

## АНАЛІЗ АРХІТЕКТУРНИХ ОСОБЛИВОСТЕЙ ОБЧИСЛЮВАЛЬНИХ СИСТЕМ З ПРОГРАМОВАНОЮ СТРУКТУРОЮ

*Стаття присвячена побудові мультиархітектурних кластерних та мультикластерних обчислювальних систем. В роботі проведено аналіз обчислювальних систем з програмованою структурою. Дослідження показали, що сучасні розподілені ОС є мультиархітектурними. Залежно від рівня розгляду їх функціональних структур, вони можуть виглядати і як MISD, і як SIMD, і як MIMD системи. Основна функціонально-структурна одиниця обчислювальних ресурсів в системах розглянутого класу – це елементарна машина (ЕМ), яка є композицією з обчислювального модуля і системного пристрою. Встановлено, що ОС з програмованою структурою відносяться до засобів обробки інформації що масштабуються та допускають формування конфігурацій з масовим паралелізмом (Scalable Massively Parallel Architecture Computing Systems). Взаємодія між ЕОМ здійснюється за допомогою програмного налаштування мережі зв'язку. Структура ОС описується графом, множина вершин якого – кількість елементарних машини (або системних пристроїв, або локальних комутаторів), а множина ребер - лінії міжмашинного зв'язку. До структур сучасних ОС пред'являється ряд вимог: простота вкладення паралельного алгоритму, зручність адресації елементарних машин, здійсненність принципу близькодії і мінімуму затримок при міжмашинних передачах інформації, масштабованість і великомасштабні структури, комутація, живучість, технологічність. Найбільш повно перерахованим вище вимогам задовольняють однорідні структури (тобто описувані однорідними графами). Такі структури є перспективними для формування масштабованих і великомасштабних обчислюваних систем (зокрема, ОС з програмованою структурою). Для оцінки продуктивності каналів зв'язку між ОС використано такі показники як пропускна здатність і латентність, що дозволяють здійснити досить повний аналіз структурних можливостей ОС. В роботі наведено приклад побудови обчислювального кластера, укомплектованого N обчислювальними вузлами, які об'єднані мережею зв'язку стандарту Gigabit Ethernet.*

*Ключові слова: обчислювальна система, елементарна машина, система обробки даних, структура, однорідна структура, комунікаційні середовища, кластер, мультикластер, алгоритми вкладення.*

**Вступ.** Зростаюча потреба у вирішенні складних завдань науки і техніки привела до створення розподілених обчислювальних систем (ОС). В архітектурному плані розподілена ОС це група взаємодіючих елементарних машин (ЕМ), оснащених засобами комунікацій і зовнішніми пристроями. Всі основні ресурси розподілених ОС (арифметико логічні пристрої, пам'ять, засоби управління і комунікацій) являються логічно і технічно розосередженими. Число ЕМ в розподілених ОС допускає варіювання від декількох одиниць до сотень тисяч (наприклад, в МВС-15000ВМ кількість процесорів дорівнює 1148, в IBM Roadrunner – 122400, а в системах IBM BlueGene другого покоління до 884736).

**Постановка задачі.** Сучасні розподілені ОС є мультиархітектурними. Залежно від рівня розгляду їх функціональних структур, вони можуть виглядати і як MISD, і як SIMD, і як MIMD системи [1 – 7]. Для таких систем характерні ієрархічна організація і різні пропускні спроможності каналів зв'язку між їхніми ресурсами (обчислювальними вузлами, ЕМ, процесорами і їх ядрами).

Час виконання паралельних програм на розподілених ОС істотно залежить від того наскільки вони ефективно вкладені в систему [8 – 11]. Під ефективним вкладенням розуміється такий розподіл гілок паралельної програми між ЕМ системи, при якому

досягаються мінімуми накладних витрат на міжмашинний обмін інформацією і дисбаланс завантаження ЕМ.

При організації ефективного функціонування розподілених ОС попереду стоїть завдання розробки моделей і алгоритмів вкладення па паралельно програм, які враховують архітектурні особливості сучасних систем.

**Виклад основного матеріалу дослідження.** Обчислювальні системи з програмованою структурою - це розподілені засоби обробки інформації з архітектурою MIMD (Multiple Instruction Multiple Data)- системи із множинним потоком команд і множинним потоком даних. До цього класу належать як векторні суперЕОМ, так і всі багатопроцесорні системи обробки даних (СОД). Загальна структурна схема таких систем показана на рис. 1. Ця архітектура включає всі рівні паралелізму, від конвеєра операцій до незалежних операцій і команд. Вживаючи термін MIMD треба мати на увазі не тільки роботу багатьох процесорів, але і безліч обчислювальних процесів, що виконуються одночасно в ОС.

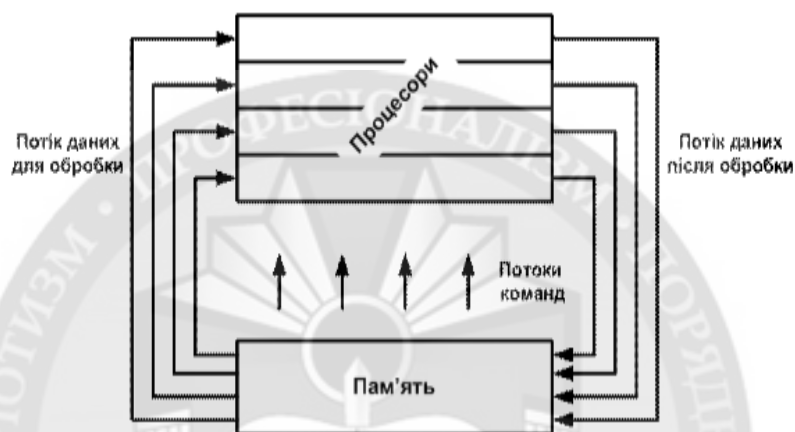


Рис. 1. Загальна структура системи класу MIMD

Особливість систем в тому, що в них закладена можливість програмного переналаштування архітектури MIMD в архітектури MISD або SIMD. Основна функціонально-структурна одиниця обчислювальних ресурсів в системах розглянутого класу - це елементарна машина, яка є композицією з обчислювального модуля і системного пристрою.

