

# Evaluarea performanței software in sistemele distribuite

Stefan Adrian-Catalin

Grupa: I.S. 1 C

## Tema 8

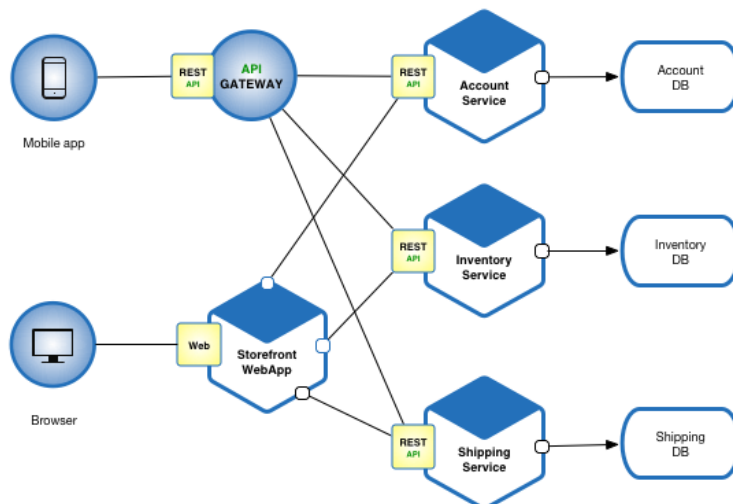
### Optimizarea performanței hardware-software (de ex. GPU-CUDA) pentru aplicații paralele obișnuite

#### 1. Introducere

Performanța software este un aspect esențial în dezvoltarea sistemelor moderne, mai ales când vine vorba de sisteme distribuite, care implică multiple resurse hardware și software interconectate. Acest raport analizează metodele de evaluare a performanței software, metricile utilizate și tehnicile specifice pentru analiza sistemelor distribuite.

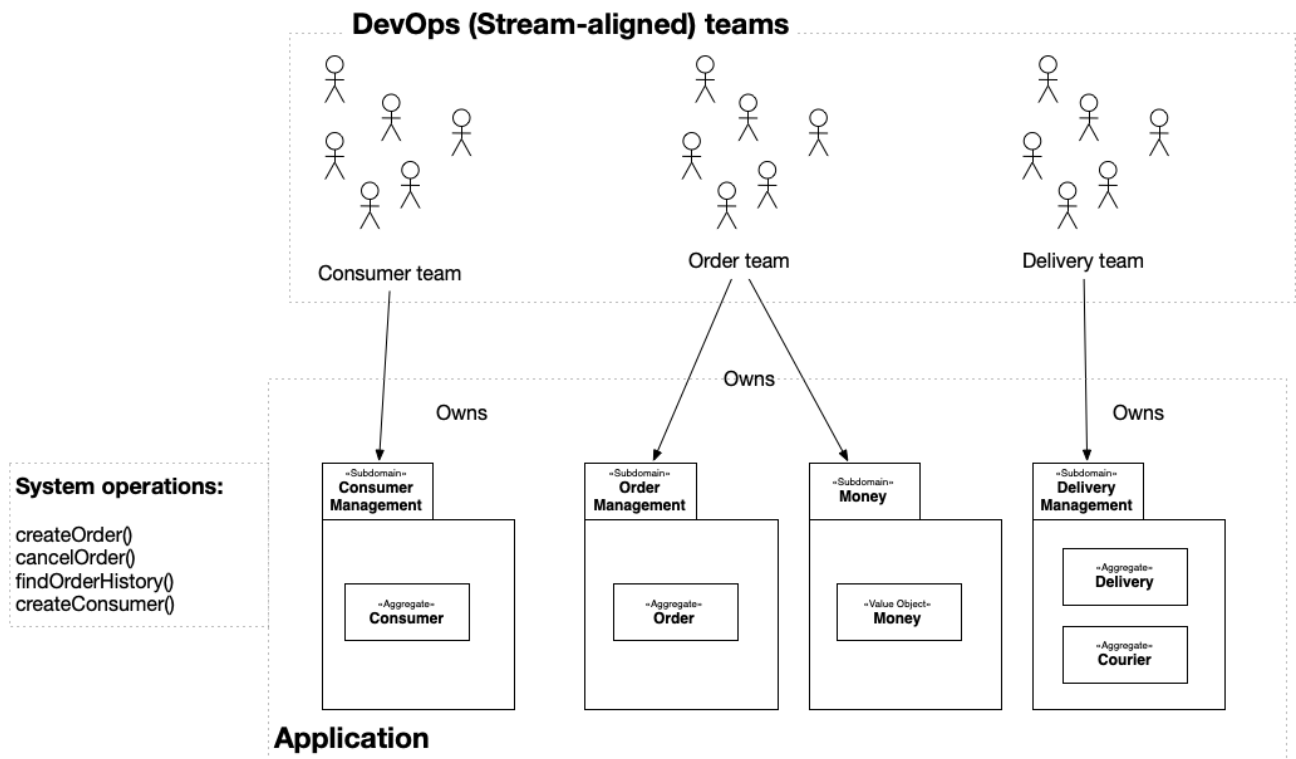
Performanța este crucială în aplicațiile critice, precum cele din domeniul financiar, medical sau al rețelelor de telecomunicații, unde latența redusă și eficiența resurselor sunt factori determinanți. Creșterea cerințelor pentru scalabilitate și eficiență necesită metode avansate de analiză și optimizare. În plus, dezvoltarea sistemelor bazate pe cloud și a arhitecturilor microserviciilor a adăugat un nou strat de complexitate în evaluarea performanței.

