

DOI: <https://doi.org/10.36910/6775-2524-0560-2021-43-28>

УДК 681.3.082.5

Арпентій Сергій Петрович, старший науковий співробітник

<https://orcid.org/0000-0003-3326-3942>

Український науково-дослідний інститут спеціальної техніки та судових експертиз Служби безпеки України.

ОСОБЛИВОСТІ ЗАСТОСУВАННЯ РОЗПОДІЛЕНИХ ОБЧИСЛЕНЬ ПРИ ОБРОБЦІ ПОТОКОВИХ ДАНИХ

Арпентій С. П. Особливості застосування розподілених обчислень при обробці поточкових даних. Проведено аналіз сучасних алгоритмів поточної обробки масивів цифрових даних та методик формалізації зазначених процедур з метою побудови відповідного математичного апарату. Побудовано узагальнену схемупотокової передачі масивів даних та схему апаратно-програмного комплексу хмарного сервісу. Вказано на особливості організації апаратно-програмного комплексу мережевого вузла відповідно до архітектури розподіленої інформаційної системи. Вказано на задачі, що мають бути вирішені з метою оптимізації зазначеної структури, зокрема задачу налаштування графіка обробки запитів відповідно до особливостей роботи загального комплексу та задачу налаштування алгоритмів паралельної обробки. Розроблено спеціалізовану математичну модель розподіленої інформаційної системи апаратно-програмного комплексу хмарного сервісу, що складається з центрального обчислювального вузла і периферійних обчислювальних вузлів, та включає у себе параметри відповідних компонент, функції відображення процедури розгортання завдання і функції маршрутизації потоку вхідних даних. На базі побудованої математичної моделі розроблено методику проведення розрахунку показників пропускної здатності і часу затримки при обробці запитів користувачів хмарного сервісу, що запропоновано розглядати як показники цільових функцій.

Ключові слова: хмарний сервіс, апаратно-програмний комплекс, розподілена інформаційна система, маршрутизація потоку даних, розгортання завдання, пропускна здатність каналу, графік завдань.

Арпентій С. П. Особенности применения распределенных вычислений при обработке поточковых данных. Проведен анализ современных алгоритмов потоковой обработки массивов цифровых данных и методик формализации данных процедур с целью построения соответствующего математического аппарата. Построена обобщенная схема потоковой передачи массивов данных и схема аппаратно-програмного комплекса облачного сервиса. Указаны особенности организации аппаратно-програмного комплекса сетевого узла согласно архитектуры распределенной информационной системы, а также задачи, которые должны быть решены в целях оптимизации указанной структуры. В частности рассмотрена задача настройки графика обработки запросов в соответствии с особенностями работы общего комплекса и задача настройки алгоритмов параллельной обработки. Разработана специализированная математическая модель распределенной информационной системы аппаратно-програмного комплекса облачного сервиса, которая состоит из центрального вычислительного узла и периферийных вычислительных узлов, и включает в себя параметры соответствующих компонент, функции отображения процедуры развертывания задачи и функции маршрутизации потока входных данных. На базе построенной математической модели разработана методика проведения расчета показателей пропускной способности и времени задержки при обработке запросов пользователей облачного сервиса, которые рассматриваются как показатели целевых функций.

Ключевые слова: облачный сервис, аппаратно-программный комплекс, распределенная информационная система, маршрутизация потока данных, развертывание задачи, пропускная способность канала, график задач.

Arpentii Sergii. Peculiarities of the application of distributed computing for processing streaming data. The analysis of modern algorithms for streaming processing of digital data arrays and methods for formalizing the procedures in order to build an appropriate mathematical apparatus is provided. The generalized scheme of streaming data arrays and the scheme of the hardware-software complex of the cloud service are built. The features of the organization of the hardware and software complex of the network node according to the architecture of the distributed information system, as well as the tasks that must be solved in order to optimize the specified structure are indicated. In particular, the problem of optimizing the schedule for processing requests in accordance with the peculiarities of the operation of the general complex and the problem of optimizing algorithms for parallel processing are considered. A specialized mathematical model of a distributed information system of a hardware-software complex of a cloud service has been developed. The model consists of a central computing node and peripheral computing nodes, and includes the parameters of the corresponding components, functions for displaying the task deployment procedure and routing functions for the input data flow. On the basis of the constructed mathematical model, the method for calculating the indicators of throughput and delay time when processing requests from users of the cloud service has been developed, which are considered as indicators of objective functions.

