



PARALELNO PRETRAŽIVANJE GRAFA NA ARHITEKTURAMA SA
DISTRIBUIRANOM I DELJENOM MEMORIJOM

PARALLEL BREADTH-FIRST SEARCH GRAPH TRAVERSAL ON COMPUTER
ARCHITECTURES WITH DISTRIBUTED AND SHARED MEMORY

Stefan Aleksić, Fakultet tehničkih nauka, Novi Sad

Oblast – ELEKTROTEHNIKA I RAČUNARSTVO

Kratak sadržaj – Glavni fokus ovog rada jeste pokušaj unapređenja algoritama pretrage grafova, pogotovo za grafove sa velikim brojem čvorova. Ovaj rad pruža informacije vezane za algoritme paralelne pretrage grafa, na računarskim arhitekturama sa distribuiranom i deljenom memorijom. Svaka od predstavljenih implementacija analizirana je za različite parametre, kako bi se testirale njene performanse. Zaključak samog istraživanja jeste da paralelna pretraga grafa može biti performantna na računarskim arhitekturama sa deljenom memorijom, dok implementacija na arhitekturama sa distribuiranom memorijom jedino ima smisla za grafove koji ne mogu da se smeste u radnoj memoriji jednog procesora.

Ključne reči: pretraživanje grafa, distribuirani sistemi, paralelni sistemi

Abstract – The main focus of this research paper is an attempt to improve upon graph traversal algorithms, especially in the cases of large graphs. This paper provides information about parallel graph traversal algorithms, as well as their implementation on distributed and parallel computer systems. Every type of implementation is analyzed through different parameters to test its achieved performance. The outcome of this research concluded that parallel graph traversal is performant in the case of computer architectures with shared memory, unfortunately only case where it becomes useful for architectures with distributed memory is when graph size becomes too large to fit on one node.

Keywords: graph traversal, distributed systems, parallel systems

1. UVOD

Upotreba apstrakcije uz pomoć grafova za analizu i razumevanje raznih vrsta podataka dobija sve veći značaj. Neki od primera podataka koji mogu da se apstrahuju koristeći grafove podrazumevaju: podatke o interakcijama na društvenim mrežama, podatke bankarskih transakcija, podatke o preporuci raznih reklama korisnicima aplikacija na osnovu njihovih interakcija, komunikacionih podataka poput elektronske pošte i telefonskih mreža, podatke bioloških sistema i različitih oblika relacionih podataka generalno.

NAPOMENA:

Ovaj rad proistekao je iz master rada čiji mentor je bio dr Veljko Petrović, docent.

Kada se govori o veštačkoj inteligenciji, apsolutno je neophodno uvesti neku vrstu grafa i primenjivati raznovrsne algoritme nad njim. Zajednički problemi u matematičkoj oblasti teorije grafova i u oblastima primene uključuju identifikaciju i rangiranje važnih entiteta, otkrivanje anomalija u obrascima ili iznenadnih promena u mrežama, pronalaženje čvrsto povezanih klastera entiteta, itd. Rešenja ovih problema obično uključuju klasične algoritme nad grafovima, kao što su: prebacivanje grafa u strukturu stabla (uklanjanje ciklusa), pronalaženje najkraćih putanja, pronalaženje dvopovezanih komponenti, uparivanja, proračune zasnovane na protoku itd.

Kako bi se zadovoljile potrebe teorijske analize grafova za nove aplikacije koje zahtevaju rad sa strukturama velikih skupova podataka, od suštinskog je značaja ubrzati osnovne probleme grafova koristeći aktuelne paralelne sisteme [1].

U ovom radu će se detaljnije razmatrati problem pronalaženja najkraćeg puta, odnosno najmanjeg broja skokova, između dva proizvoljna čvora u povezanom, neusmerenom grafu i to primenom algoritma za pretragu grafa po širini, poznatog kao *breadth-first traversal*. Biće prikazani primeri već postojećih ideja paralelizacije ovog algoritma, kao i implementacija koje se zasnivaju na pomenutim idejama, u okviru arhitektura sa distribuiranom i deljenom memorijom. Takođe, biće dat i prikaz rezultata dobijenih primenjujući ove implementacije na nasumično generisanim grafovima, kao i nad grafovima koji su specifično pogodni za ovu vrstu obrade.