Обчислювальний модуль (ОМ) служить як для обробки та зберігання інформації, так і для виконання функцій з управління системою в цілому.

Системний пристрій (СП) - це та апаратна частина ЕМ, яка призначається тільки для забезпечення взаємодії даного ЕМ з найближчими сусідніми машинами (точніше, з системними пристроями, з якими є безпосередній зв'язок).

Допускається конфігурація ОС з довільним числом ЕМ. Отже, ОС з програмованою структурою відносяться до засобів обробки інформації що допускають формування конфігурацій з масовим паралелізмом (Scalable Massively Parallel Architecture Computing Systems).

Взаємодія між ЕМ здійснюється через програмно налаштовану мережу зв'язку. Структура ОС описується графом  $G = (C, E)$ , множина вершин  $C$  якого - кількість елементарних машини (або системних пристроїв, або локальних комутаторів), а множина ребер  $E$  - лінії міжмашинного зв'язку.

До структур сучасних ОС висувається ряд вимог [3].

1) Простота вкладення паралельного алгоритму розв'язання складної задачі в структуру ОС. Структура ОС повинна бути адекватна досить широкому класу вирішуваних задач; настройка проблемно-орієнтованих віртуальних змін не повинна бути пов'язана зі значними накладними витратами.

2) Зручність адресації елементарних машин і «перенесення» підсистем в межах обчислювальної системи. Обчислювальна система повинна надавати можливість користувачам створювати паралельні програми з віртуальними адресами ЕМ. Отже, структура ОС повинна дозволяти реалізувати найпростіший «механізм» перетворення віртуальних адрес ЕМ в реальні (фізичні) адреси машин системи. Необхідність організації одночасного вирішення декількох задач на ОС (тобто необхідність поділу простору елементарних машин між задачами) обґрунтовує вимога простоти переміщення підсистем в межах системи (при збереженні їх топологічних властивостей).

3) Здійсненність принципу близькодії і мінімуму затримок при міжмашинних передачах інформації в ОС. Принцип близькодії перевизначає реалізацію обмінів інформацією між «віддаленими» один від одного ЕМ через проміжні машини системи. Отже, в умовах обмеженості числа зв'язків, структура ОС повинна забезпечувати мінімум затримок при «транзитних» передачах інформації.

4) Масштабованість і великомасштабні структури ОС. Для формування конфігурацій ОС із заданою ефективністю потрібно, щоб структура була здатна до нарощування і скорочення числа вершин (машин). Зміна числа ЕМ в ОС не повинна призводити до корінних перекомутацій між машинами і (або) до необхідності зміни числа зав'язків для будь-яких ЕМ.

Для досягнення високої продуктивності ОС при існуючих можливостях мікропроцесорної техніки потрібно число ЕМ близько  $10-10^6$ . Для підтримки великомасштабних ОС (масового паралелізму) необхідно, щоб структура ОС була здатна ефективно здійснювати міжмашинний обмін інформацією в умовах неможливості реалізації зв'язків повного графу (наприклад, через обмеженість числа виводів в корпусах ВІС).

5) Комутація структури ОС. Обчислювальна система повинна бути пристосована до реалізації групових міжмашинних обмінів інформацією. Отже, структура ОС повинна мати здатність здійснювати вказану кількість одночасних непересічних взаємодій між елементарними машинами.

6) Живучість структури ОС. Важливою вимогою до ОС в цілому являється забезпечення працездатності при відмові її компонентів або навіть підсистем. Основою функціональної цілісності ОС, як колективу елементарних машин, є живучість структури. Під останньою розуміють здатність структури ОС забезпечити зв'язність необхідного числа робочих ЕМ в системі при ненадійних лініях міжмашинних зв'язків.

7) Технологічність структур ОС. Структура мережі міжмашинних зв'язків ОС не повинна пред'являти особливих вимог до елементної бази, до технології виготовлення мікропроцесорних ВІС. Системи повинні бути сприйнятливі до масової технології, їх «обчислювальне ядро» має формуватися з масових мікропроцесорних ВІС. Останнє дозволить досягти прийнятних значень техніко-економічних показників ОС.

Структурні затримки при передачах інформації між машинами ОС визначаються відстанню (в сенсі теорії графів) між вершинами структури, зіставленими взаємодіючими машинами. Для оцінки структурних затримок в обчислювальних системах використовують діаметр  $d_i$  і середній діаметр  $\bar{d}$  структури.

Діаметр визначається максимальною відстанню на множині коротших шляхів між парами вершин структури ОС:

$$d = \max_{i,j} \{d_{ij}\}, \quad (1)$$

а середній діаметр

$$\bar{d} = (N-1)^{-1} \sum l \cdot n_i, \quad (2)$$

де  $d_{ij}$  - відстань, тобто мінімальне число ребер, що утворюють шлях з вершини  $i$  в вершину  $j$ ;  $i, j \in \{0, 1, \dots, N-1\}$ ;  $n_l$  - число вершин, що знаходяться на відстані  $l$  від будь-якої виділеної вершини (однорідного) графа  $G$ .

Показником що оцінює структуру комутованої ОС, є вектор-функція:

$$\chi(G, s, s^i) = \chi_h(G, s, s^i), \quad h \in \{1, 2, \dots, N/2\}, \quad (3)$$

в якій координата  $\chi_h(G, s, s^i)$ , є ймовірність реалізації в системі при заданій структурі структури  $G$  і коефіцієнтах готовності  $s, s^i$ , відповідно, одній ЕМ і лінії зв'язку  $h$ .

Структурна живучість ОС оцінюється вектор-функцією

$$\chi(G, s, s^i) = \{ \chi_r(G, s, s^i) \}, \quad r \in E_2^N \{1, 2, \dots, N\}, \quad (4)$$

де  $\chi_r(G, s, s^i)$ , є ймовірністю існування підсистеми рангу  $r$  (підмножиною з  $r$  працездатних ЕМ, зв'язність яких встановлюється через працездатні лінії зв'язку) при заданій структурі  $G$ , коефіцієнтах готовності  $s$  і  $s^i$  елементарної машини і лінії зв'язку, відповідно.

