

Δ -stepping: a parallelizable shortest path algorithm

Ulrich Meyer
& Peter Sanders
Presented by Lucas Bautista



Ulrich Meyer
Goethe University Frankfurt



Peter Sanders
Karlsruhe Institute of
Technology

Parallel Single Source Shortest Path (SSSP)

- Processing large graphs quickly is important
- SSSP is a common subroutine
- Betweenness Centrality, Route Finding, generally anything that needs to use the shortest distance between two nodes



SSSP

- Dealing with nonnegative directed graphs
- Label Setting Algorithms
- Label Correcting Algorithms
- Label setting has better worst-case bounds, but label-correcting is often better in practice

Label Setting Algorithms

- The algorithm finalizes the label for one node at a time
- Once a node's shortest distance is known, it never needs to change.
- The algorithm processes nodes in non-decreasing order of distance.

Example: Dijkstra's algorithm

Dijkstra's algorithm

Start with $d[s]=0$, For all $v \neq s$, $d[v] = \infty$

Maintain a priority queue of nodes keyed by tentative distance.

- Extract the node u with smallest $d[u]$
- Its label is now set (and final)
- Relax all outgoing edges (u,v) :
if $d[v] > d[u] + w(u,v)$ then update $d[v]$
- Repeat until queue is empty.

Complexity: $O(|E|\log(|V|))$

Label Correcting Algorithms

- Nodes' labels can be revised multiple times
- The algorithm keeps correcting them until no shorter paths exist
- There's no strict order in which nodes are processed
- We can relax edges arbitrarily and repeatedly.

Example: Bellman Ford Algorithm

Bellman Ford

Initialize $d[s]=0$, all others $d[v]=\infty$

Repeat $n-1$ times:

- For every edge (u,v) :
if $d[v] > d[u] + w(u,v)$ then $d[v] = d[u] + w(u,v)$

Complexity: $O(|E| * |V|)$

Each iteration “corrects” labels based on previously found paths.

No order for shortest path convergence

Δ -Stepping

- Buckets of vertices grouped by their temporary distance labels
- $B[i]$ contains vertices with labels in $[i \cdot \Delta, (i+1) \cdot \Delta]$
- Inner loop deals with relaxing light edges
- Outer edge deals with relaxing heavy edges

```

foreach  $v \in V$  do  $\text{tent}(v) := \infty$ 
 $\text{relax}(s, 0)$ ;
while  $\neg \text{isEmpty}(B)$  do
     $i := \min\{j \geq 0: B[j] \neq \emptyset\}$ 
     $R := \emptyset$ 
    while  $B[i] \neq \emptyset$  do
         $\text{Req} := \text{findRequests}(B[i], \text{light})$ 
         $R := R \cup B[i]$ 
         $B[i] := \emptyset$ 
         $\text{relaxRequests}(\text{Req})$ 
     $\text{Req} := \text{findRequests}(R, \text{heavy})$ 
     $\text{relaxRequests}(\text{Req})$ 

    (* Insert source node with distance 0 *)
    (* A phase: Some queued nodes left (a) *)
    (* Smallest nonempty bucket (b) *)
    (* No nodes deleted for bucket  $B[i]$  yet *)
    (* New phase (c) *)
    (* Create requests for light edges (d) *)
    (* Remember deleted nodes (e) *)
    (* Current bucket empty *)
    (* Do relaxations, nodes may (re)enter  $B[i]$  (f) *)
    (* Create requests for heavy edges (g) *)
    (* Relaxations will not refill  $B[i]$  (h) *)

Function  $\text{findRequests}(V', \text{kind} : \{\text{light}, \text{heavy}\})$  : set of Request
    return  $\{(w, \text{tent}(v) + c(v, w)) : v \in V' \wedge (v, w) \in E_{\text{kind}}\}$ 

Procedure  $\text{relaxRequests}(\text{Req})$ 
    foreach  $(w, x) \in \text{Req}$  do  $\text{relax}(w, x)$ 

Procedure  $\text{relax}(w, x)$ 
    if  $x < \text{tent}(w)$  then
         $B[\lfloor \text{tent}(w) / \Delta \rfloor] := B[\lfloor \text{tent}(w) / \Delta \rfloor] \setminus \{w\}$ 
         $B[\lfloor x / \Delta \rfloor] := B[\lfloor x / \Delta \rfloor] \cup \{w\}$ 
         $\text{tent}(w) := x$ 

    (* Insert or move  $w$  in  $B$  if  $x < \text{tent}(w)$  *)
    (* If in, remove from old bucket *)
    (* Insert into new bucket *)
  
```

Fig. 1. A sequential variant of Δ -stepping (with cyclical bucket reuse). The sets of light and heavy edges are denoted by E_{light} and E_{heavy} , respectively. Requests consist of a tuple (node, weight).