2. PRETHODNA ISTRAŽIVANJA

Najpoznatiji algoritmi za pretraživanje grafova su algoritam **obilaska grafa po širini** (engl. *breadth-first search*) i algoritam **obilaska grafa po dubini** (engl. *depth-first search*). Ideja ovih algoritama je da pronađu najkraći put od izvornog čvora do proizvoljnog ciljnog čvora u grafu. Od interesa za ovaj rad jeste algoritam *BFS*. Pseudokod višeizvornog *BFS* algoritma za obilazak neusmerenog grafa je prikazan u okviru [Pseudokod 1].

Kako bi se podelio posao obrade na entitete u paralelnom okruženju, neophodno je obaviti neku vrstu particionisanja grafa kao celine. Ideja na koji način je moguće izvršiti ovakvo particionisanje je predstavljena u okviru rada [2]. Predloženi algoritam je *BFS* algoritam sinhronizovan po nivoima koji napreduje nivo po nivo, počevši od izvornog čvora, gde je nivo temena definisan kao njegov graf udaljenosti od izvora. **Razmatrani su samo usmereni grafovi.**

Алгоритам Вишеизворни обилазак графа по шириниУлаз: Граф G , skup izvornih temena S Изаз: Листа путања за свако теме P

```

1: procedure multisourceBreadthFirstSearch( $G, S$ )
2:   create queue  $Q$ 
3:   create list  $P$  where  $P[v_i] \leftarrow \emptyset$ 
4:   for all sources  $s$  in  $S$  do
5:     enqueue( $Q, s$ )
6:      $P[s] \leftarrow s$ 
7:   while  $Q$  is not empty do
8:      $w \leftarrow$  dequeue( $Q$ )
9:     mark  $w$ 
10:    for all edges  $e$  in adjacentEdges( $G, w$ ) do
11:       $x \leftarrow$  adjacentVertex( $w, e$ )
12:      if  $x$  is not marked then
13:         $P[x] \leftarrow$  append( $x, P[w]$ )
14:        enqueue( $Q, x$ )

```

Pseudokod 1 Вишеизворни обилазак neusmerenog графа по ширини

Jednodimenziono particionisanje графа podrazumeva raspodelu temena графа, tako da svaki чвор i svi potezi koji proizilaze iz tog чвора pripadaju samo jednom procesoru. Skup temena koji pripadaju jednom procesoru q će dalje biti nazivan skup *lokalnih temena*. U nastavku je prikazana ilustracija jednodimenzionog particionisanja temena графа na P particija korišćenjem matrice susedstva A , koja je podeljena tako da su lokalna temena za neki procesor q kontinualna.

$$\begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_p \end{bmatrix}$$

Индекси подматрица суседства A_i представљају индикаторе процесора којима су дodelјене. Potezi који произилазе од темена v_j формирају листу потега, која је представљена листом темена у оквиру j -ог реда матрице суседства A . У оквиру слике [Pseudokod 2] је дат pseudokod алгоритма за дистрибуирање претраживање графа са једnodimenzionim particionisanjem, počevši од чвора v_s . У оквиру алгоритма, свако теме v_i бива означено својим нивоом, у оквиру низа $L_{v_s}(v_i)$, што означава растојање чвора v_s од чвора v_i . Низ L_{v_s} је дистрибуиран у складу са дистрибуцијом темена, тако да процесор P_i поседује растојања почетног темена v_s до својих локалних темена $\{v_{(i,0)}, v_{(i,1)}, \dots, v_{(i,m)}\}$ где је $m = \frac{|V|}{|P|}$.