Для оцінки продуктивності каналів зв'язку між ЕМ системи використовують показники: пропускна здатність і латентність.

Пропускна здатність каналу зв'язку (bandwidth, throughput, channel capacity) - найбільша кількість інформації, що передається по каналу зв'язку в одиницю часу.

Латентність (latency) - це затримка при передачі інформації між ЕМ системи, викликана програмними і апаратними накладними витратами.

Наведені показники дозволяють здійснити досить повний аналіз структурних можливостей ОС.

Найбільш повно перерахованим вище вимогам задовольняють однорідні структури (тобто ті що описуються однорідними графами). Такі структури є перспективними для формування масштабованих і великомасштабних обчислюваних систем (зокрема, ОС з програмованою структурою).

У комп'ютерній індустрії набули поширення  $n$ -мірні структури обчислювальних систем, відомі зараз як циркулянтні (Circulant Structures). Вперше вони були визначені і досліджені у відділі обчислювальних систем Інституту математики АН СРСР на початку 70-х років і спочатку називалися -  $D_n$  графами [14].

За визначенням  $D_n$  - граф або циркулянтна структура є граф  $G$  виду:  $\{N; \omega_1, \omega_2, \dots, \omega_n\}$ , в якому:

$N$  число вершин або порядок графа;

вершини позначені цілими числами  $i$  по модулю  $N$ , отже,  $i \in \{0, 1, \dots, N-1\}$ ;

вершина  $i$  з'єднана ребром (або є суміжною) з вершинами  $i \pm \omega_1, i \pm \omega_2, \dots, i \pm \omega_n \pmod{N}$ ;

$\{\omega_1, \omega_2, \dots, \omega_n\}$  множина цілих чисел, яких називають утворюючими, такими, що,  $0 < \omega_1 < \omega_2 < \dots < \omega_n < (N+1)/2$ , а для чисел  $N; \omega_1, \omega_2, \dots, \omega_n$  найбільшим спільним дільником є 1;

$n$  - розмірність графа;

$2n$  - ступінь вершини в графі.

Як приклад розглянемо  $D$  -граф (рис. 2) або двовимірний циркулянт виду  $\{12; 3, 4\}$ :

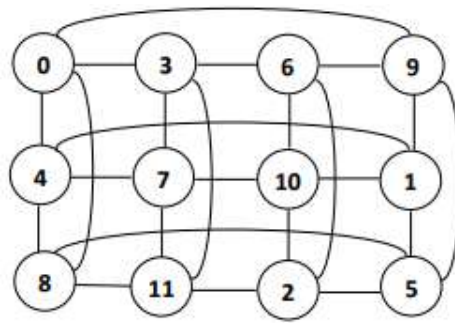


Рис. 2.  $D$  - граф  $\{ 12; 3, 4 \}$

Графи  $G$  виду  $\{ N; \omega_1, \omega_2, \dots, \omega_n \}$  тобто  $D$ - графи або циркулянт з одиничною утворюючою (Loop Networks - петльові структури) інтенсивно вивчаються останнім часом. Циркулянтні структури  $\{ N; 1, \omega_2 \}$ ; широко внесені в практику обчислювальних систем.

Цілі числа  $i \in \{ 0, 1, \dots, N-1 \}$ ; , що визначають вершини  $D$ - графа, називаються адресами. Адресація вершин в таких структурах називається діофантовою (в честь давньогрецького математика з Олександрії Діофанта, Diophantos, 3 століття).

У циркулянтних структурах при повному перенесенні будь-якої підструктури (всіх вершин підструктури на одну і ту ж відстань в одному з напрямків) зберігаються всі її властивості і адресація вершин.

Отже, при діофантовій адресації елементарних машин ОС можна простими засобами реконфігурації здійснити віртуальну адресацію вершин-машин і, отже, створювати відмовостійкі паралельні програми, реалізовувати мультипрограмні режими обробки інформації, виключати вершини-машини, що відмовили, з підсистем, а значить забезпечити живучість ОС.

При цьому алгоритм роботи реконфігуратора структури ОС зводиться до зміни адреси всіх машин підсистеми за формулою:

$$\alpha := [\alpha + (j - i)](\text{mod } N), \alpha \in \{ 0, 1, \dots, N-1 \}, \quad (5)$$

де  $i$  - номер ЕМ, виключеного з підсистеми, а  $j$  - номер машини, підключеної в підсистему,  $i, j \in \{ 0, 1, \dots, N-1 \}$ .

У якості структур ОС, що допускають масштабування (зміна числа машин) без корінної перекомутації вже наявних міжмашинних зв'язків, використовуються  $K(N, v, g)$  -графи [6]. У такі граfi вкладаються  $D$  - граfi; граф  $K(N, v, g)$ - це неорієнтований однорідний граф з числом  $i$  степенями вершин, відповідно,  $N$  і  $v$ , та значенням обхвату  $g$  (рис. 3).

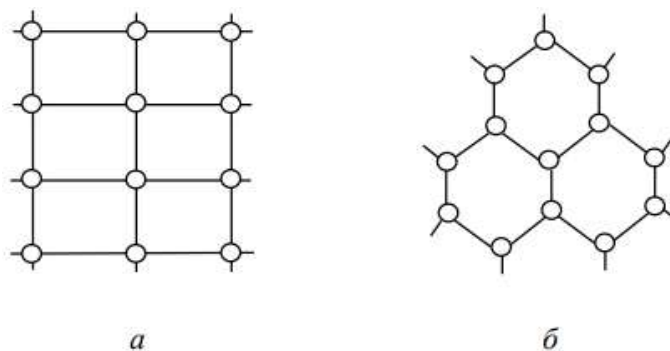


Рис. 3. Фрагменти  $K(N, v, g)$ - графів

В  $K(N, v, g)$  - графах кожна вершина при  $v \geq 3$  входить в не менше  $v$  найкоротших простих циклів довжиною  $g$  (довжина найкоротшого циклу в графі називається обхватом). При  $v = 2 \cdot K(N, v, g)$  -граф є простим циклом з  $N$  вершинами.