Bucket Processing

- Each vertex in the bucket has outgoing edges which are either “light” (weight $\leq \Delta$) or “heavy” (weight $> \Delta$)
- Relaxing an edge can cause the destination vertex to be inserted into the current bucket
- Process bucket until it is empty, then relax its heavy edges

```

foreach  $v \in V$  do  $\text{tent}(v) := \infty$ 
 $\text{relax}(s, 0);$ 
while  $\neg \text{isEmpty}(B)$  do
     $i := \min\{j \geq 0: B[j] \neq \emptyset\}$ 
     $R := \emptyset$ 
    while  $B[i] \neq \emptyset$  do
         $\text{Req} := \text{findRequests}(B[i], \text{light})$ 
         $R := R \cup B[i]$ 
         $B[i] := \emptyset$ 
         $\text{relaxRequests}(\text{Req})$ 
     $\text{Req} := \text{findRequests}(R, \text{heavy})$ 
     $\text{relaxRequests}(\text{Req})$ 

Function  $\text{findRequests}(V', \text{kind} : \{\text{light}, \text{heavy}\})$  : set of Request
    return  $\{(w, \text{tent}(v) + c(v, w)) : v \in V' \wedge (v, w) \in E_{\text{kind}}\}$ 

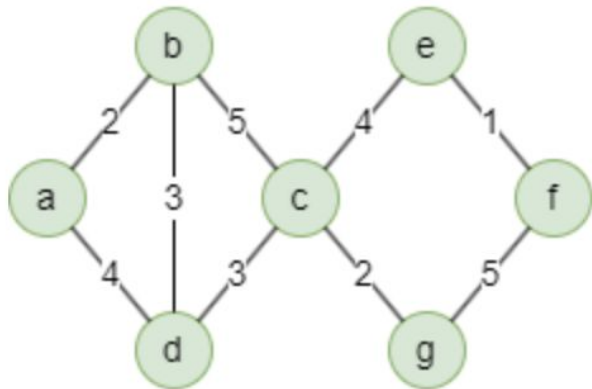
Procedure  $\text{relaxRequests}(\text{Req})$ 
    foreach  $(w, x) \in \text{Req}$  do  $\text{relax}(w, x)$ 

Procedure  $\text{relax}(w, x)$ 
    if  $x < \text{tent}(w)$  then
         $B[\lfloor \text{tent}(w) / \Delta \rfloor] := B[\lfloor \text{tent}(w) / \Delta \rfloor] \setminus \{w\}$ 
         $B[\lfloor x / \Delta \rfloor] := B[\lfloor x / \Delta \rfloor] \cup \{w\}$ 
         $\text{tent}(w) := x$ 

```

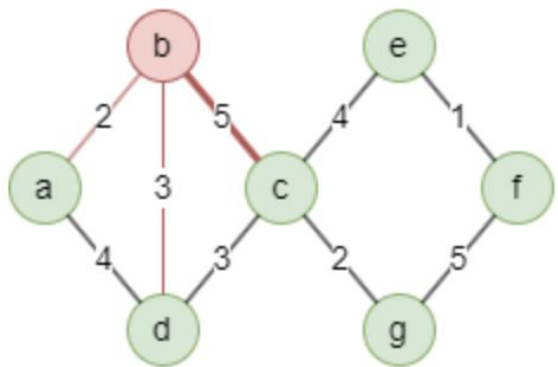
(* Insert source node with distance 0 *)
 (* A phase: Some queued nodes left (a) *)
 (* Smallest nonempty bucket (b) *)
 (* No nodes deleted for bucket $B[i]$ yet *)
 (* New phase (c) *)
 (* Create requests for light edges (d) *)
 (* Remember deleted nodes (e) *)
 (* Current bucket empty *)
 (* Do relaxations, nodes may (re)enter $B[i]$ (f) *)
 (* Create requests for heavy edges (g) *)
 (* Relaxations will not refill $B[i]$ (h) *)

Fig. 1. A sequential variant of Δ -stepping (with cyclical bucket reuseage). The sets of light and heavy edges are denoted by E_{light} and E_{heavy} , respectively. Requests consist of a tuple (node, weight).