Key words: cloud service, hardware and software complex, distributed information system, routing of data flow, task deployment, channel capacity, task schedule.

Вступ. Протягом останніх двох десятиріч значною мірою актуалізувалися мегатренди цифровізації та віртуалізації процедур обробки великих масивів даних [1-4]. Розвиток глобальної мережі Інтернет та локальних мереж і способів мобільної передачі даних, а також поява мережевих сервісів (зокрема, хмарних платформ) призвів до появи парадигми «Інтернету речей» (Internet of Things, IoT) та «Інтернету всього» (Internet of Everything, IoE). Концепції, що лежать в основі зазначених парадигм, з одного боку, принципово розширюють функціональні можливості мережевих інформаційних систем (Network Information Systems, NIS), а з іншого боку — призводять до експоненційного росту вимог до параметру перепускності мережевих каналів та вимог до обчислювального ресурсу комплексу [5, 6], що, зокрема, пов'язано з необхідністю розгортання

комплексної та цілісної методології захисту «чутливих» даних (Sensitive Information, SI), що формалізується через політику захисту даних від втрат (Data Leakage Prevention, DLP). Зазначені тенденції вказують на необхідність оптимізації процесу передачі даних шляхом застосування математичних алгоритмів потокової обробки даних (Data Stream Processing, DSP). Запропонований підхід дозволяє формалізувати задачу обробки даних на рівні побудови комплексної схеми оптимізації відповідних алгоритмів, які базуються на визначенні екстремумів цільових функцій, що вказує на *актуальність* даного дослідження.

Аналіз сучасних досліджень у галузі потокової обробки даних включав у себе роботу з публікаціями у фахових журналах та масиви статистичних даних [1-6]. У результаті аналізу було визначено переваги системної обчислювальної моделі процедури потокової обробки масивів даних, що працює зі вхідними запитами у режимі реального часу [7-11]. Пріоритет було надано підходам, що включають у себе централізацію компонент потокової обробки даних (Data Streaming Centralized Components, DSCC), зокрема роботу з мережевими середовищами структурованої потокової передачі даних «Spark Streaming» [9-10]. Особливості зазначеного підходу полягають у застосуванні алгоритмів обробки вхідних запитів з більшою пропускнуою здатністю і, відповідно, мінімізацією показника затримки [9-11]. Аналіз також включав у себе роботу з моделями граничних мережових середовищ, зокрема такими як то оверлейні мережі, а відповідно і моделях розгортання вузла обробки поточкових даних та його оптимізації через зменшення навантаження на обчислювальний ресурс загальної системи і перепускність каналу передачі даних та застосування запиту подібності [12, 13]. Крім цього у аналіз було включено методи організації багатоетапної процедури роботи з поточковими даними [14-16] і методики стабілізації виконання мережових алгоритмів шляхом включення принципу реплікації даних [9-11, 17, 18].

Проведений аналіз вказав на необхідність загальнення методик побудови мережових алгоритмів для роботи з поточковими даними шляхом проведення розподілених обчислень, що було виділено як *невирішену частину загальної проблеми*.

Таким чином, *метою роботи* стала розробка цілісної методології оптимізації процедури організації розподілених обчислень при роботі з поточковими даними шляхом обчислення екстремумів цільових функцій пропускнуої здатності і часу затримки при обробці вхідних запитів.

1. Загальні принципи потокової передачі масивів даних у рамках хмарного сервісу

Математичне моделювання процедури потокової передачі масивів даних базується на визначенні форматів та об'ємів даних, з якими працює мережовий ресурс, актуальних підходів побудови програмних алгоритмів і апаратних засобів реєстрації, передачі, обробки та збереження даних. Так, аналіз показує зростання загальних об'ємів цифрових даних, що передаються через глобальні і локальні мережі пов'язано з наступними ключовими факторами у сфері інформаційних технологій (ІТ):

- розробка і активне впровадження дешевих, компактних та високоякісних систем аудіо-, фото- та відеореєстрації;
- експоненційний ріст тактової частоти центральних процесорів (central processor, CP), поява мультипроцесорних систем та багатоядерних процесорів;
- експоненційний ріст щільності запису носіїв інформації (оптичних, магнітних та твердотільних накопичувачів);
- експоненційний ріст перепускності мережових каналів;
- поява форматів стиснення цифрових мультимедійних даних (аудіофайлів, зображень та відеоданих).