Performanța software nu depinde doar de viteza de execuție, ci și de gestionarea eficientă a resurselor hardware. Optimizarea performanței poate implica reducerea consumului de memorie, distribuirea eficientă a sarcinilor de procesare sau îmbunătățirea algoritmilor pentru a reduce latența.



Un exemplu concret este utilizarea Amazon Web Services (AWS) pentru rularea aplicațiilor la scară globală. Evaluarea performanței unei astfel de infrastructuri implică măsurarea latenței între centrele de date, analiza utilizării CPU și RAM și optimizarea fluxurilor de date între microservicii. De asemenea, serviciile de tip Content Delivery Network (CDN), cum ar fi Cloudflare sau Akamai, sunt utilizate pentru a reduce latența și a îmbunătăți experiența utilizatorilor finali.

Pentru a înțelege mai bine importanța performanței software, putem lua exemplul rețelelor 5G, unde latența trebuie să fie minimă pentru a permite comunicarea în timp real între dispozitive IoT (Internet of Things). În acest context, optimizarea software și distribuția eficientă a sarcinilor între noduri devin esențiale. În plus, inteligența artificială și învățarea automată sunt din ce în ce mai utilizate pentru predictibilitatea performanței și pentru ajustarea dinamică a resurselor în funcție de cerințele de utilizare.



## 2. Conceptul de performanță software

Performanța software se referă la modul în care un program sau un sistem răspunde cerințelor utilizatorilor în ceea ce privește viteza de execuție, consumul de resurse și scalabilitate. O aplicație performantă trebuie să ofere timpi de răspuns rapizi, să utilizeze eficient resursele disponibile și să fie scalabilă, astfel încât să poată gestiona creșterea numărului de utilizatori sau a volumului de date procesate. Evaluarea acesteia presupune identificarea factorilor critici care influențează eficiența unei aplicații.

Factorii care influențează performanța includ:

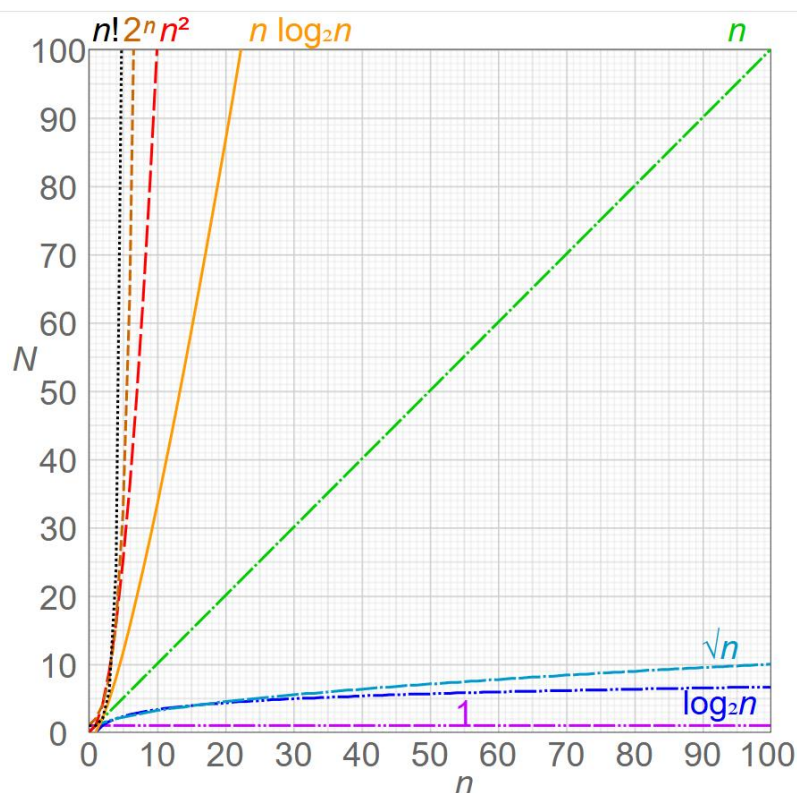
- **Eficiența codului sursă** – algoritmi folosiți și optimizările implementate au un impact direct asupra performanței. De exemplu, un algoritm  $O(n \log n)$  este considerabil mai eficient decât unul  $O(n^2)$  pentru sortarea unor volume mari de date.
- **Hardware-ul pe care rulează aplicația** – performanța unui program poate fi limitată de viteza procesorului (CPU), capacitatea memoriei RAM, viteza de citire/scriere a SSD-ului sau lățimea de bandă a rețelei.