Паралельні алгоритми і програми. Поняття паралельного алгоритму (Parallel Algorithm) відноситься до фундаментальних в теорії обчислювальних систем [8]. Паралельний алгоритм - це опис процесу обробки інформації, орієнтований на реалізацію в колективах обчислювачів. Такий алгоритм, на відміну від послідовного, передбачає одночасне виконання множини операцій в межах одного кроку обчислень і як послідовний алгоритм зберігає залежність наступних етапів від результатів попередніх.

Паралельний алгоритм розв'язання задачі становить основу паралельної програми, яка, в свою чергу, впливає на алгоритм функціонування колективу обчислювачів. Запис паралельного алгоритму на мові програмування, доступному колективу обчислювачів, і називають паралельною програмою, а сама мова - паралельною. Паралельні алгоритми і програми слід розробляти для тих задач, які недоступні для вирішення на засобах, заснованих на моделі обчислювача. Ці задачі прийнято називати складними або трудомісткими.

Методи і алгоритми обробки інформації, вирішення задач, як правило, - послідовні. Процес "приспосовування" методів до реалізації на колективі обчислювачів або процес "розщеплення" послідовних алгоритмів вирішення складних задач називається розпаралелюванням (Paralleling).

Теоретична і практична діяльність по створенню паралельних алгоритмів і програм обробки інформації називається паралельним програмуванням (Parallel Programming).

Якість паралельного алгоритму (або його ефективність) визначається методикою розпаралелювання складних задач. Виділяють два основних підходи до розпаралелювання задач: локальне і глобальне (великоблочне) розпаралелювання. Перший підхід орієнтований на розщеплення алгоритму розв'язання складної задачі на гранично прості блоки (операції або оператори) і вимагає виділення для кожного етапу обчислення максимально можливої кількості одночасно виконуваних блоків. Він не призводить до паралельних алгоритмів, ефективно реалізованим колективом обчислювачів. Справді, процес такого розпаралелювання досить трудомісткий, а одержування паралельних алгоритмів характеризуються не тільки структурною неоднорідністю, а й істотно різними обсягами операцій на різних етапах обчислень. Останнє є серйозною перешкодою на шляху (автоматизації) розпаралелювання і забезпечення ефективної експлуатації ресурсів колективу обчислювачів. Локальне розпаралелювання дозволяє оцінити граничні можливості колективу обчислювачів при вирішенні складних задач, отримати граничні оцінки по розпаралеленню складних задач.

Другий підхід орієнтований на розбиття складної задачі на великі блоки-підзадачі, між якими існує слабкий зв'язок. Тоді в алгоритмах, побудованих на основі великоблочного розпаралелювання, операції обміну між підзадачами становитимуть незначну частину в порівнянні із загальним числом операцій в кожній підзадачі. Такі підзадачі називають гілками паралельного алгоритму, а відповідні їм програми - гілками паралельної програми.

Кластерні та мультікластерні ОС і GRID-системи. При побудові кластерних ОС, як правило, використовуються масові апаратно-програмні засоби, що, по суті, є принципом конструювання кластерних ОС та забезпечує їх високу техніко-економічну ефективність [7].

Для створення обчислювальних кластерів використовуються MISD-, SIMD-, MIMD-архітектури, різні функціональні структури і конструктивні рішення.

Основна функціонально-структурна одиниця обчислювальних ресурсів в кластерних ОС - це елементарна машина. Конфігурація ЕМ допускає варіювання в широких межах - від процесорного ядра до ЕОМ, оснащеної засобами комунікацій і зовнішніми пристроями.

Сучасні кластерні ОС [9] переважно конфігуруються з багатопроцесорних вузлів і багатоядерних процесорів. Комунікаційні середовища кластерних ОС будуються на базі комутаторів і мають ієрархічну організацію. На рис. 4 наведено приклад обчислювального кластера, укомплектованого  $N$  обчислювальними вузлами, які об'єднані мережею зв'язку стандарту Gigabit Ethernet. Кожен вузол укомплектований парою двоядерних процесорів AMD

Opteron 275, взаємодіють по шині HyperTransport. Кожен процесор має інтегрований контролер пам'яті, через який його ядра мають доступ до їх загальної пам'яті. Видно, що комунікаційне середовище кластера має ієрархічну організацію, в якій перший рівень - це мережа зв'язку стандарту Gigabit Ethernet; другий - шина HyperTransport; третій - спільна пам'ять ядер. Швидкість передачі інформації на різних рівнях комунікаційного середовища істотно різні.

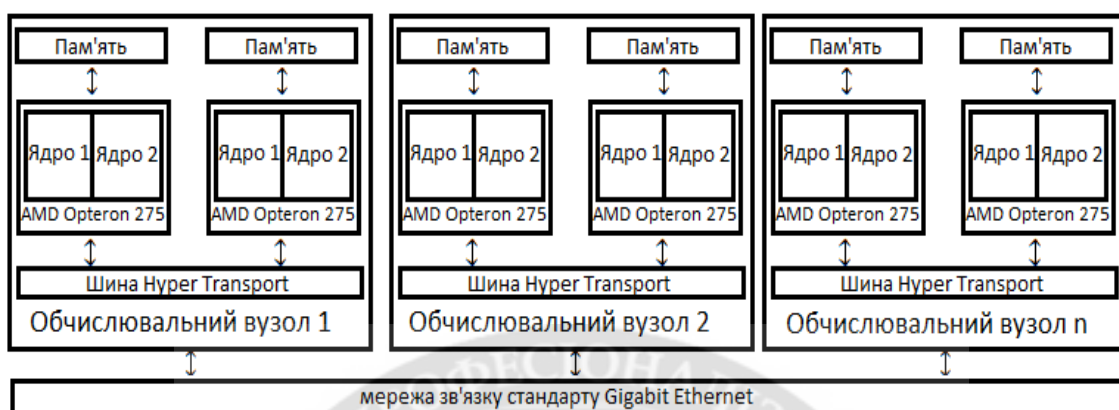


Рис. 4. Приклад ієрархічної організації комунікаційного середовища кластерної ОС

Комунікаційні середовища сучасних кластерних ОС будуються на базі технологій Gigabit Ethernet, InfiniBand, Myrinet [11, 12]. Для об'єднання обчислювальних вузлів в систему використовуються комутатори (switch) з фіксованою кількістю входних портів для підключення кінцевих пристроїв (обчислювальних вузлів, комутаторів, систем зберігання даних і т.д.). У цих умовах для побудови великомасштабних кластерних ОС використовуються різні схеми формування комутаторів. Основна мета - комплексування заданого числа обчислювальних вузлів і забезпечення високої продуктивності комунікаційного середовища.