Вказаний перелік включає у себе фактори розширення можливостей для передачі і обробки даних та водночас і фактори появи критичних обмежень на перепускність інформаційних каналів загальної системи. Для визначення «вузьких місць» по відношенню до конкретної задачі необхідно побудувати базову схему потокової передачі і обробки вхідних даних. Таким чином, з метою формалізації зазначених процедур на математичному рівні пропонується ввести наступні позначення (для прикладу представлено схему роботи з графічними даними та відеоданими):

- множина робочих станцій представлена набором $A: \{A_n\}$, де $n \in [1; N]$, кожен з елементів якого, у свою чергу, представляє собою набір параметрів $\{A_n^k\}$ (середня тактова частота CP, кількість CP, кількість ядер, ядерний коефіцієнт, тощо), де $k \in [1; K]$ для $\forall n \in [1; N]$;
- множина об'єктів масиву поточкових даних $B: \{B_l\}$, де $l \in [1; L]$, кожен з елементів якого, у свою чергу, представляє собою набір параметрів $\{B_l^m\}$ (кількість кольорових каналів, динамічний діапазон кольорового каналу, розмір кадру, розмір відеофрагменту, кількість кадрів у відеофрагменті, тощо), де $l \in [1; L]$ для $\forall m \in [1; M]$;

- функція що моделює завантаження робочої станції при обробці вхідних запитів $F_A(A_n)$;
- функція стиснення відповідно форматуваних вхідних даних алгоритму стиснення $F_B(B_l)$.