- **Sistemul de operare și mediul de execuție** – gestionarea eficientă a proceselor, firelor de execuție (threads) și resurselor partajate are un impact semnificativ. De exemplu, în Linux, scheduler-ul CPU decide care procese primesc mai mult timp de execuție.
- **Interacțiunea cu baza de date și sistemele de stocare** – citirea și scrierea datelor poate deveni un bottleneck major. Utilizarea indexării eficiente, caching-ului și tehnicilor de optimizare SQL poate îmbunătăți semnificativ performanța.

Un exemplu de aplicație cu cerințe ridicate de performanță este un sistem de procesare a tranzacțiilor financiare în timp real.

Un alt exemplu este optimizarea unui server web care trebuie să gestioneze un număr mare de utilizatori simultan. Prin utilizarea unui balansor de sarcină (load balancer) și caching inteligent, se poate reduce timpul de răspuns și optimiza fluxul de date.

- Grafic despre performanța algoritmilor în funcție de complexitatea lor:

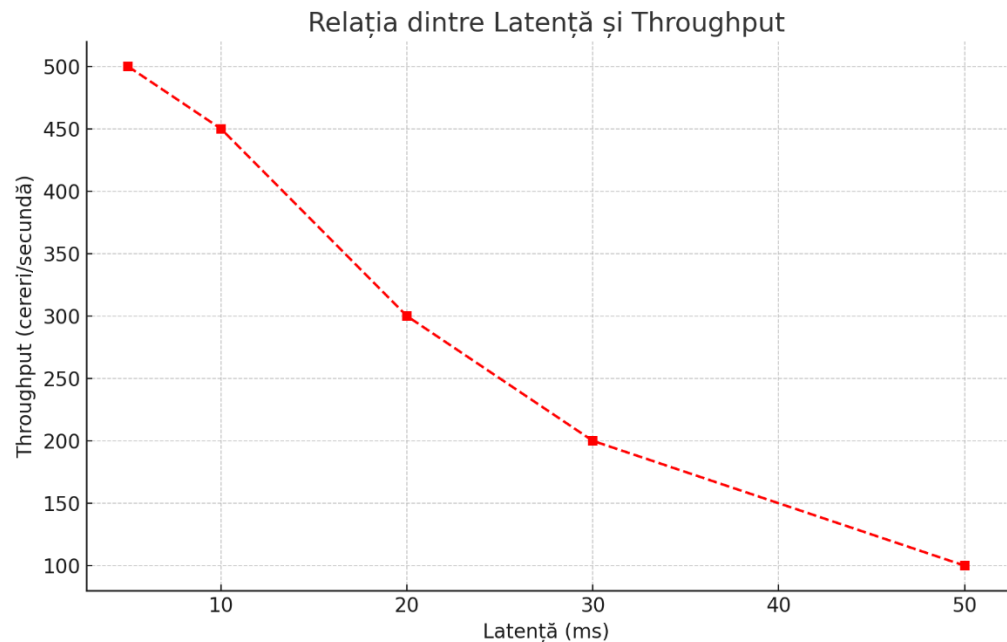


### 3. Metrice pentru evaluarea performanței

Evaluarea performanței software implică utilizarea mai multor metrice cantitative care oferă informații despre eficiența și comportamentul unui sistem. Aceste metrice ajută la identificarea problemelor și la optimizarea performanței prin ajustarea codului, resurselor hardware și configurațiilor software:

- **Timp de execuție (Execution Time):**
  - Reprezintă durata totală necesară pentru finalizarea unei sarcini.
  - Un timp de execuție mai mic indică o performanță mai bună.
  - Se măsoară în milisecunde (ms), secunde (s) sau nanosecunde (ns).
  - Ex: Într-un sistem de plăți online, timpul de execuție trebuie să fie minim pentru a nu afecta experiența utilizatorului.
- **Latența (Latency):**
  - Reprezintă timpul de întârziere între inițierea unei sarcini și finalizarea acesteia.

- Poate fi afectată de suprasarcina rețelei, viteza procesorului sau latența accesului la baze de date.
- Ex: Într-un joc online, latența mare (lag) afectează sincronizarea acțiunilor între jucători.
- **Throughput:**
  - Măsoară numărul de operații procesate într-un anumit interval de timp.
  - Cu cât throughput-ul este mai mare, cu atât un sistem poate gestiona mai multe cereri.
  - Ex: Un server web performant trebuie să poată gestiona milioane de cereri HTTP pe secundă.



