

B3CC: Concurrency

07: Delta-stepping

Ivo Gabe de Wolff

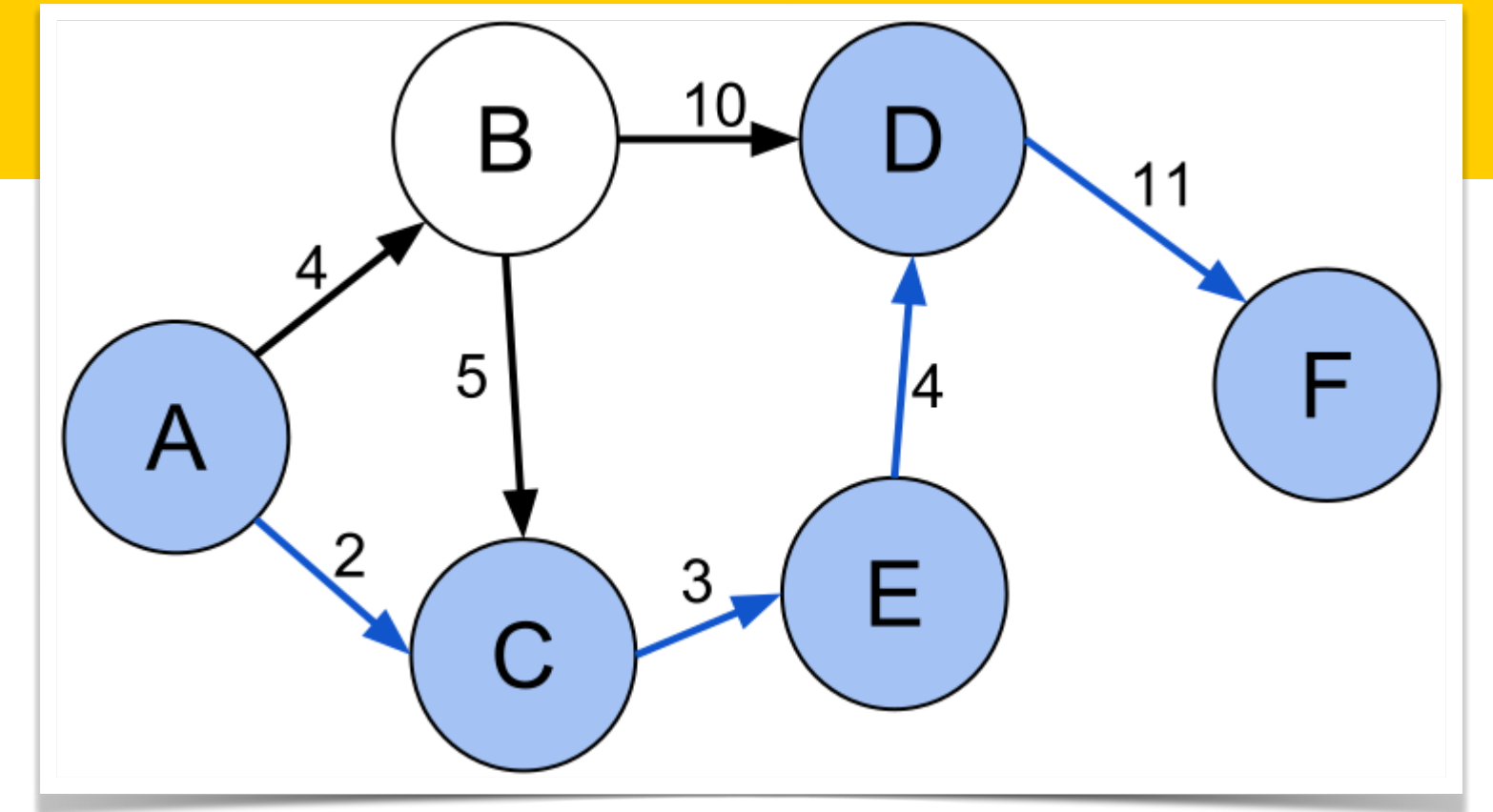
Announcement

- The second practical assignment has been released
 - <https://ics-websites.science.uu.nl/docs/vakken/b3cc/assessment.html>
 - You may work in pairs, if you wish
 - Deadline: 20-12-2023 @ 23:59