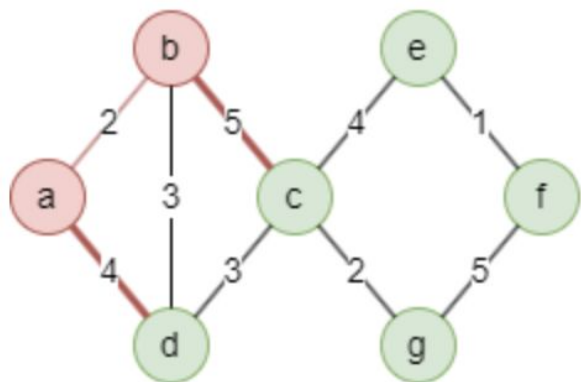
Node	a	b	c	d	e	f	g
dist()	∞	0	∞	∞	∞	∞	∞

Bucket	Distance Range	Node(s)
B[0]	[0, 3)	b



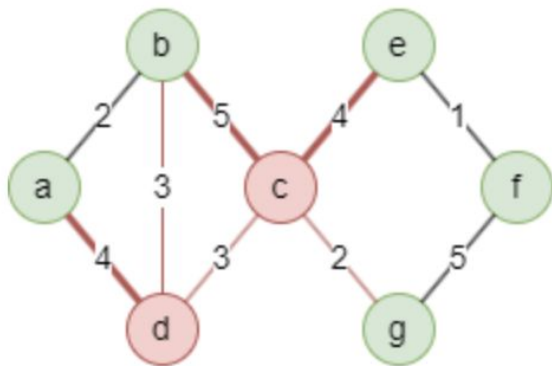
Node	a	b	c	d	e	f	g
dis()	2	0	∞	3	∞	∞	∞

Bucket	Distance Range	Node(s)
B[0]	[0, 3)	a
B[1]	[3, 6)	d



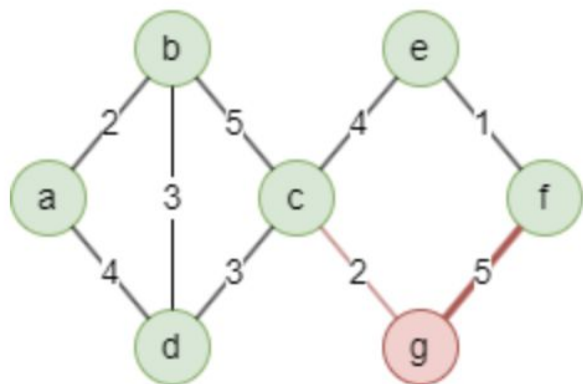
Node	a	b	c	d	e	f	g
dist()	2	0	5	3	∞	∞	∞

Bucket	Distance Range	Node(s)
B[0]	[0, 3)	
B[1]	[3, 6)	d c



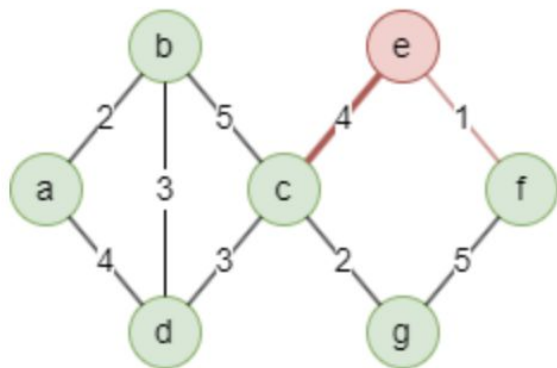
Node	a	b	c	d	e	f	g
dist()	2	0	5	3	9	∞	7

Bucket	Distance Range	Node(s)
B[0]	[0, 3)	
B[1]	[3, 6)	
B[2]	[6, 9)	g
B[3]	[9, 12)	e



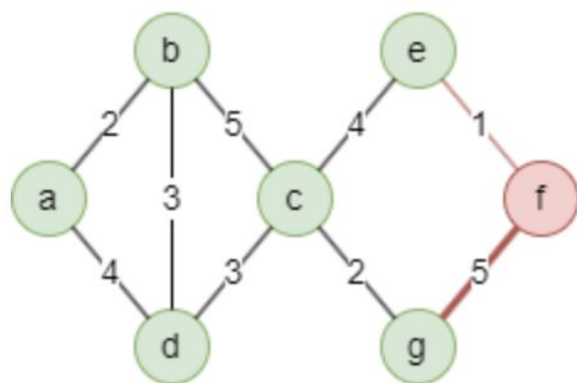
Node	a	b	c	d	e	f	g
dist()	2	0	5	3	9	12	7

Bucket	Distance Range	Contain Node(s)
B[0]	[0, 3)	
B[1]	[3, 6)	
B[2]	[6, 9)	
B[3]	[9, 12)	e
B[4]	[12, 15)	f



Node	a	b	c	d	e	f	g
dis()	2	0	5	3	9	10	7

Bucket	Distance Range	Contain Node(s)
B[0]	[0, 3)	
B[1]	[3, 6)	
B[2]	[6, 9)	
B[3]	[9, 12)	f
B[4]	[12, 15)	