- **Utilizarea resurselor:**
  - Analizează cât de eficient sunt folosite CPU, RAM, stocarea și rețeaua.
  - Un sistem bine optimizat folosește resursele minim necesare pentru a-și îndeplini sarcinile.
  - Ex: Un browser web eficient consumă mai puțină memorie RAM și optimizează încărcarea paginilor.
- **Scalabilitate:**
  - Măsoară capacitatea unui sistem de a gestiona creșterea numărului de utilizatori și a volumului de date.
  - Un sistem scalabil își poate mări sau reduce resursele în funcție de cerere.
  - Ex: Serviciile cloud precum AWS și Azure oferă scalabilitate automată pentru aplicațiile mari.

## 4. Metode de evaluare a performanței software

### 4.1 Benchmarking

Benchmarking-ul presupune testarea performanței unui sistem prin compararea rezultatelor obținute cu un set de referințe standard. Această metodă este utilizată atât pentru hardware (procesoare, GPU, SSD, RAM), cât și pentru software (aplicații, baze de date, sisteme de operare).

Exemple:

- SPEC CPU
- LINPACK pentru calcul numeric

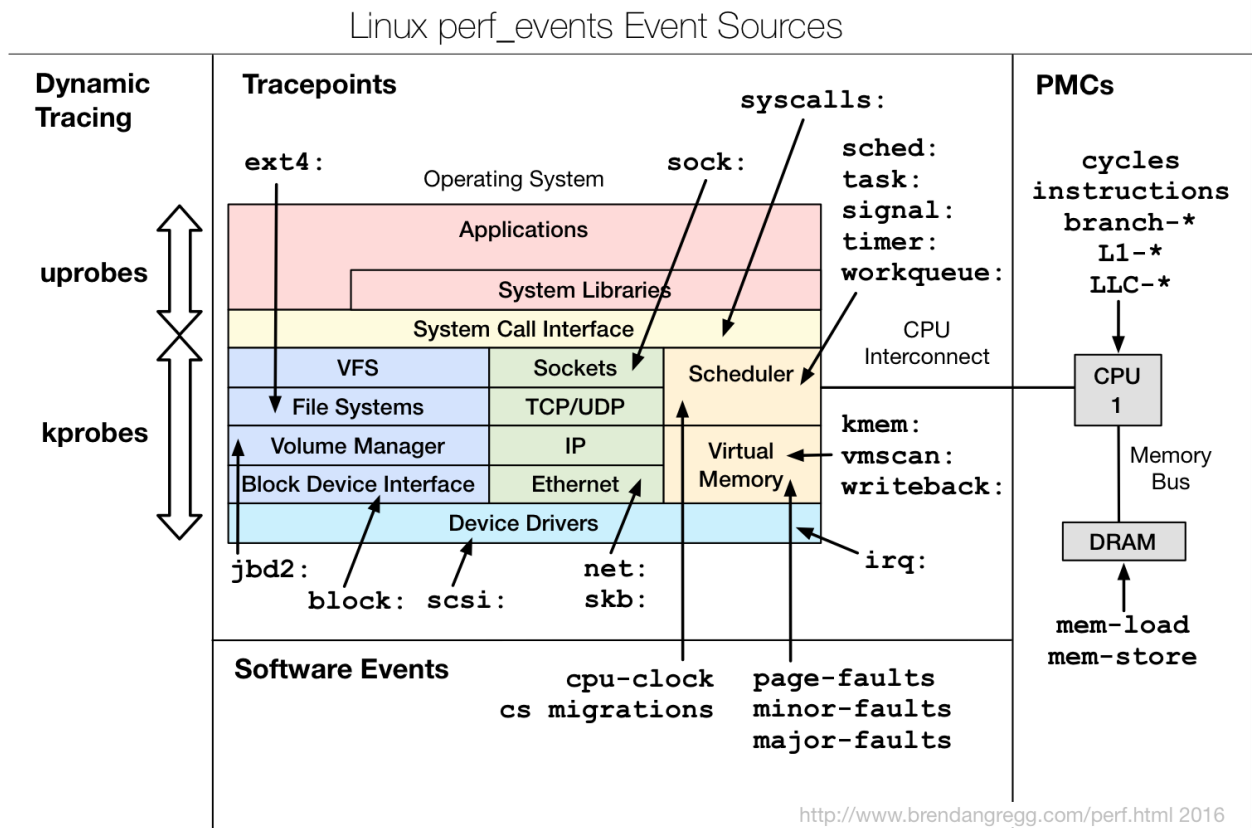
- TPC pentru baze de date

Benchmarking-ul este utilizat frecvent în evaluarea performanței hardware, cum ar fi compararea procesoarelor Intel și AMD. Un exemplu practic este testarea cu Geekbench sau PassMark.