Parallel algorithms

- You have seen many sequential algorithms
 - In “datastructuren”
- Can we convert a sequential algorithm to a parallel algorithm?
 - No automatic approach
 - Sequential code typically has a long, sequential, critical path
- Today: Converting Dijkstra’s shortest-path algorithm to Delta-stepping

SSSP



- A central problem in algorithmic graph theory is the *single-source shortest path* problem
 - e.g. starting at vertex A, what is the shortest path to reach vertex F?
 - Many practical and theoretical applications
 - One of the benchmarks used in the Graph500 supercomputer ranking
 - The smallest problem size uses 2^{26} vertices, requiring 17 GB RAM

SSSP

- Given...
 - A directed graph $G(V, E)$ with $n = |V|$ nodes (or vertices) and $m = |E|$ edges
 - A distinguished node in the graph s : the “source”
 - A function c that returns the (non-negative) weight of a given edge in G
- Objective:
 - For each node v reachable from s ,
compute the weight of the minimum-weight (i.e. shortest) path from s to v
 - The *weight of the path* is the sum of the weights of its edges, denoted $\text{dist}(s, v)$ or $\text{dist}(v)$
 - If v is not reachable from s , then $\text{dist}(s, v) := \infty$

Dijkstra's algorithm

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take the unvisited node v with the shortest tentative distance
 - Its tentative distance is now fixed
 - Look at its neighbors: we may have a shorter path to them, via v
 - Update tentative distances of neighbors accordingly

Dijkstra's algorithm

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take the unvisited node v with the shortest tentative distance
 - Its tentative distance is now fixed
 - Look at its neighbors: we may have a shorter path to them, via v
 - Update tentative distances of neighbors accordingly

$O(n)$ iterations

$O(\log n)$ with
proper data structure
(min heap)

Towards parallel shortest path

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take some set of nodes
 - ~~Their tentative distances are now fixed~~
 - Look at their neighbors: we may have a shorter path to them
 - Update tentative distances of neighbors accordingly

Towards parallel shortest path

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take some set of nodes
 - ~~Their tentative distances are now fixed~~
 - Look at their neighbors: we may have a shorter path to them
 - Update tentative distances of neighbors accordingly

Work may be redundant:
Later iterations may need to look at
this node again.

Trade-off between parallelization
and work overhead

Towards parallel shortest path

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take some set of nodes
 - ~~Their tentative distances are now fixed~~
 - Look at their neighbors: we may have a shorter path to them
 - Update tentative distances of neighbors accordingly

What set?

Only the first unvisited node:
Dijkstra's algorithm

All nodes:
Bellman-Ford

Δ -stepping

- *Delta-stepping* is a parallelisable single-source shortest path algorithm
 - Algorithm stores an array of buckets B
 - Nodes are grouped by tentative distance in *buckets*
 - The range of distances in a bucket is parameter Δ (Greek letter Delta)
 - Bucket $B[i]$ stores the set of unsettled nodes v with
$$i \cdot \Delta \leq \text{tent}(v) < (i + 1) \cdot \Delta$$

Towards parallel shortest path

- Keep track of *tentative distance* per vertex
 - The distance of the shortest known path
- Repeatedly,
 - Take all nodes from the first non-empty bucket
 - ~~Their tentative distances are now fixed~~
 - *Find requests*: In parallel, look at their neighbors: we may have a shorter path to them
 - *Relax requests*: In parallel, update tentative distances of neighbors accordingly

This could add some of the same nodes back into the same bucket, or new nodes into this bucket.

Light and heavy edges

- A node may repeatedly be in the same bucket.
- This gives redundant work: we repeatedly look at its neighbors.
- To reduce this, we handle *light* and *heavy* edges separately.
 - Light edges ($c(e) \leq \Delta$) may cause that nodes are added back to the same bucket.
 - Heavy edges ($c(e) > \Delta$) can only affect later buckets.

Towards parallel shortest path

- Repeatedly,
 - Find index i of the first non-empty bucket.
 - Repeatedly handle all outgoing light edges from nodes $B[i]$:
 - Remove all nodes from $B[i]$
 - Find requests of light edges
 - Relax requests
 - Keep track of all nodes that have been in this bucket
 - When the bucket remains empty, handle all outgoing heavy edges of nodes that have been in $B[i]$:
 - Find requests of light edges
 - Relax requests

This could add some of the same nodes back into the same bucket, or new nodes into this bucket.

