Systems*. — 2020. — P. 1–17.
- Lin H., Li M., Jia C., Liu J., An H.* Degree-of-node task scheduling of fine-grained parallel programs on heterogeneous systems // *Journal of Computer Science and Technology*. — 2019. — Vol. 34. — P. 1096–1108.
- Mao H., Alizadeh M., Menache I., Kandula S.* Resource management with deep reinforcement learning / *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*. — 2016. — P. 50–56.
- Murad S., Abdal R., Salih N., Alsandi A., Faraj R., Muhammed A., Derahman M., Ahmed R.* Optimized min-min task scheduling algorithm for scientific workflows in a cloud environment // *Journal of Theoretical and Applied Information Technology*. — 2022. — Vol. 100. — P. 480–506.
- Open Source Software. Daggen. — <https://github.com/frs69wq/daggen> (accessed: 22 October 2024) // *github*. — 2024.
- Quintin J., Wagner F.* WSCOM: online task scheduling with data transfers / *2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012)*. — 2012. — P. 344–351.
- Ren W., Beard R., Atkins E.* Information consensus in multivehicle cooperative control // *IEEE Control Systems Magazine*. — 2007. — Vol. 27, No. 2. — P. 71–82.
- Roeder J., Pimentel A. D., Grelck C.* GCN-based reinforcement learning approach for scheduling DAG applications / *IFIP International Conference on Artificial Intelligence Applications and Innovations*. — 2023. — P. 121–134.
- Sih G., Lee E.* A compile-time scheduling heuristic for interconnection-constrained heterogeneous processor architectures // *IEEE Transactions on Parallel and Distributed Systems*. — 1993. — Vol. 4, No. 2. — P. 175–187.
- Soualhia M., Khomh F., Tahar S.* Predicting scheduling failures in the cloud: A case study with Google clusters and hadoop on Amazon EMR / *2015 IEEE 17th International Conference on High Performance Computing and Communications, 2015 IEEE 7th International Symposium on Cyberspace Safety and Security, and 2015 IEEE 12th International Conference on Embedded Software and Systems*. — 2015. — P. 58–65.
- Suman C., Kumar G.* Analysis of process scheduling algorithm for multiprocessor system / *2018 7th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO)*. — 2018. — P. 564–569.
- Topcuoglu H., Hariri S., Wu M.* Task scheduling algorithms for heterogeneous processors / *Proceedings. Eighth Heterogeneous Computing Workshop (HCW'99)*. — 1999. — P. 3–14.
- Vergados D., Amelina N., Jiang Y., Kravlevska K., Granichin O.* Toward optimal distributed node scheduling in a multihop wireless network through local voting // *IEEE Transactions on Wireless Communications*. — 2018. — Vol. 17, No. 1. — P. 400–414.
- Virágh C., Vásárhelyi G., Tarcai N., Szörényi T., Somorjai G., Nepusz T., Vicsek T.* Flocking algorithm for autonomous flying robots // *Bioinspiration and Biomimetics*. — 2013. — Vol. 9.
- Yan Y., Gao Y., Chen Y., Guo Z., Chen B., Moscibroda T.* TR-Spark: transient computing for big data analytics / *Proceedings of the Seventh ACM Symposium on Cloud Computing*. — 2016.
- Yu W., Chen G., Cao M.* Distributed leader-follower flocking control for multi-agent dynamical systems with time-varying velocities // *Systems and Control Letters*. — 2010. — Vol. 59. — P. 543–552.
- Zhang Z., Hua Q.-S., Zhang X., Liao X.* DAG scheduling with communication delays based on graph convolutional neural network // *Wireless Communications and Mobile Computing*. — 2022. — No. 9013361.