## 4.2 Profilarea performanței

Profilarea implică analiza detaliată a codului pentru a identifica punctele critice care afectează performanța. Instrumente utilizate:

- Gprof (GNU Profiler)
- Perf (Linux Performance Profiler)
- Valgrind



## 4.3 Simulări de performanță

Simulările permit testarea comportamentului unui sistem în diferite scenarii fără a afecta mediul de producție. Exemple:

- OMNeT++ pentru rețele
- SimPy pentru modelarea sistemelor cu cozi de așteptare

## 5. Evaluarea performanței în sisteme distribuite

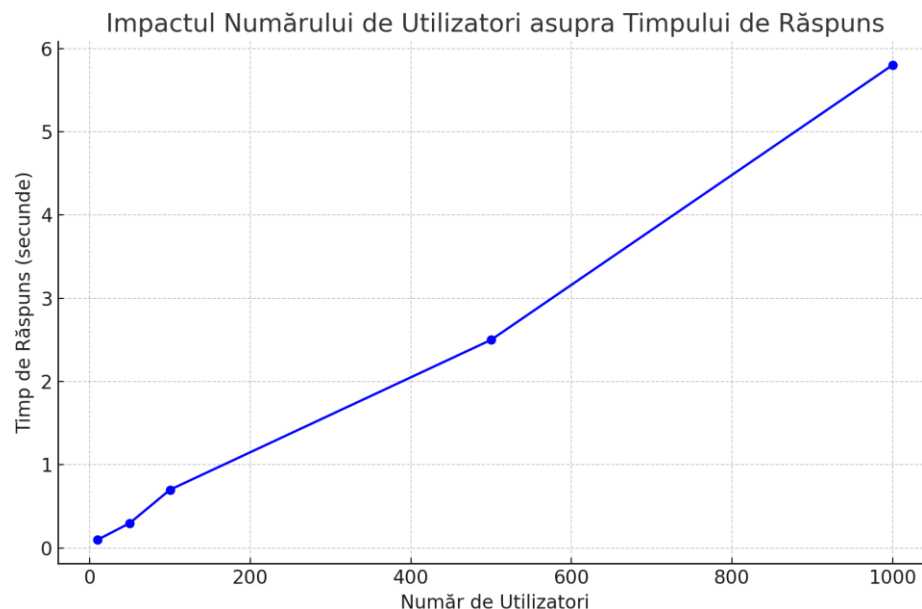
Sistemele distribuite adaugă un nivel suplimentar de complexitate în evaluarea performanței. Acestea sunt caracterizate prin existența mai multor noduri care trebuie să colaboreze pentru a executa sarcini comune, iar acest lucru implică o serie de factori specifici care influențează performanța globală.

Factorii specifici de analizat includ:

- **Comunicarea între noduri** (latența și overhead de rețea): reprezintă unul dintre principalii factori care afectează performanța unui sistem distribuit. Latența dintre noduri este influențată de distanța fizică dintre servere, de viteza conexiunilor de rețea și de mecanismele de retransmisie a datelor în cazul erorilor. Un alt aspect important este overhead-ul de rețea, care

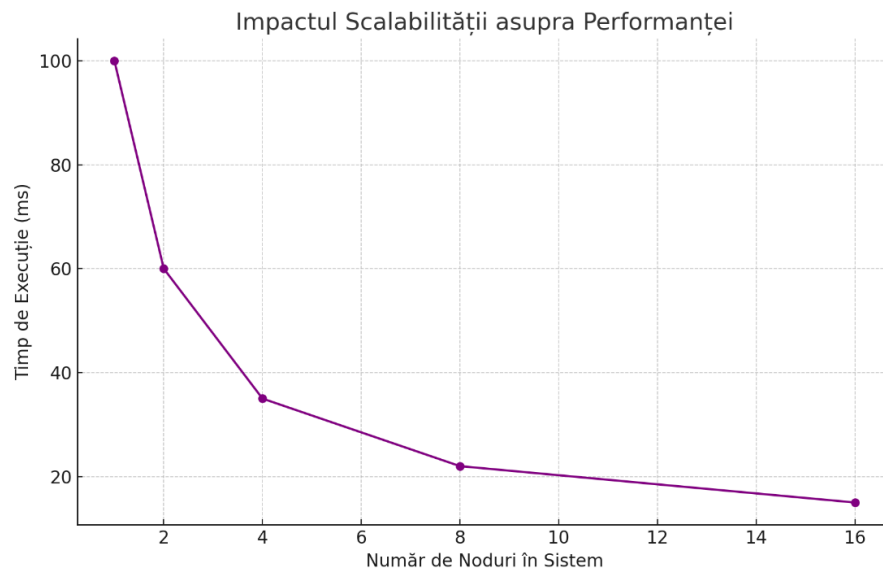
se referă la resursele suplimentare consumate pentru a gestiona traficul dintre noduri, incluzând protocoalele de comunicare și mecanismele de sincronizare. În sistemele mari, cum sunt centrele de date ale marilor companii de cloud computing, optimizarea rețelei interne prin rutare eficientă și utilizarea unor algoritmi avansați de transmisie a datelor poate reduce latența și îmbunătăți performanța generală.