Initialisation

- All vertices have infinite tentative distance, except s which has distance zero
- All buckets are empty, except $B[0]$ which contains s

Basic algorithm

```

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 $\text{findRequests}(V', \text{kind} : \{\text{light}, \text{heavy}\}) : \text{set of Request}$
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 $\text{relax}(w, x)$ (* Insert or move w in B if $x < \text{tent}(w)$ *)
if $x < \text{tent}(w)$ **then**
 $B[\lfloor \text{tent}(w) / \Delta \rfloor] := B[\lfloor \text{tent}(w) / \Delta \rfloor] \setminus \{w\}$ (* If in, remove from old bucket *)
 $B[\lfloor x / \Delta \rfloor] := B[\lfloor x / \Delta \rfloor] \cup \{w\}$ (* Insert into new bucket *)
 $\text{tent}(w) := x$

Example

On the blackboard

Number of buckets

- How many buckets do we need?
 - Depends on the longest path (which we don't know yet)
 - If we use a cyclic or circular buffer, depends on the longest edge.
 - with a cyclic buffer we can reuse the buckets we no longer need.

Cyclic array of buckets

- Cyclic array (or circular buffer, ring buffer)
 - Only a small number of buckets are in use (not empty) at a time:
 - When we work with bucket i , we don't need index $i + length$ yet.
 - and we don't need $i - 1$ any more.
 - Reuse index i for value $i + length$, (and later $i + 2 * length, \dots$)
 - i th value is stored at index $i \% length$