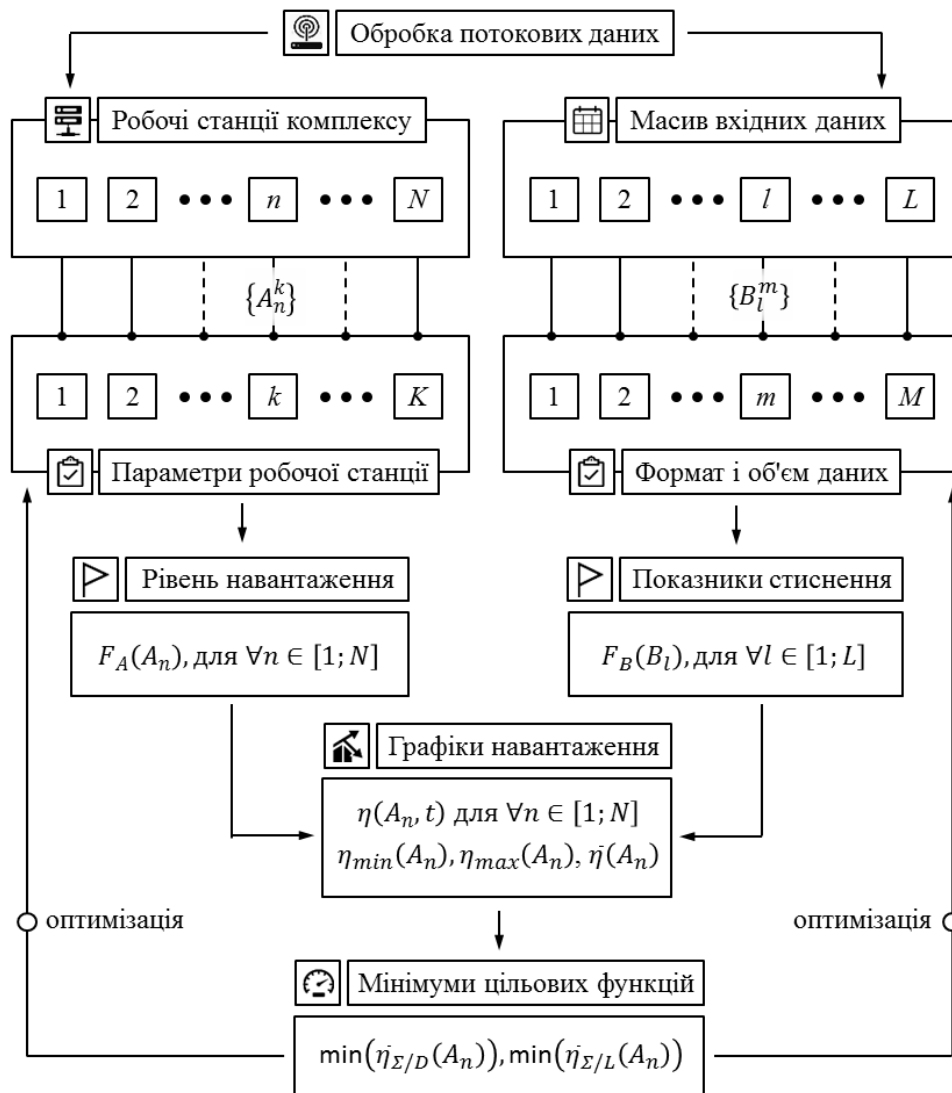


Рис. 1. Базова схема організації та оптимізації процедури потокової передачі масивів даних

Відповідно до побудованого математичного апарату можна визначити графіки навантаження для кожною з робочих станцій $\eta(A_n, t)$ від часу, а також отримати набори мінімальних, максимальних і середніх значень. Відповідно до цього необхідно врахувати усередненні значення затримки $\bar{\eta}_{\Sigma/D}(A)$ та втрат $\bar{\eta}_{\Sigma/L}(A)$ для загальної системи аналізу, пошук мінімумів яких, відповідно, надалі дозволить оптимізувати процедуру обробки вхідних поточкових даних (рис. 1).

Аналіз сучасних мережевих систем обробки поточкових даних, що працюють з великими масивами вхідних даних, вказує на надмірний рівень затримки і втрат. Це пов'язано з нестабільністю навантаження на вузол координації обробки запитів, що, у свою чергу пов'язано з архітектурою апаратно-програмної платформи типового хмарного сервісу. На сьогоднішній день найбільш ефективним підходом вирішення поставленого завдання є організація розподіленої інформаційної системи (Distributed Information System, DIS) шляхом застосування стандартної методики периферійних обчислень (Edge Computing, ES).

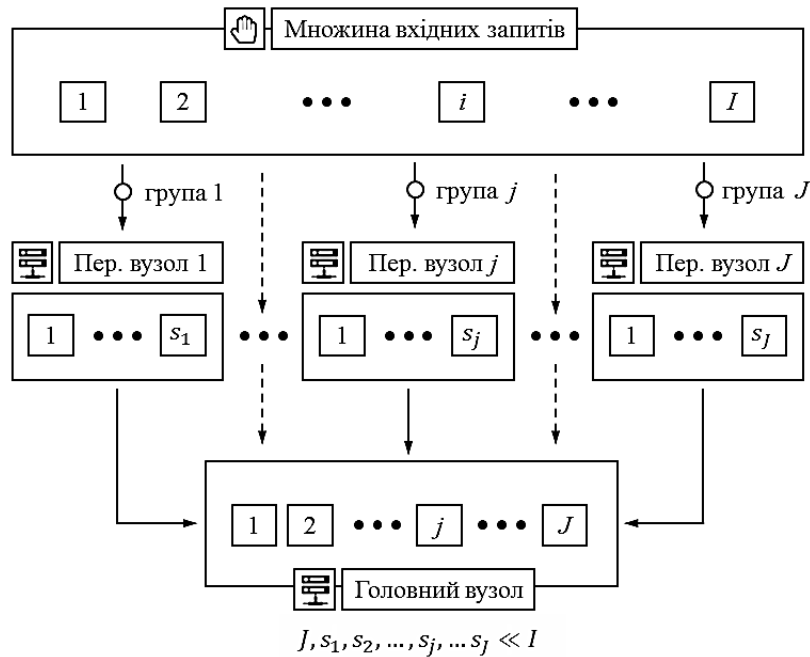


Рис. 2. Узагальнена схема організації розподіленої інформаційної системи обробки поточкових даних

Відповідно методики DIS повний набір вхідних запитів $\{q_i\}$, де $i \in [1; I]$, має бути поділено на J груп запитів (набір $\{Q_j\}$), причому кількість елементів у кожній групі $\forall j \in [1; J]$ визначається через набір $\{s_j\}$. І якщо у випадку класичної організації інформаційної системи перепускність каналу визначається через максимально можливе значення I , для DIS шляхом ускладнення архітектури і організації периферійних вузлів перепускність окремого каналу можна зменшити до $\max(J, s_1, s_2, \dots, s_j, \dots, s_J) \ll I$, як це показано на рис. 2.