- **Consistența datelor** (impactul mecanismelor de replicare): este un alt factor esențial în evaluarea performanței sistemelor distribuite. În astfel de sisteme, datele sunt replicate pe mai multe noduri pentru a asigura disponibilitatea și toleranța la erori. Totuși, această replicare poate introduce întârzieri și conflicte în actualizarea informațiilor. Mecanismele de menținere a consistenței, cum ar fi modelele de consistență strictă sau eventuală, influențează performanța generală a sistemului. În baze de date distribuite, alegerea între un model de consistență puternică (strong consistency), în care toate nodurile trebuie să fie sincronizate instantaneu, și un model mai flexibil, care permite un anumit grad de inconsistență temporară (eventual consistency), depinde de cerințele aplicației și de toleranța la latență.
- **Distribuția sarcinilor** (balanța între workload-ul nodurilor): joacă un rol crucial în menținerea performanței sistemului. Un sistem distribuit eficient trebuie să echilibreze workload-ul astfel încât niciun nod să nu fie supraîncărcat, în timp ce altele rămân inactive. Tehnicile de load balancing sunt utilizate pentru a redirecționa solicitările utilizatorilor către nodurile care au cea mai mică încărcare, astfel încât să fie maximizată utilizarea resurselor. Load balancing-ul poate fi realizat prin algoritmi statici, unde sarcinile sunt distribuite în mod prestabilit, sau prin algoritmi dinamici, care analizează în timp real starea fiecărui nod și redistribuie sarcinile în funcție de utilizarea resurselor.



Evaluarea performanței într-un sistem distribuit poate implica analiza traficului de rețea, monitorizarea consumului de resurse pe fiecare nod și testarea mecanismelor de load balancing. O metodă frecvent utilizată este colectarea de metrice în timp real, cum ar fi timpul de răspuns, rata de eroare și utilizarea CPU, pentru a identifica punctele slabe ale sistemului. De asemenea, simulările de performanță sunt utilizate pentru a testa comportamentul sistemului în diferite scenarii, cum ar fi creșteri bruște ale numărului de utilizatori sau defectarea unor noduri. Analiza acestor factori permite optimizarea

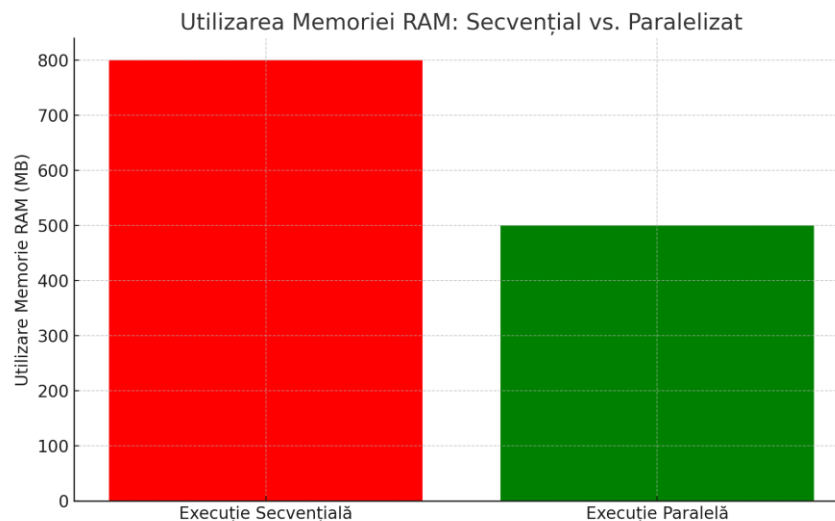
sistemelor distribuite astfel încât acestea să funcționeze eficient și să răspundă cerințelor de scalabilitate și fiabilitate.

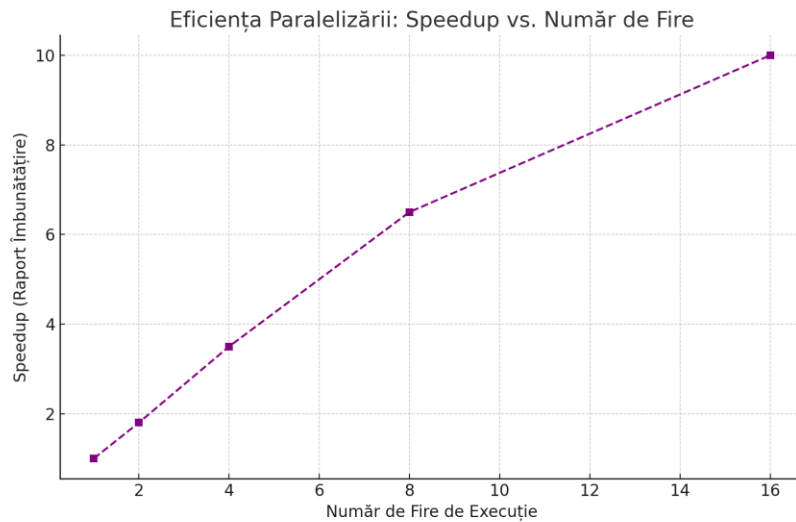


## 6. Metode de optimizare a performanței

Optimizarea performanței software este esențială pentru îmbunătățirea timpilor de execuție, reducerea consumului de resurse și scalarea eficientă a aplicațiilor. În special pentru sistemele cu cerințe mari de procesare, cum sunt aplicațiile științifice, inteligența artificială și simulările complexe, utilizarea tehnicilor avansate de optimizare este indispensabilă. Printre cele mai eficiente metode se numără paralelizarea prin threading, distribuirea sarcinilor între procese și accelerarea hardware cu GPU.