З розвитком обчислювальної техніки і телекомунікаційних технологій з'явилася тенденція до побудови розподілених ОС шляхом об'єднання за допомогою локальних і глобальних мереж зв'язку розосереджених обчислювальних ресурсів починаючи від простих ПК і закінчуючи ОС з масовим паралелізмом. Так в кінці XX століття в індустрії обробки інформації отримали розвиток GRID-технології (Global Resource Information Distribution) [4]. На основі цих технологій створюються більш масштабні просторово-розподілені системи обробки інформації, здатні реалізувати паралельні алгоритми вирішення суперскладних задач на своїх розосереджених ресурсах. Така GRID-система являє собою композицію множини ЕОМ і ОС, просторово-розподіленої комунікаційної мережі і програмних компонентів для здійснення паралельних обчислень. В GRID-системі можуть використовуватися гетерогенні і не сумісні обчислювальні засоби, однак для їх сумісної роботи над паралельними алгоритмами повинні бути застосовані спеціальні механізми (наприклад, стандартизовані протоколи) узгодженої взаємодії. GRID-система - це ні що інше як розподілена ОС в сенсі введеного вище визначення.

Що стосується мультикластерних ОС, то це не окремий клас ОС, а просторово-розподілені системи, що формуються шляхом об'єднання ресурсів розосереджених кластерів.

Обчислювальні кластери належать до класу систем з розподіленою пам'яттю. Паралельні програми для таких ОС розробляються в моделі передачі повідомлень (message passing). Широке розповсюдження отримали стандарти, бібліотеки і середовища для створення паралельних програм в даній моделі, зокрема стандарт Message Passing Interface (MPI), система Parallel Virtual Machine (PVM), IBM Message Passing Library (MPL), P4 і ін.

В рамках даної моделі паралельні гілки програми синхронізують свою роботу шляхом обміну інформацією по каналах між машинних зв'язків.

Залежно від складності задач і характеру їх надходження в теорії обчислювальних систем виділяють два основні режими функціонування ОС з програмованою структурою [12]: моно- і мультипрограмний режими.

У монопрограмному режимі для вирішення задачі використовуються всі ресурси ОС. Задача представляється у вигляді паралельної програми, число гілок в якій або фіксоване, або допускає варіювання в заданому діапазоні. В якості одиниці ресурсу виступає елементарна машина ОС. Всі машини використовуються для вирішення задачі. Якщо максимальне число гілок в паралельній програмі менше загального числа ЕМ в системі, то «надлишкові» машини використовуються для підвищення надійності функціонування ОС.

До мультипрограмних відносять режими:

- 1) обробки наборів задач;
- 2) обслуговування потоків задач.

При роботі ОС в цих режимах одночасно вирішується кілька задач, отже, ресурси системи діляться між ними.

При організації функціонування ОС у разі обробки набору задач враховується не тільки кількість задач в наборі, але їх параметри: число гілок в програмі (точніше, число машин, на яких вона буде виконуватись), час рішення або імовірнісний закон розподілу часу рішення і ін. Алгоритми організації функціонування ОС задають розподіл задач по машинам і послідовність виконання задач на них. В результаті стає відомим, в якому проміжку часу і на яких машинах (або на якій підсистемі) буде вирішуватися будь яка задача набору.

Обслуговування потоку задач на ОС - принципово відрізняється від вибору наборів: задачі надходять у випадкові моменти часу, їх параметри випадкові, отже, детермінований вибір підсистем для вирішення тих чи інших задач виключений. Для режима потоку задач використовуються методи і алгоритми, що забезпечують оптимальне функціонування обчислювальних систем.

При роботі ОС в будь-якому з мультипрограмних режимів система представляється у вигляді композиції підсистем різних рангів. У міру рішення задачі ця композиція «автоматично» (за допомогою операційної системи) реконфігурується так, щоб забезпечити її адекватність поточної мультипрограмної ситуації. Будь-яка підсистема має всі архітектурні властивості системи, тому її організація при вирішенні виокремленої їй задачі може здійснюватися тими ж методами, якими була організована робота всієї ОС в першому режимі.

Час виконання паралельних програм на розподілених ОС суттєво залежить від того наскільки вони ефективно вкладені в систему [3].

Під ефективним вкладенням (Task Map, Task Allocation, Task Assignment) розуміється такий розподіл гілок паралельної програми між ЕМ системи, при якому досягаються мінімуми накладних витрат на міжмашинному обміні інформацією і дисбалансу завантаження ЕМ.

Залежно від того, чим характеризується паралельна програма, яка надійшла в систему, виникають різні постановки задачі вкладення.

Якщо для паралельної програми заданий інформаційний граф, то говорять про задачу вкладення інформаційного графа (структури) програми в ОС [13,14]. У разі, якщо про паралельну програму відомо лише кількість ЕМ, які необхідно для її реалізації, то розглядається задача формування в межах ОС підсистеми ЕМ, що задовольняє заданим критеріям якості.

Задача оптимального вкладення паралельної програми виникає при організації функціонування ОС в будь-якому з основних режимів: рішення складної задачі, обробки наборів задач, обслуговування потоків задач.

В режимі обслуговування потоків задач в межах ОС формуються підсистеми ЕМ. Для них і надходять паралельні програми для виконання. Елементарні машини, що входять в підсистему, можуть належати різним обчислювальним системам, розміщуватися в різних комутаційних шафах і т.і. Наслідком цього є те, що швидкості передачі інформації між ЕМ можуть бути істотно різними. Тому потрібно не тільки формувати підсистеми з необхідним числом ЕМ, а й ефективно вкладати паралельні програми в них.