Алгоритам Паралелни BFS са једnodimenzionim particionisanjemУлаз: Граф $G(V, E)$, полазно теме s Изаз: Низ d чији елементи представљају број скокова од почетног темена s до темена $v_i \in V$

```

1: procedure bfsDistributedSearch( $G(V, E), s$ )
2:   Initialize  $L_s(v) \leftarrow \begin{cases} \infty, & v \neq s \\ 0, & v = s \end{cases}$ 
3:   for level  $\leftarrow 0$  to  $\infty$  do
4:      $F \leftarrow \{v | L_s(v) = \text{level}\}$ , skup локалних темена за процес са нивоом level
5:     if  $F$  is empty за све процесне then
6:       Терминирај спољашњу петљу свим процесима
7:      $N \leftarrow$  {суседна темена из skupa  $F$  (нису нужно локална темена за тренутни процес)}
8:     for all processes  $q$  do
9:        $N_q \leftarrow$  {темена из  $N$  која припадају процесу  $q$ }
10:      Send  $N_q$  to process  $q$ 
11:      Receive  $\tilde{N}_q$  from process  $q$ 
12:       $\tilde{N} \leftarrow \cup_q \tilde{N}_q$  ( $\tilde{N}_q$  је могуће да има дубликате)
13:      for  $v \in \tilde{N}$  and  $L_s(v) = \infty$  do
14:         $L_s(v) \leftarrow \text{level} + 1$ 
15:   return  $L_s$ 

```

Pseudokod 2: Paralelni BFS algoritam sa јednodimenzionim particionisanjem

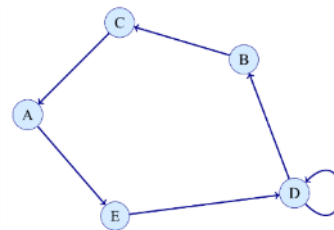
Алгоритам функционише на следећи начин. На сваком нивоу (дубине графа), сваки од процесора има skup F који представља skup локалних темена који се налазе на растојању *level* од почетног чвора. Листа потега која се добија од svakог од темена из skupa F бива унификована у skup потега N који представља skup суседних темена. Поједина темена у оквиру skupa N неће бити из skupa локалних темена који је дodelјен процесору. За ова

temena, врши се razmena sa procesorom q , kome се temena koja njemu pripadaju šalju, a temena koja pripadaju trenutnom procesoru primaju od istog procesora q . Svaki od procesora врши ovu razmenu temena i formira konačan skup \tilde{N} , skup temena do kojih се stiglo u trenutnom nivou, a koji су из skupa локалних темена за trenutni procesor.

3. TEORIJSKE OSNOVE

Граф (engl. *graph*), posmatrano iz ugla diskretne matematike, odnosno preciznije **теорије граfoва** (engl. *graph theory*), представља структуру која се састоји од skupa objekata koji међу собом могу имати специфичну povezanost, odnosno relaciju. Sami objekti koji се povezuju се називају **теменима** (engl. *vertex*), dok су veze između tih objekata **пoteзи, odnosno иvice** (engl. *edge*) [3]. Posmatrano iz ugla računarstva, граф представља apstraktnu структуру podataka koja služi за samo modelovanje граfoва из matematičke oblasti теорије граfoва на računaru [4].

Путања (engl. *path*) представља секвенцу или низ потега, takvih да се одредишно теме i -ог потега poklapa са polaznim чвором $i+1$ -ог потега [4]. Kod neusmerenih граfoва će zbog njihove prirode važiti да ukoliko postoji put od чвора A до чвора B , postojaće i put od чвора B до чвора A , dok kod usmerenih граfoва ovo неће бити правило. Put у графу otvara novу taksonomiju у свету граfoва, на **povezane** (engl. *connected*) i **nepovezane** (engl. *disconnected*) граfoве. Povezan граф представља структуру графа koja је povezana у smislu *topološkog prostora*, odnosno govori да postoji put od proizvoljnog izvornog temena до proizvoljnog одредишног темена. Nasuprot ovome, ukoliko се за граф не може tvrditi да postoji путања između два proizvoljna temena, takav граф је nepovezan [4]. **Најкраћи пут** (engl. *shortest path*) između два чвора у графу, odnosno путању са најmanjom kardinalnošću од свих mogućih путања između два задата чвора у графу [5].