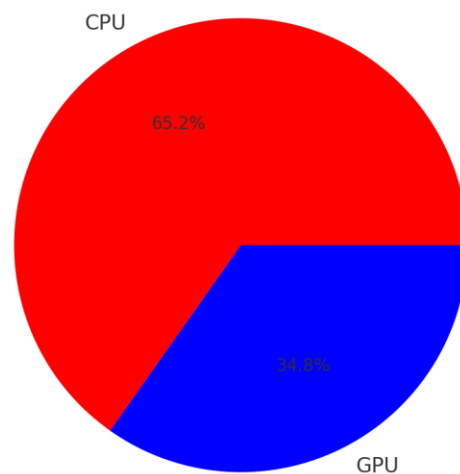
- **Paralelizare** prin threading (Pthreads, Java Threads, OpenMP): este una dintre cele mai utilizate tehnici pentru optimizarea software, permițând executarea simultană a mai multor sarcini. Aceasta presupune împărțirea unei operațiuni în mai multe fire de execuție (threads), care rulează în paralel pe un procesor multicore. Un exemplu de implementare este utilizarea Pthreads (POSIX Threads) pentru programarea multi-threaded în C și C++, Java Threads pentru aplicațiile Java și OpenMP pentru paralelizarea automată a codului. Prin folosirea acestor metode, un program poate beneficia de o reducere semnificativă a timpului de execuție, deoarece mai multe operații sunt procesate simultan în loc să fie executate secvențial.





- **Distribuirea sarcinilor** folosind MPI: este esențială în aplicațiile de calcul paralel distribuit, unde mai multe noduri colaborează pentru a rezolva probleme complexe. MPI permite comunicarea între procese aflate pe sisteme diferite, fiind utilizat în simulările științifice, modelarea climatică și supercomputing. Avantajul principal al utilizării MPI constă în scalabilitatea ridicată, deoarece aplicațiile pot fi distribuite pe clustere mari de calculatoare, fiecare procesor contribuind la rezolvarea unei părți a problemei.
- **Accelerare hardware** cu GPU (CUDA, OpenCL): este o metodă modernă și extrem de eficientă pentru optimizarea performanței. Spre deosebire de procesoarele tradiționale (CPU), care sunt optimizate pentru sarcini secvențiale, plăcile grafice (GPU) sunt capabile să ruleze mii de fire de execuție în paralel, făcându-le ideale pentru aplicații de deep learning, procesare video, simulări și randare grafică. NVIDIA CUDA este un framework utilizat pentru programarea GPU-urilor NVIDIA, în timp ce OpenCL este o soluție mai generală, compatibilă cu diverse arhitecturi hardware, inclusiv AMD și Intel. Folosirea acestor tehnologii poate crește semnificativ viteza de execuție pentru algoritmi care necesită procesare intensivă.

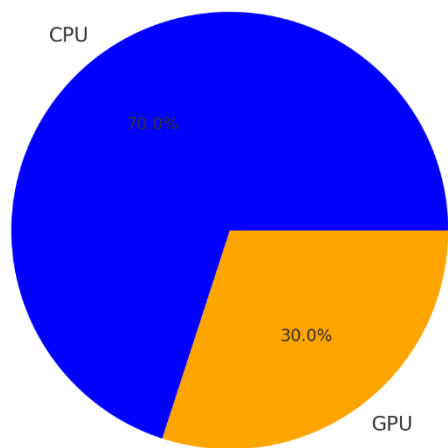
Consum de Energie CPU vs. GPU în Execuție Paralelă