При цьому має бути вирішена задача рівномірного розподілу периферійних вузлів відносно обчислювальних ресурсів робочих станцій та задача узгодження процедур обробки запитів при обробці поточкових даних. Слід зазначити, що у рамках DIS-архітектури один запит може оброблятися кількома периферійними вузлами, що характеризуються різними показниками обчислювальної потужності і, відповідно, часу затримки. Це зумовлює впровадження алгоритмів паралельної обробки поточкових даних. При цьому узгодження графіка виконання процедур обробки запитів включає у себе етап розгортання завдань з фіксованою маршрутизацією потоку та введення політики потокової маршрутизації з фіксованим розташуванням окремих завдань, причому аналіз складності завдання враховує необхідність виконання зазначених етапів відповідно мети мінімізації параметри затримки обробки вхідного запиту.

2. Математичне моделювання та оптимізація виконання процедури маршрутизації потоків даних

Математичне моделювання апаратно-програмного комплексу з архітектурою DIS, що складається з одного центрального вузла ($j = 0$) та $j \in [1; J]$ периферійних вузлів, включає у себе формалізацію наступних параметрів:

- набір запитів, що підлягають обробці (множина $\{q_i\}$, де $i \in [1; I]$);
- потік даних, що включає у себе набір елементів $\{d_i\}$ (де $d_i \in [d_i^1; d_i^{N_i}]$), який має бути оброблено при виконанні завдання q_i для всіх $i \in [1; I]$;
- копія потоку даних, що включає у себе набір елементів $\{c_i\}$, для яких $c_i \in [c_i^1; c_i^{K_i}]$, яка може бути розгорнута при виконанні завдання q_i для всіх $i \in [1; I]$;
- перепускність $\sigma(d_i)$, що характеризує затримку при обробці потоку даних $\{d_i\}$ та перепускність окремого вузла обробки даних σ_j (для $\forall j \in [0; J]$);
- матриця $M(i, j)$ затримки міжвузлом потоку даних $\{d_i\}$ та іншими вузлами;
- функція $\vartheta(i, j, c_i)$ відображення процедури розгортання копії c_i на окремому вузлі j ;
- функція $\theta(i, d_i, c_i)$ відображення процедури маршрутизації поточкових даних d_i через копію c_i .

Відповідно, розгортання копії c_i визначається за допомогою функції $\vartheta(i, j, c_i)$, що співвідносить завдання, копію потоку та вузол. Якщо функція приймає значення $\vartheta(i, j, c_i) = 1$ на відповідному вузлі розгортається одна копія, у протилежному випадку, коли функція $\vartheta(i, j, c_i) = 0$ — копія не розгортається. Формалізувати множину значень функції та визначення K_i можна за допомогою наступної системи рівнянь:

$$\begin{cases} \vartheta(i, j, c_i) \in \{0; 1\} \\ K_i = \sum_{j=0}^J \left(\sum_{i=1}^I \vartheta(i, j, c_i) \right) \end{cases} \quad (1)$$

Аналогічно, функція $\theta(i, d_i, c_i)$ застосовується для відображення процедури маршрутизації поточкових даних d_i через копію c_i завдання, що підлягає обробці. Якщо функція приймає значення $\theta(i, d_i, c_i) = 1$ відбувається маршрутизація потоку d_i , а коли функція $\theta(i, d_i, c_i) = 0$ — маршрутизація потоку d_i не відбувається. На рівні побудови математичного апарату це формалізується як:

$$\theta(i, d_i, c_i) \in \{0; 1\} \text{ для } \forall q_i, d_i, c_i \text{ де } i \in [1; I]. \quad (2)$$

Запропонований підхід на базовому рівні дозволяє провести розрахунок та, відповідно, оптимізувати параметри обробки поточкових даних відповідно споживання обчислювальних ресурсів, перепускності обчислювальних вузлів а також рівню затримки під час обробки вхідних запитів.