Slika 1: Primer usmerenog графа

Kada би се граф са слике [Slika 1] predstavio matrično, dobila би се slabo posednuta matriца, te је najpogodniji način predstavjanja граfoва koji nemaju veliki broj poteга, korišćenjem nazubljenih matriца. U nastavku sledi postupak kako граф са слике [Slika 1] превести у nazubljenu matriцу susedstva.

$$A = [a_{ij}]_{n \times n}, \quad a_{ij} = \begin{cases} 1, & \text{if } e_{ij} = (v_i, v_j) \in E \\ 0, & \text{if } e_{ij} = (v_i, v_j) \notin E \end{cases}$$

Sada, ukoliko би се umesto 1 на odgovarajućoj poziciji у matriци iskoristili identifikator samog чвора ka kome postoji poteг, dobija се matriца koja izgleda ovako.

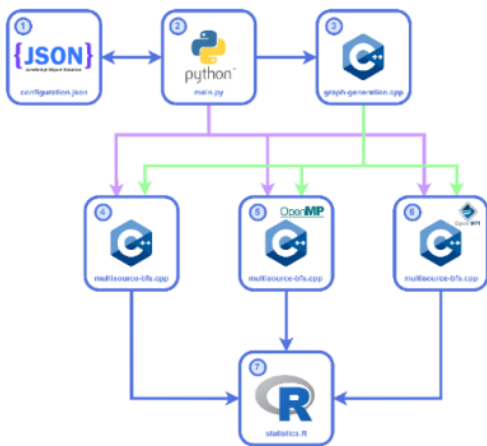
$$A = \begin{bmatrix} & B & & & & & \\ & C & & & & & \\ & D & & & & & \\ & E & & & & & \\ A & & E & & & & \end{bmatrix}$$

U okviru implementacije će se upravo koristiti ovakav način reprezentacije, kako bi se posao obrade mogao pravilno rasporediti na entitete za obradu.

4. IMPLEMENTACIJA

Implementacija je u radu realizovana kako za sekvencijalnu obradu koja će biti referentna, tako i za obradu u okviru distribuiranih arhitektura, kao i arhitektura sa deljenom memorijom. Ideja paralelne, odnosno distribuirane obrade je data u okviru [Pseudokod 2]. Implementacija samih algoritama je odrađena u okviru programskog jezika C++, dok su za paralelno, odnosno distribuirano programiranje iskorišćeni alati *OpenMP* i *OpenMPI*, respektivno.

Kako bi se mogla testirati obrada različitih implementacija, kreiran je odgovarajući tok podataka (eng. *pipeline*) koji generiše rezultate, vrši obradu i na kraju analizira dobijene rezultate. Ovaj tok podataka je prikazan na slici [Slika 2].



Slika 2: Arhitektura toka podataka

Sam tok podataka teče od (1) *JSON* konfiguracionog fajla u kome su postavljani svi neophodni parametri za pokretanje *Python* skripte (2). Skripta pokreće generisanje nasumičnih, neusmerenih, povezanih grafova (3), a nakon toga svaka od implementacija (4, 5, 6) vrši njihovu obradu. Na kraju, uz pomoć (7) *R* skripti se vrši analiza performansi svake od implementacija. Sama definicija metode koja generiše nasumičan, neusmeren, povezan graf je prikazana na slici [Slika 3].

Što se tiče paralelne implementacije *BFS* algoritma za arhitekture sa deljenom memorijom, ona je prikazana u okviru slike [Slika 4]. U okviru implementacije se mogu odvojiti regioni od interesa. U okviru linija 46 – 60 se izvršava inicijalizacija svih promenljivih koje se koriste pri obradi i ovaj segment je sekvencijalni deo, odnosno deo koji izvršava glavna master nit. Nakon toga sledi paralelni region u okviru linija 61 – 125, gde se vrši inicijalizacija privatnih promenljivih, a kasnije se otpočine i sama obrada u okviru *while* petlje (linije 83 – 124). Petlja se izvršava sve dok u globalu postoje čvorovi koji nisu obrađeni, a u okviru nje se raspoređuju čvorovi do kojih se trenutno došlo ostalim nitima (linije 85 – 95),