Node	a	b	c	d	e	f	g
dist()	2	0	5	3	9	10	7

Bucket	Distance Range	Node(s)
B[0]	[0, 3)	
B[1]	[3, 6)	
B[2]	[6, 9)	
B[3]	[9, 12)	
B[4]	[12, 15)	

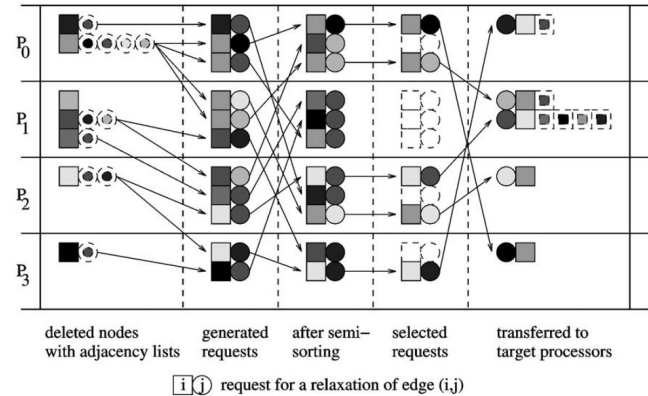
Δ -Stepping in Parallel

- Each PU is assigned nodes **randomly**. Adjacency lists are organized into heavy and light edges
- Each PU maintains a bucket structure with the queued nodes that it is responsible for
- Every $\Theta(\log n)$ iterations it is checked whether any PU has found a nonempty bucket and, if so, the globally smallest index with a nonempty bucket is found, is processed.
- Each PU scans its request buffer and sequentially performs the relaxations assigned to it. Since no other PUs work on its nodes the relaxations will be atomic.

Δ -Stepping in Parallel

- Instead of one PU handling all edges of a large node, multiple PUs cooperate to generate requests for its neighbors.
- If node u has 10 million outgoing edges \rightarrow split its adjacency list across 4 PUs; each generates part of the requests.
- The requests are grouped (batched) by which PU they'll be sent to (to reduce communication overhead) \rightarrow semisorted
- Before sending, we can filter out duplicates or non-improving requests locally to reduce communication overhead.

Fig. 6. Load balanced edge relaxation using semi-sorting.



Delta Stepping in a Distributed Setting

- The hash function $\text{ind}(w)$ replaces the index array used in the PRAM algorithm.
- Processors are assigned to groups of size 2^i . these groups share adjacency lists.
- The group cooperates via collective communication (broadcasts, reductions)

Choosing Δ

- Choosing $\Delta = 1$, makes our problem basically Dijkstra (with an array instead of a heap) (which has better runtime for sparsely bounded degrees)
- Choosing $\Delta \geq n$, makes our problem Bellman Ford
- We want to choose Δ so that we can take advantage of the parallelism but also want to reduce the amount of unnecessary work we do

Sequential Analysis

- Sequential Δ -stepping can be implemented in $O(n + m + L/\Delta + n\Delta + m\Delta)$
- On random edge weights $O(n + m + d * L)$
- $L := \max\{\text{dist}(v) : \text{dist}(v) < \infty\}$

Parallel Analysis

- Simple parallelization runs in $O(L/\Delta * d * l_{\Delta} * \log n)$
- Can take advantage of shortcut edges

$$l_{\Delta} = 1 + \max_{\langle u, v \rangle \in C_{\Delta}} \min \{ |A| : A = (u, \dots, v) \text{ is a minimum-weight } \Delta\text{-path} \}.$$

Performance

- The tested graphs ranged from 10^3 up to 10^6 nodes and comprised up to $3 * 10^6$ edges
- Tests were run on an INTEL Paragon with 16 processors
- Parallel/Distributed: for $n = 2^{19}$ nodes and $d = 3$, speedup **9.2** was obtained against the sequential Δ -stepping approach
- Sequential Δ -stepping approach is **3.1** times faster than an optimized implementation of Dijkstra's algorithm
- Results on dense graphs are slightly worse [Communication costs]

Strengths, Weaknesses, Thoughts

Strengths:

- I liked the connection between Dijkstra & Bellman Ford (and how delta stepping is something in the middle)
- Load balancing requests among processors

Weaknesses:

- Slightly hard to follow proofs

Future work:

- Better load balancing strategies
- Better use of communication (lower messages)
- Maybe extending it to negative edges as well using ideas from Bellman Ford

Discussion

- The experiments show $O(\log n)$ phases even though theory guarantees only $O(\log^2(n)/\log \log n)$
- Why might the empirical performance be better than the worst-case bound?
- Is there a better way to choose Δ , depending on the graph we are processing