З метою розширення функціональності побудованої моделі слід ввести поняття ресурсу каналу, що відображатиметься функцією $R(i, j, c_i)$. Відповідна величина характеризує ресурс, що використовується для виконання системою функцій прийому та передачі поточкових даних. Якщо копія c_i запиту q_i по обробці поточкових даних d_i розгорнута на вузлі j' , функція $R(i, j, c_i)$, що обчислюється через функцію $\vartheta(i, j, c_i)$ і перепускність $\sigma(d_i)$, визначатиме ресурс, що споживається копією c_i на всіх вузлах $\forall j \neq j'$:

$$\begin{cases} R(i, j, c_i) = \sum_{d_i} (\sigma(d_i) \cdot \vartheta(i, j, c_i)) \\ j \neq j' \end{cases} \quad (3)$$

Повне значення спожитого ресурсу каналу при передачі поточкових даних на вузлі j у рамках обробки запиту q_i має бути меншим за перепускність вузла. Тому необхідно визначити максимальне значення ресурсу каналу $R_{max}(i, j, c_i)$, при цьому сумарне значення ресурсу, що споживається запитом вузлом j — $R_{\Sigma}^j(i, j, c_i)$ та іншими вузлами — $R_{\Sigma}^{\pm}(i, j, c_i)$, має бути меншим або дорівнювати $R_{max}(i, j, c_i)$:

$$\begin{cases} R_{\Sigma}^j + R_{\Sigma}^{\pm} \leq R_{max}(i, j, c_i) \\ R_{\Sigma}^j(i, j, c_i) = \sum_{d_i \in \{d_i^j\}} \left(\sum_{c_i \in [c_i^{\pm}; c_i^{K_i}]} (\vartheta(i, j, c_i) \cdot R(i, j, c_i)) \right) \\ R_{\Sigma}^{\pm}(i, j, c_i) = \sum_{d_i \notin \{d_i^j\}} \left(\sum_{c_i \in [c_i^{\pm}; c_i^{K_i}]} (\vartheta(i, j, c_i) \cdot R(i, j, c_i)) \right) \end{cases} \quad (4)$$

Аналогічним чином на основі функції $\theta(i, d_i, c_i)$ і матриці $M(i, j)$ затримки міжвузлом потоку даних $\{d_i\}$ та іншими вузлами можна розрахувати значення затримки при виконанні запитів:

$$S(i, d_i, c_i) = \sum_{d_i \notin \{d_i^j\}} (\theta(i, d_i, c_i) \cdot M(i, j)). \quad (5)$$

Сумарна затримка при обробці всіх запитів системи, що може бути використана як цільова функція оптимізація процесу обробки поточкових даних, відповідно розраховуватиметься як:

$$S_{\Sigma} = \sum_{i \in [1; I]} \left(\sum_{c_i \in [c_i^1; c_i^{K_i}]} \left(\sum_{j \in [0; J]} (\vartheta(i, j, c_i) \cdot S(i, d_i, c_i)) \right) \right). \quad (6)$$

Таким чином, задача оптимізація роботи системи обробки поточкових даних апаратно-програмним комплексом хмарного сервісу може бути вирішена через вирішення математичної задачі пошуку мінімуму цільової функції $\min(S_{\Sigma})$ для всіх $i \in [1; I]$, $c_i \in [c_i^1; c_i^{K_i}]$ та $j \in [0; J]$

Висновки. В результаті проведеного дослідження було визначено сучасні тенденції у області реєстрації, обробки, передачі та зберігання великих об'ємів даних. Проведено аналіз сучасних алгоритмів поточної обробки масивів цифрових даних та методик формалізації процедур обробки і передачі даних з метою побудови ефективної математичної моделі. Побудовано узагальнену схему поточної передачі великих масивів графічних даних і відеоданих, а також схему апаратно-програмного комплексу хмарного сервісу, що базується на архітектурі розподіленої інформаційної системи. Вказано на особливості організації апаратно-програмного комплексу мережевого вузла відповідно до схеми розподіленої інформаційної системи, вказано на задачі, що мають бути вирішені з метою оптимізації зазначеної архітектури. Проведено аналіз оптимізації налаштування графіка обробки запитів відповідно до особливостей роботи загального комплексу обробки поточкових даних та вирішення задачі налаштування алгоритмів паралельної обробки. Розроблено математичну модель розподіленої інформаційної системи апаратно програмного комплексу хмарного сервісу, що складається з центрального обчислювального вузла і периферійних обчислювальних вузлів, та включає у себе параметри складових архітектури, функції відображення процедури розгортання завдання і функції маршрутизації потоку вхідних даних. На базі побудованої математичної моделі розроблено методику проведення розрахунку часу затримки при обробці запитів користувачів хмарного сервісу, що запропоновано розглядати як показники цільових функцій.