upisuju obrađeni primljeni čvorovi u parcijalne putanje (linije 99 – 111) i na kraju vrši obrada uslova da li se petlja završava (linije 113 – 123).

```

3 void GenerationGraphHandler::generate_graph(map<string, Object*> parameters)
4 {
5     srand(parameters[SEED]);
6
7     long node_number = parameters[NODE_NUMBER],
8         max_degree = parameters[MAX_DEGREE],
9         min_degree = parameters[MIN_DEGREE],
10    available_nodes,
11    random_choice,
12    neighbour;
13
14    graph.resize(node_number);
15
16    vector<long> degrees(node_number);
17    vector<long> possibilities(node_number);
18
19    for(long i = 0; i < node_number; i++)
20        degrees[i] = 0;
21    possibilities[i] = 1;
22    graph[i] = {};
23
24
25    available_nodes = node_number;
26    for(long l = 0; l < node_number; l++)
27    {
28        long degree = rand() % (max_degree - min_degree) + min_degree;
29        for(long j = degrees[l]; j < degree; j++)
30        {
31            do
32            {
33                random_choice = rand() % available_nodes;
34                neighbour = possibilities[random_choice];
35            } while(neighbour == l);
36            graph[l].push_back(neighbour);
37            degrees[l]++;
38            graph[neighbour].push_back(l);
39            degrees[neighbour]++;
40
41            if(degrees[neighbour] == max_degree)
42                possibilities[random_choice] = possibilities[--available_nodes];
43        }
44    }
45
46
47

```

Slika 3: Generisanje grafa

```

44 void ParallelGraphHandler::multisource_bfs(map<string, Object*> parameters)
45 {
46     int thread_number = parameters[THREAD_NUMBER];
47
48     long node_number = parameters[NODE_NUMBER],
49         work_load = node_number / thread_number;
50
51    vector<long> sources = parameters[SOURCES].get_long_array();
52    vector<vector<vector<two_longs>>> next(thread_number);
53
54    paths.resize(node_number);
55
56    fill(paths.begin(), paths.end(), two_longs(-1, -1));
57
58    long level = 0,
59        frontiers_global_count = 0;
60
61    #pragma omp parallel num_threads(thread_number) firstprivate(level)
62    {
63        int tid = omp_get_thread_num();
64        long frontiers_count;
65        vector<long> frontiers;
66
67        for(long source : sources)
68        {
69            if(tid == source / work_load)
70            {
71                paths[source] = { 0, -1 };
72                frontiers.push_back(source);
73            }
74
75
76            #pragma omp master
77            frontiers_global_count = sources.size();
78            next[tid].resize(thread_number);
79
80            #pragma omp barrier
81
82            while(frontiers_global_count > 0)
83            {
84                for(long current_node : frontiers)
85                {
86                    for(long next_node : graph[current_node])
87                    {
88                        long thread = next_node / work_load;
89                        if(paths[next_node].first == -1)
90                            next[thread][tid].push_back({ current_node, next_node });
91                    }
92                }
93
94                #pragma omp barrier
95
96                frontiers.clear();
97
98                for(long i = 0; i < thread_number; i++)
99                {
100                    for(two_longs next_node : next[tid][i])
101                    {
102                        if(paths[next_node.second].first == -1)
103                            paths[next_node.second] = { level + 1, next_node.first };
104                        frontiers.push_back(next_node.second);
105                    }
106                }
107
108                next[tid][tid].clear();
109
110            }
111
112            frontiers_count = frontiers.size();
113
114            #pragma omp single
115            frontiers_global_count -= frontiers_count;
116
117            #pragma omp reduction (+:frontiers_global_count)
118            frontiers_global_count += frontiers_count;
119
120            level++;
121
122            #pragma omp barrier
123
124        }
125
126

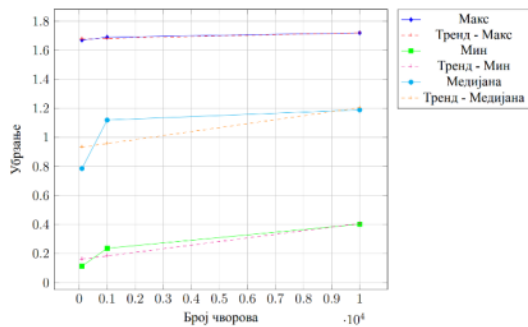
```

Slika 4: Paralelni BFS za arhitekture sa deljenom memorijom

Implementacija algoritma paralelne obrade grafa pomoću *BFS* za arhitekture sa distribuiranom memorijom prikazana na slici [Slika 5]. Ova implementacija je dosta složenija, u odnosu na implementaciju na arhitekturama sa deljenom memorijom, odnosno sekvencijalnu, ali se mogu podeliti slični regioni u kodu. Na linijama 127 – 155 se vrši inicijalizacija promenljivih koje se koriste pri obradi, ovo podrazumeva i sam tip podataka koji se razmenjuje (linije 153 – 155).