У режимі обробки наборів задач будується розклад виконання паралельних програм. В даному режимі задача вкладення виникає або на етапі запуску програми на виділеній підсистемі, або ж на етапі синтезу розкладу, в разі, якщо алгоритми вкладення інтегровані в підсистему підтримки мультипрограмних режимів функціонування ОС.

Існуючі методи і алгоритми вкладення паралельних програм в ОС розрізняються за типом систем, на які вони орієнтовані, по використовуваних моделях паралельних програм, показниками оптимальності вкладення, точності, централізованості або децентралізованості алгоритмів, по статичній або динамічній схемі формування вкладень.

Запропоновано алгоритми вкладення паралельних програм, в яких ОС представлена у вигляді графа міжмашинних зв'язків [11]. Передбачається, що інформація між ЕМ передається по найкоротшому шляху і всі канали зв'язку однорідні. Як показує продуктивність каналу зв'язку між ЕМ використовується довжина шляху (в сенсі теорії графів). Застосування подібних алгоритмів ускладнено тим, що сучасні ОС мультиархітектурні, їх комунікаційні середовища мають ієрархічну організацію і, як наслідок, не однорідні канали зв'язку між ЕМ. Крім того, більшість систем будується на базі комутаторів, в яких застосовуються алгоритми динамічної маршрутизації. Тому доцільно використання моделей обчислювальних систем, в яких продуктивність каналів зв'язку характеризується такими показниками як пропускна здатність і латентність.

Ефективне вкладення в ОС паралельної програми вимагає врахування структури інформаційних обмінів і обсягів даних, що передаються між її гілками. При цьому використовуються алгоритми, які основані лише на структурі інформаційного графа програми і ігнорують обсяги та інтенсивність міжмашинних обмінів [11]. Однак облік таких параметрів практично необхідний. Аналіз популярних вільно розповсюджених паралельних програм, бібліотек і тестів продуктивності High-Performance Linpack, NAS Parallel Benchmarks і SPEC MPI2007 показує, що більшість з них характеризується неоднорідними інформаційними графами. Неоднорідність виражається в обсягах даних, що передаються між гілками, розмірах використовуваних повідомлень, а також в схемах реалізованих обмінів.

Як показники оптимальності вкладення використовують час виконання паралельної програми або функцію, мінімізація якої доставляє мінімум часу виконання програми.

Алгоритми вкладення спираються на припущення про те, що час виконання паралельної програми - це сума часу виконання усіх її гілок. Однак практика рішення промислових задач і аналізу їх продуктивності показує, що час виконання програми визначається максимальним часом виконання її гілок. Крім того, в ряді робіт [8, 10, 11] використовуються складові показники у вигляді лінійної комбінації двох функцій, наприклад, середня завантаженість елементарних машин в системі і час передачі даних по каналах зв'язку. Використання подібних показників ускладнюється вибором коефіцієнтів пропорційності, якими пов'язані функції і найчастіше їх не сумісними розмірностями. Наприклад, розмірність першої функції - кількість операцій, а другої - одиниці часу. Вкладення паралельних програм в просторово розподілені ОС вимагає обліку в показнику часу доставки паралельних програм до елементарних машин системи.

Вкладення в ОС паралельних програм, для яких не задані інформаційні графи, здійснюється шляхом формування підсистеми елементарних машин і розподілу по ній паралельних гілок програми. Сформована підсистема повинна забезпечувати ефективну реалізацію паралельних програм і забезпечувати мінімум часу для їх виконання.

Більшість алгоритмів формування в межах ОС підсистем орієнтовані на системи з певними структурами мереж міжмашинного зв'язку (наприклад, на 2D-решітки, гіперкубічні і тороїдальні структури) і прагнуть зберігати топологічну схожість підсистеми структурі ОС. Можливості застосування таких алгоритмів для ОС з ієрархічною організацією, канали зв'язку в яких мають різну продуктивність, істотно обмежені. Тому актуальною є розробка алгоритмів формування підсистем, що допускають застосування як в системах зі статичною структурою мережі між машинного зв'язку, так і в системах з динамічною структурою.

**Висновки.** Сучасні ОС функціонують під управлінням систем пакетної обробки завдань (наприклад, TORQUE / OpenPBS, IBM LoadLeveler, Sun Grid Engine, Altair PBS Pro, SLURM, Cleo та ін.), які реалізують підтримку мультипрограмних режимів функціонування. Такі системи пакетної обробки завдань (СПОЗ) підтримують черги завдань і здійснюють планування обчислень. Для кожної паралельної програми за допомогою алгоритмів підтримки в межах системи формується підсистема ЕМ.

Поширення отримали елементарні алгоритми формування підсистем: PAL (Processor Allocation Linear) - здійснює вибір перших  $M$  вільних ЕМ, і PAR (Processor Allocation Random) - генерує випадкову підмножину з  $M$  вільних ЕМ.

Початкове вкладення паралельної програми і виділення підсистем також здійснюється СПОЗ. Дане завдання найчастіше вирішується циклічним (Task Map Round Robin, TMRR) або лінійним (Task Map Linear) алгоритмом формування вкладень. Алгоритм TMRR розподіляє паралельні гілки по ЕМ системи з  $N$  обчислювальних вузлів в наступному порядку: перша гілка на перший обчислювальний вузол, друга на другий, ..., гілка  $N$  на вузол  $N$ , гілка  $N + 1$  на вузол  $N + 1$  і т.і. Алгоритм TML розподіляє гілки на перші  $M$  вільних ЕМ.

Засоби оптимізації вкладення паралельних програм присутні і в комунікаційних бібліотеках стандарту MPI (MPICH2, OpenMPI, MVAPICH2, Intel MPI, HP-MPI). Вкладення реалізується утилітою mtrixes, що підтримує в більшості випадків лише алгоритми TMRR і TML.

Аналіз архітектурних особливостей обчислювальних систем з програмованою структурою дозволяє зробити висновок:

1. Сучасні обчислювальні системи є мультиархітектурними, їх комунікаційні середовища мають ієрархічну організацію.