References.

1. Hummer, W., Satzger, B., & Dustdar, S. (2013). Elastic stream processing in the cloud. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 3(5), 333-345. doi:10.1002/widm.1100.
2. Zeadally, S., Das, A. K., & Sklavos, N. (2019). Cryptographic technologies and protocol standards for Internet of Things. *Internet of Things*, 100075. doi: 10.1016/j.iot.2019.100075.
3. Li, J., Pu, C., Chen, Y., Gmach, D., & Milojevic, D. (2016). Enabling elastic stream processing in Shared Clusters. *2016 IEEE 9th International Conference on Cloud Computing (CLOUD)*. doi:10.1109/cloud.2016.0024.
4. Akidau, T., et al.: The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing. In: *Very Large Data Bases 2015*, vol. 8, pp. 1792–1803 (2015).
5. Kersten, H., & Klett, G. (2014). *Data leakage prevention*. Heidelberg: Mitp.
6. Löhel, J. (2014). *Data Leakage Prevention: Der Einsatz von Datenträger-verschlüsselung zur erweiterten Absicherung der mobilen IT-Nutzung*. Erscheinungsort nicht ermittelbar: Verlag nicht ermittelbar.
7. Soyata T., et al.: Combat: mobile cloud-based compute/communications infrastructure for battlefield applications. In: *Proceedings of SPIE*, vol. 8403, pp. 1–13. <https://doi.org/10.1117/12.919146>.
8. Mayer-Schonberger, V. (2013). *Big Data*. London: John Murray General Publishing Division..
9. Dixit, A., Choudhary, J., & Singh, D. P. (2018). Survey of Apache Storm Scheduler. *SSRN Electronic Journal*. doi: 10.2139/ssrn.3168564.
10. Ganesan, D. *Apache Spark: Einführung zu Technologie und Anwendung*. (2017). Troisdorf: SIGS DATACOM GmbH..
11. Chintapalli, S., et al.: Benchmarking streaming computation engines: storm, flink and spark streaming. In: *International Parallel and Distributed Processing Symposium 2016*, pp. 1789–1792 (2016).
12. Rahman, A., Liu, X., & Kong, F. (2014). A survey on geographic load balancing based data center power management in the smart grid environment. *IEEE Communications Surveys & Tutorials*, 16(1), 214-233. doi:10.1109/surv.2013.070813.00183.
13. Jonathan, A., Chandra, A., Weissman, J.B.: Multi-query optimization in wide area streaming analytics. In: *Symposium on Cloud Computing 2018*, pp. 412–425 (2018).
14. Femminella, M., Pergolesi, M., & Reali, G. (2016). Performance evaluation of edge cloud computing system for big data applications. *2016 5th IEEE International Conference on Cloud Networking (Cloudnet)*. doi:10.1109/cloudnet.2016.56.
15. Heintz, B., Chandra, A., Sitaraman, R.K.: Optimizing grouped aggregation in geo-distributed streaming analytics. In: *High Performance Distributed Computing 2015*, pp. 133–144 (2015).
16. Barretto, W., B. Kochem Vendramin, A. C., & Fonseca, M. (2019). RW-Through: A data replication protocol suitable FOR GeoDistributed And Read-intensive workloads. *Workshop Em Clouds E Aplicações*. doi:10.5753/wcga.2019.7592.
17. Yin, F., Li, X., Li, X., & Li, Y. (2019). Task Scheduling for Streaming Applications in a Cloud-Edge System. *Security, Privacy, and Anonymity in Computation, Communication, and Storage Lecture Notes in Computer Science*, 105–114. doi: 10.1007/978-3-030-24900-7_9.
18. Hwang, J., Cetintemel, U., Zdonik, S.B.: Fast and highly-available stream processing over wide area networks. In: *International Conference on Data Engineering 2008*, pp. 804–813 (2008).