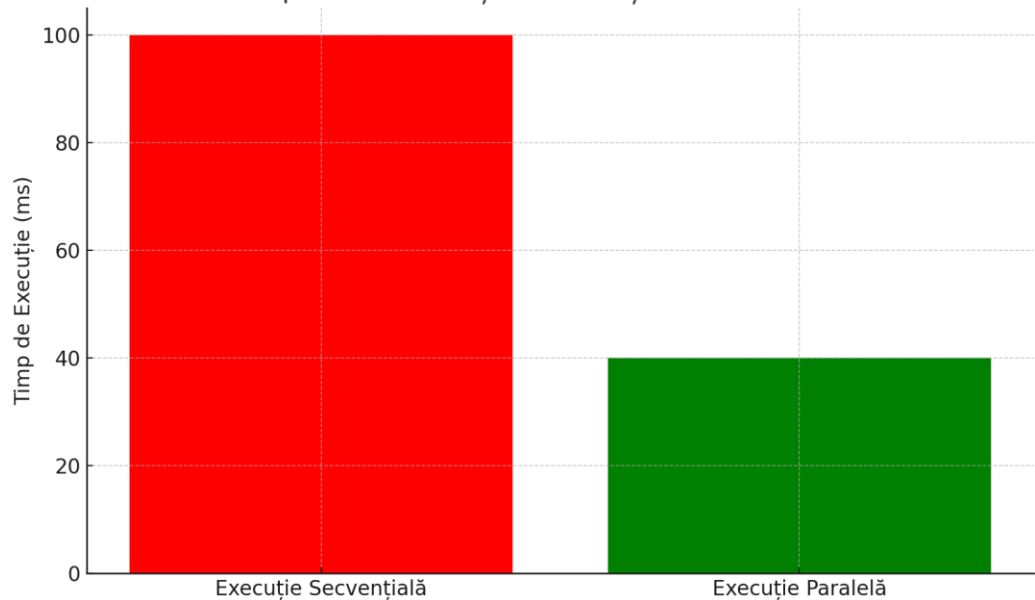


Această secțiune reprezintă un punct de plecare pentru proiectul final, care va analiza optimizarea performanței software utilizând paralelizare și accelerare hardware. Proiectul va implementa și compara diferite metode de paralelizare, precum fire de execuție, MPI, OpenMP și CUDA, evidențiind impactul fiecărei tehnici asupra performanței generale a aplicației.

Utilizarea Resurselor în Procesare Paralelă



Compararea Execuției Secvențiale vs. Paralelizate



7. Concluzii

Evaluarea performanței software și a sistemelor distribuite este esențială pentru optimizarea aplicațiilor moderne și pentru asigurarea unei experiențe eficiente utilizatorilor. Într-un mediu tehnologic aflat într-o continuă dezvoltare, sistemele trebuie să fie rapide, scalabile și eficiente din punct de vedere al utilizării resurselor. O evaluare corectă a performanței permite detectarea problemelor, identificarea gâtuirilor de procesare și aplicarea unor tehnici de optimizare care să îmbunătățească funcționarea aplicațiilor.

Utilizarea unor metode precum benchmarking-ul, profilarea performanței și simulările permite o analiză detaliată asupra modului în care un sistem gestionează resursele disponibile și răspunde la diferite scenarii de încărcare. Benchmarking-ul oferă o comparație standardizată a performanței hardware și software, permițând dezvoltatorilor să aleagă configurațiile optime. Profilarea performanței ajută la identificarea secțiunilor lente ale codului, permițând optimizarea algoritmilor și gestionarea mai eficientă a memoriei. Simulările sunt esențiale pentru anticiparea comportamentului aplicațiilor în medii distribuite, fără a afecta utilizatorii reali.

Sistemele distribuite prezintă provocări suplimentare, cum ar fi latența rețelei, consistența datelor și distribuirea optimă a sarcinilor între noduri. Evaluarea performanței în astfel de sisteme necesită tehnici avansate de monitorizare și analiza traficului, pentru a asigura un echilibru optim între scalabilitate și eficiență. Prin utilizarea unor algoritmi de load balancing și a tehnologiilor de caching, se poate reduce timpul de răspuns și îmbunătăți disponibilitatea serviciilor.

În proiectul final, se va analiza o aplicație specifică, comparând execuția sa în mod secvențial cu variante optimizate prin paralelizare și accelerare hardware. Se vor implementa și compara diferite metode de optimizare, precum fire de execuție (threads), distribuția sarcinilor folosind MPI și accelerarea GPU cu CUDA, pentru a observa impactul fiecărei tehnici asupra performanței. Prin profilarea codului și analiza metricilor de performanță, se va demonstra eficiența tehnicilor de paralelizare și se vor evidenția avantajele utilizării acestora în aplicațiile cu cerințe ridicate de procesare.

Optimizarea performanței software nu este un proces singular, ci un proces continuu, influențat de evoluția hardware-ului și a noilor tehnologii. Alegerea celor mai potrivite tehnici de optimizare depinde de natura aplicației, resursele disponibile și cerințele de performanță, iar o analiză atentă a performanței este cheia dezvoltării unor sisteme rapide, eficiente și scalabile.

## 8. Referințe

- <https://www.datadoghq.com/blog/performance-metrics/>
- <https://www.brendangregg.com/perf.html>
- <https://developer.nvidia.com/cuda-zone>
- <https://www.researchgate.net/profile/R-Boutaba>
- <https://aws.amazon.com/latency/>
- <https://microservices.io/patterns/microservices.html>
- <https://www.spec.org/>
- <https://www.cpubenchmark.net/>