Nakon ovoga, na linijama 157 – 227 se nalazi sama *while* petlja, koja se izvršava ponovo sve dok postoji u globalu barem jedan čvor do kog se stiglo do trenutne dubine

malna dobijena medijana ubrzanja za distribuiranu implementaciju iznosi 119%. Dobra strana ovoga jeste to što performanse rastu i svakako će u jednom trenutku ubrzanje postati matematički beskonačno, jer neće postojati procesor koji bi sekvencijalno mogao da obradi toliko velike grafove, pa će distribuirano rešenje ujedno postati i jedino za ovakve poslove. No, pri izradi ovog rada, gde nije bilo uslova testiranja ovoga na većem broju računara, rezultati su obeshrabrujući, ali sa poznatim shvatanjem zašto je to tako.



Grafikon 2: Zavisnost ubrzanja distribuirane obrade od broja čvorova u grafu

6. ZAKLJUČAK

U ovom radu se detaljnije razmotrio problem pronalaženja najkraćeg puta, odnosno najmanjeg broja skokova, između dva čvora u grafu i to primenom algoritma za pretragu grafa po širini, poznatog kao *breadth-first traversal*. Implementirani su algoritmi paralelne verzije ovog algoritma, koristeći sisteme sa distribuiranom i deljenom memorijom. Takođe, dat je i prikaz rezultata dobijenih primenjujući pomenute imeplementacije na nasumično generisanim grafovima.

Na osnovu analize implementacija, može se zaključiti da je realizacija paralelnog pretraživanja grafa moguća, ali ni pod raznom jednostavna i nikako apsolutna. Ubrzanja do kojih se došlo su tu isključivo jer se radi o veštački sintezovanim grafovima, kao i činjenica da ta ubrzanja postoje samo za slučajeve visokog stepena povezanosti, odnosno paralelizacije u okviru okruženja koje radi sa deljivom memorijom.

Nažalost, za arhitekture sa distribuiranom memorijom se ne može reći ni to, s obzirom da je najbrže zabeleženo ubrzanje, koje je retko dostignuto, jako malo. Međutim, kao što je već napomenuto, ukoliko se radi o veoma velikim grafovima, koji jednostavno ne mogu stati u radnu memoriju jednog računara, distribuirana obrada je jedino moguće rešenje njihovog procesiranja, s toga je ipak ne treba zanemarivati.

7. LITERATURA

- [1] D. A. Bader, Massive Graph Analytics, CRC Press, 2022.
- [2] A. Yoo, E. Chow, K. Henderson, W. McLendon, U. Catalyurek i B. Hendrickson, A Scalable Distributed Parallel Breadth-First Search Algorithm on BlueGene/L, ACM/IEEE Conference on Supercomputing, 2005.
- [3] W. T. Tutte, Cambridge mathematical library: Graph theory, Cambridge: Cambridge University Press, 2001.
- [4] K. Mehlhorn, Data structures and algorithms 2, Berlin: Springer, 2011.
- [5] M. Neerajkumar, An efficient algorithm for shortest path tree in dynamic graph, LAP Lambert Academic Publishing, 2014.

Kratka biografija



Stefan Aleksić rođen je u Zaječaru 02. marta 2000. godine. Elektronski fakultet u Nišu, studijski program Elektrotehnika i računarstvo upisao je 2018. godine. Nakon završenih osnovnih studija, 2022. godine, upisao je master akademske studije na Fakultetu tehničkih nauka u Novom Sadu, na studijskom programu Računarstvo i automatika.

kontakt: stefan.sa.aleksic@gmail.com