2. Ефективність експлуатації ресурсів розподілених ОС значно залежить від того, як організовано їх функціонування. У будь-якому з режимів функціонування виникає завдання ефективного вкладення паралельної програми в ОС. Під ефективним вкладенням розуміється такий розподіл гілок паралельної програми, при якому досягаються мінімуми накладних витрат на міжмашинний обмін і дисбаланс завантаження елементарних машин.

3. Існуючі алгоритми не дозволяють будувати ефективні вкладення паралельних програм в ОС з ієрархічною організацією комунікаційних середовищ. Актуальною є розробка не трудомістких методів і алгоритмів вкладення паралельних програм в ОС з ієрархічною організацією комунікаційних середовищ.

4. Задача оптимального вкладення є важкою. Актуальною є розробка наближених методів і алгоритмів вкладення паралельних програм в розподіленій ОС.

#### ЛІТЕРАТУРА:

1. Абрамов В.О. Архітектура електронно-обчислювальних машин: навч. посіб. К.: КМПУ імені Б.Д. Грінченка, 2007. 84 с.
2. Арсенюк І. Р. Теорія алгоритмів: навч. посібник. Вінниця: ВНТУ, 2006. – 150 с.
3. Воеводин В.В., Воеводин Вл. В. Параллельные вычисления. Спб.: БХВ-Петербург, 2002. – 608 с.
4. Дорошенко А.Е. Математические модели и методы организации высокопроизводительных параллельных вычислений. Киев: Наук. думка, 2000. 176 с.
5. Жуков І.А. Корочкін О.В. Паралельні та розподілені обчислення: навч. посібник. К.:Корнійчук, 2005. 226 с.
6. Каляев І.А., Левин І.І., Семерников Е.А., Шмойлов В.І. Реконфигурируемые мультимониторные вычислительные структуры. Ростов-на-Дону: ЮНЦ РАН, 2008. 397 с.
7. Курносів М. Г. Опыт построения кластерных вычислительных систем с удаленной загрузкой узлов. Высокопроизводительные параллельные вычисления на кластерных системах: материалы пятого Международного научно-практического семинара. Нижний Новгород, 2005. – С. 149-154.
8. Курносів М. Г. Назначение ветвей параллельной программы на процессорные ядра распределенной вычислительной системы. Многопроцессорные вычислительные и управляющие системы: матер. междуна. науч.-техн. конференции. Таганрог, 2007. – С. 227-231.

9. Левин В. К. Современные суперкомпьютеры семейства МВС. Лаборатория Параллельных Информационных Технологий, НИВЦ МГУ. URL: [http://www.computer-museum.ru/histussr/super\\_phase\\_0.htm](http://www.computer-museum.ru/histussr/super_phase_0.htm) (дата звернення: 02.06.2018).

10. Литвинов О.А., Хандецкий В.С. Розподілена обробка інформації: монографія. Д.: ТОВ «Баланс-Клуб», 2013. – 314 с.

11. Палагин А.В., Опанасенко В.Н. Реконфигурируемые вычислительные системы: Основы и приложения: монографія. Киев: Просвіта, 2006. 280 с.

12. Палагин А. В., Яковлев С. Ю. Системная интеграция средств компьютерной техники. Вінниця: Універсум-Вінниця, 2005. 680 с.

13. Хорошевский В.Г., Курносоев М.Г. Моделирование алгоритмов вложения параллельных программ в структуры распределенных вычислительных систем. Моделирование-2008 (Simulation-2008): Труды Международной научной конференции. Киев: Изд-во ИПМЭ им. Г.Е. Пухова, 2008. Т. 2. С. 435-440.

14. Ясько М.М. Паралельні та розподілені обчислення. – Д.: РВВ ДНУ, 2010. – 111с.

#### REFERENCES:

1. Abramov V.O. (2007), “Arkhitektura elektronno-obchysljuvalnykh mashyn: navch. posib.” [Architecture of electronic computing machines], KMPU imeni B.D. Hrinchenka, Kyiv, 84 p.

2. Arseniuk I.R. (2006), “Teoriia alghorytmiv: navch. posibnyk” [Theory of algorithms], VNTU, Vinnytsia, 150 p.

3. Voevodyn V. V. and Voevodyn Vl. V. (2002), “Parallelnyye vychysleniya” [Parallel computing], BKhV-Peterburgh, Peterburgh, 608 p.

4. Doroshenko A.E. (2000), “Matematicheskiye modely y metody orhanyzatsyy vysokoproizvodytelnykh parallelnykh vychysleniy” [Mathematical models and methods for organizing high-performance parallel computing], Nauk. dumka, Kyev, 176 p.

5. Zhukov I. A. and Korochkin O.V. (2005) “Paralelni ta rozpodileni obchyslennia: navch. posibnyk” [Parallel and distributed computing], Korniiichuk, Kyiv, 226 p.

6. Kaliaev Y. A., Levyn Y.Y., Semernikov E.A. and Shmoilov V.Y. (2008) Rekonfihuryruemye multikonveierynye vychyslytelnye struktury [Reconfigurable multicore computing structures], YuNTs RAN, Rostov-na-Donu, 397 p.

7. Kurnosov M. H. (2005), “Opyt postroyeniya klasternykh vychyslytelnykh system s udalenoii zahruzkoii uzlov” [Experience in building clustered computing systems with remote node load], Vysokoproizvodytelnye parallelnyye vychysleniya na klasternykh systemakh: materiyaly piatoho Mezhdunarodnoho nauchno-praktycheskoho semynara. Nyzhnyi Novhorod, pp. 149-154.

8. Kurnosov M. H. (2007) “Naznacheniye vetvei parallelnoi prohrammy na protsessornyye yadra raspredelennoi vychyslytelnoi systemy” [Assigning branches of a parallel program to the processing cores of a distributed computing system], Mnohoprotsessornyye vychyslytelnye y upravliaiushchye systemy: mater. mezhdun. nauch. -tekh. Konferentsyy, Tahanroh, pp. 227-231.

9. Levyn V. K. (2018), “Sovremennyye superkompiutery semeistva MVS” [Modern supercomputers of the MVS family], Laboratoriya Parallelnykh Ynformatsyonnykh Tekhnolohiyi, NYVTs MHU. URL: [www.computer-museum.ru/histussr/super\\_phase\\_0.htm](http://www.computer-museum.ru/histussr/super_phase_0.htm) (data zvernennia: 02.06.2018).

10. Lytvynov O.A. and Khandetskiy V.S. (2013) “Rozpodilena obrobka informatsii: monohrafiia” [Distributed processing of information], TOV «Balans-Klub», Dnepr, 314 pp.

11. Palahyn A.V. and Opanasenko V.N. (2006), “Rekonfihuryruemye vychyslytelnye systemy: Osnovy y prylozheniya: monohrafiia” [Reconfigurable computing systems: Fundamentals and applications], Prosvita, Kyev, 280 p.

12. Palahyn A. V. and Yakovlev S.Yu. (2005), “Systemnaia yntehratsiia sredstv kompiuternoii tekhniky” [System integration of computer equipment], Universum-Vinnytsia, Vinnytsia, 680 p.

13. Khoroshevskiy V. H. and Kurnosov M. H. (2008), “Modelyrovaniye alhorytmov vlozheniya parallelnykh prohramm v struktury raspredelennykh vychyslytelnykh system” [Modeling of algorithms for nesting parallel programs in the structures of distributed computing systems], Modelyrovaniye-2008 (Simulation-2008): Trudy Mezhdunarodnoi nauchnoi konferentsyy, Yzd-vo YPMЭ ym. H.E. Pukhova, Kyev, T. 2. pp. 435-440.

14. Yasko M.M. (2010), “Paralelni ta rozpodileni obchyslennia” [Parallel and distributed computing], D.: RVV DNU, Dnepr, 111p.

д.т.н., проф. Ленков С.В., к.т.н., доц. Огневой А.В., Коваль Б. А., Присяжнюк В.В.  
АНАЛИЗ АРХИТЕКТУРНЫХ ОСОБЕННОСТЕЙ ВЫЧИСЛИТЕЛЬНЫХ СИСТЕМ С  
ПРОГРАММИРУЕМОЙ СТРУКТУРОЙ

*Статья посвящена построению мультиархитектурных кластерных и мультикластерных вычислительных систем. В работе проведен анализ вычислительных систем с программируемой структурой. Исследования показали, что современные распределенные ВС являются мультиархитектурными. В зависимости от уровня рассмотрения их функциональных структур, они могут выглядеть и как MISD, и как SIMD, и как MIMD системы. Основная функционально-структурная единица вычислительных ресурсов в системах рассматриваемого класса - это элементарная машина (ЭМ), которая представляет собой композицию из вычислительного модуля и системного устройства. Установлено, что ВС с программируемой структурой относятся к масштабируемым средствам обработки информации и допускают формирование конфигураций с массовым параллелизмом (Scalable Massively Parallel Architecture Computing Systems). Взаимодействие между ЭВМ осуществляется с помощью программной настройки сети связи. Структура ВС описывается графом, множество вершин которого - количество элементарных машины (или системных устройств, или локальных коммутаторов), а множество ребер - линии межмашинной связи. К структурам современных ВС предъявляется ряд требований: простота вложения параллельного алгоритма, удобство адресации элементарных машин, осуществимость принципа близкодействия и минимума задержек при межмашинных передачах информации, масштабируемость и крупномасштабные структуры, коммутация, живучесть, технологичность. Наиболее полно перечисленным выше требованиям удовлетворяют однородные структуры (то есть описываемые однородными графами). Такие структуры являются перспективными для формирования масштабируемых и крупномасштабных вычислительных систем (в частности, ВС с программируемой структурой). Для оценки производительности каналов связи между ВС использованы следующие показатели как пропускная способность и латентность, позволяющие осуществить достаточно полный анализ структурных возможностей ВС. В работе приведен пример построения вычислительного кластера, укомплектованного N вычислительными узлами, которые объединены сетью связи стандарта Gigabit Ethernet.*

*Ключевые слова: вычислительная система, элементарная машина, система обработки данных, структура, однородная структура, коммуникационные среды, кластер, мультикластер, алгоритм вложения.*

Prof. Lenkov S.V., Ph.D. Ognevy O.V., Koval B. A., Prysyzhnyuk V.V.  
ANALYSIS OF ARCHITECTURAL PECULIARITIES OF COMPUTER SYSTEMS WITH  
PROGRAMMABLE STRUCTURE

*The article is devoted to the construction of multi-architectural cluster and multi-cluster computing systems.*

*The paper analyzes computational systems with programmable logic. Research have shown that modern distributed OSs are multiarchitectural. Depending on the level of consideration of their functional structures, they may look both as MISD, both as SIMD, and as MIMD systems.*

*The main functional and structural unit of computing resources in the systems of the class under consideration is an elementary machine (EM), which is a composition of a computing module and a system device.*

*It has been established that OS systems with a programmable structure are scalable information processing tools and allow the formation of scalable massively parallel architecture computing systems.*

*Interaction between the computer is carried out through the software configuration of the communication network. The structure of the OS is described by a graph, the set of vertices of which is the number of elementary machines (or system devices, or local switches), and the set of edges - the line of interconnection.*

*The structure of the modern OS presents a number of requirements: the simplicity of embedding a parallel algorithm, the convenience of addressing elementary machines, the feasibility of the principle of near-the-go and the minimum delays in inter-machine transmissions of information, scalability and large-*

*scale structures, commutation, survivability, technology. The most complete requirements listed above satisfy the homogeneous structures (that is, they are described by homogeneous graphs). Such structures are promising for the formation of scalable and large-scale computing systems (in particular, OS with a programmable structure).*

*To assess the performance of the communication channels between the EM system used indicators of bandwidth and latency, to allow a fairly complete analysis of the structural capabilities of the OS.*

*In this paper, we present an example of constructing a computing cluster, complete with  $N$  computing nodes, which are connected by a Gigabit Ethernet connection.*

*Keywords: computing systems, elementary machines, data processing systems, homogeneous structures, communication environments, cluster and multiclaser, embedding algorithms.*