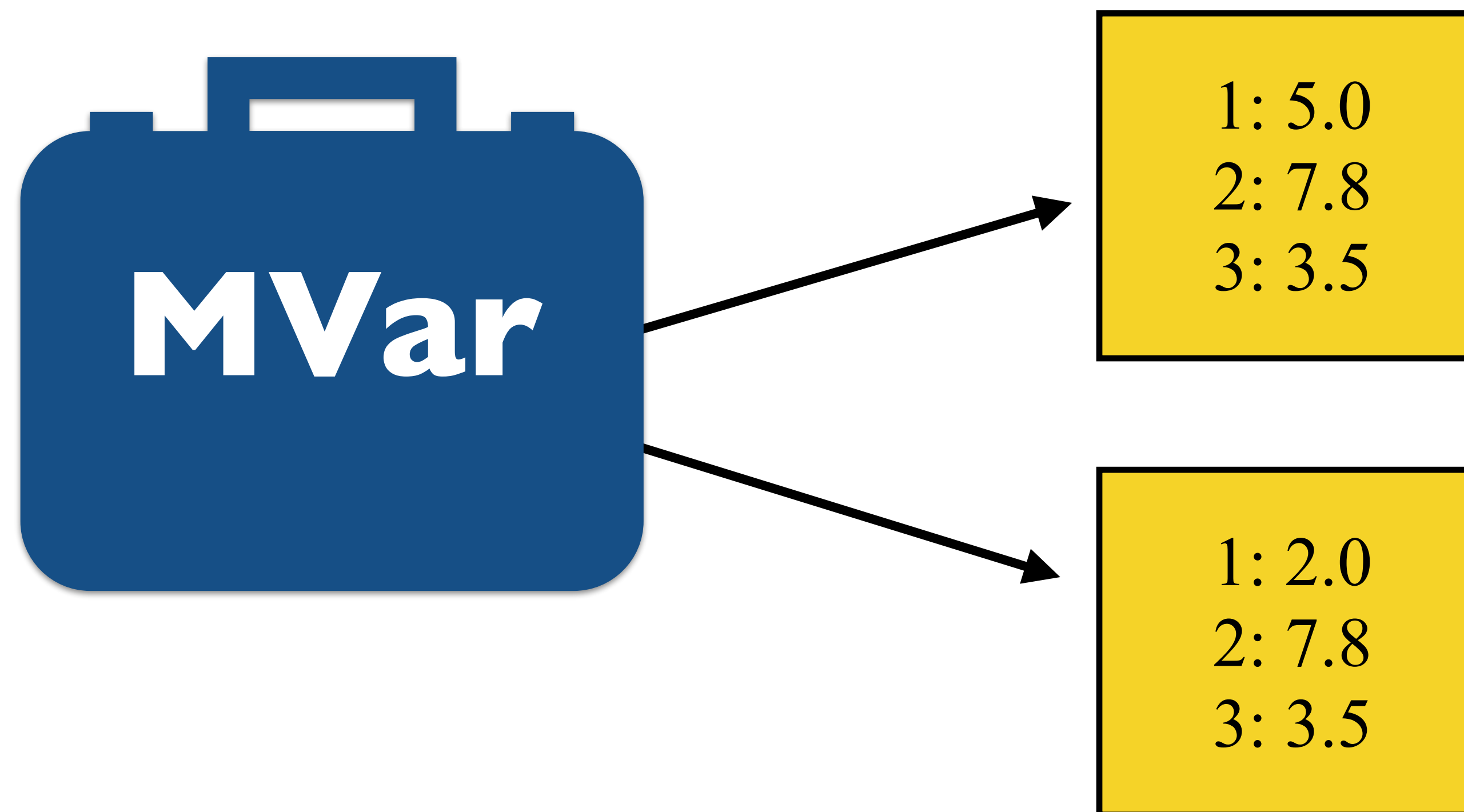
Utilities

- Inductive representation for graph structures
 - Data.Graph.Inductive.Graph contains functions for querying the given graph
 - Number of nodes / vertices in the graph
 - Return the in / out edges for the given node

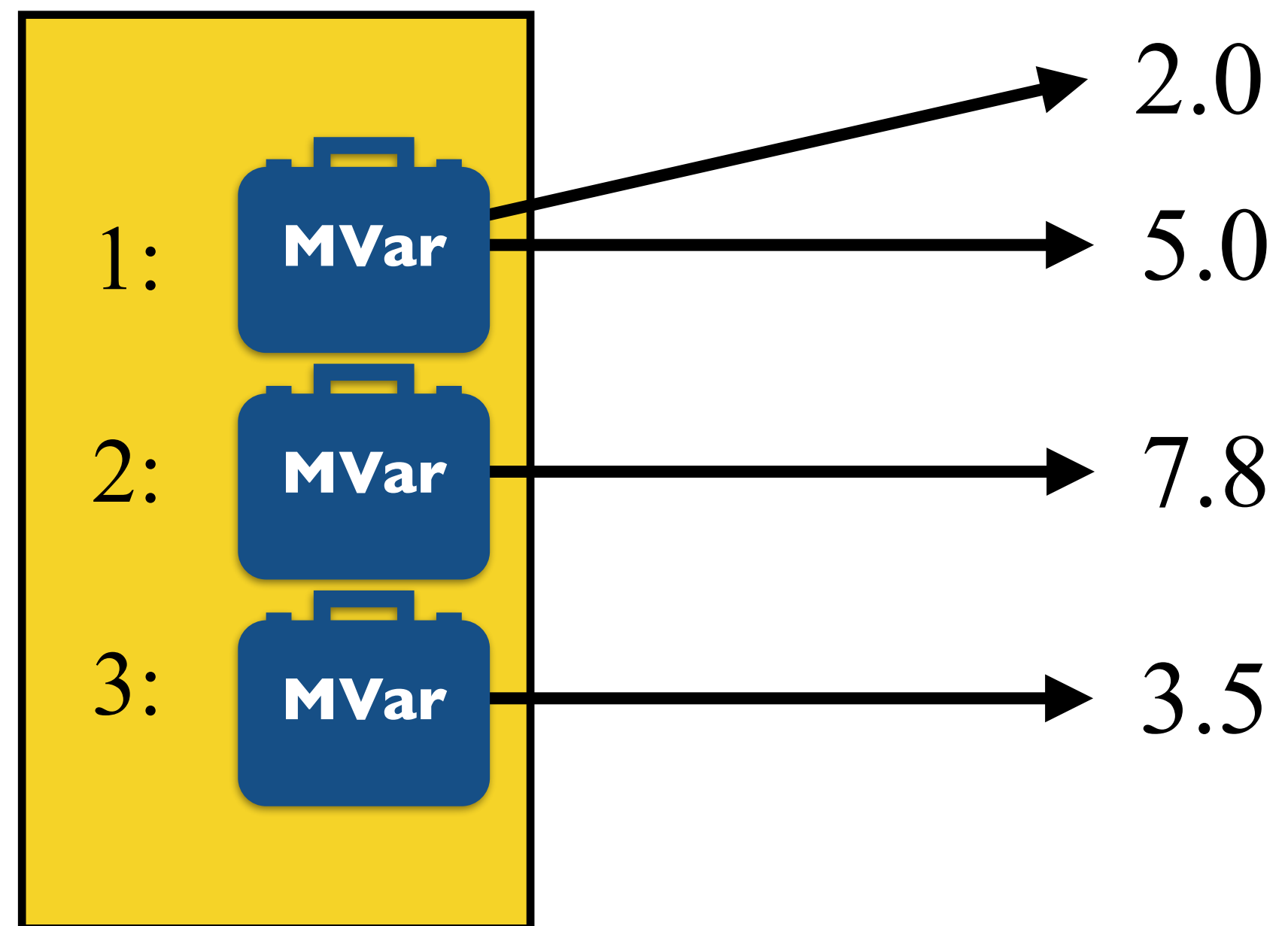
containers

- Assorted immutable data structures
 - `[Int]Map` (dictionary), `[Int]Set`, etc.
 - Put in an `IORef/MVar/etc.` to create a simple (non-concurrent) mutable container

MVar (IntMap Float)



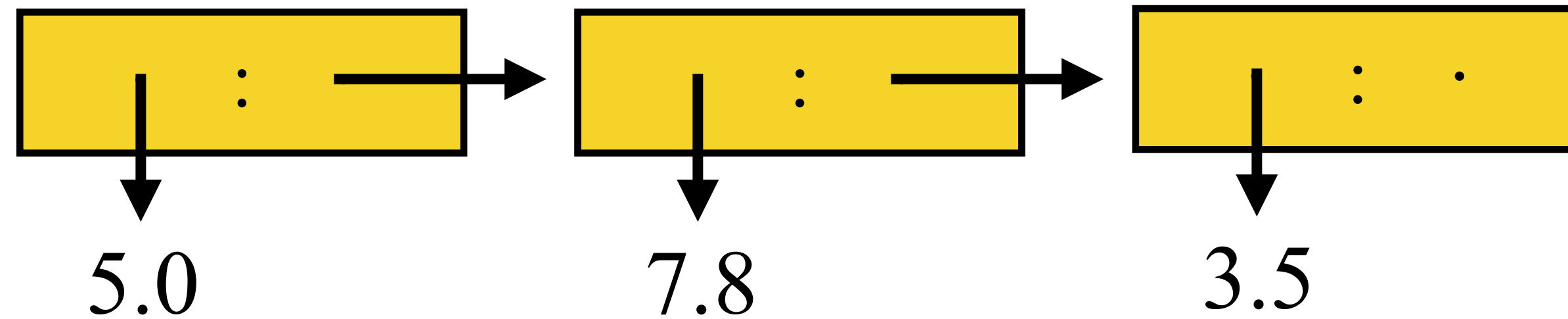
IntMap (MVar Float)



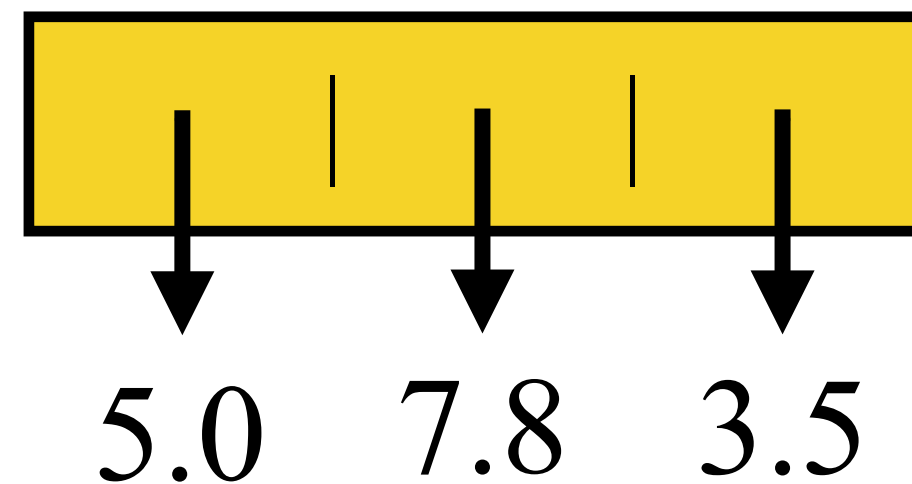
Lists vs Vectors

- Lists in Haskell are linked-lists:

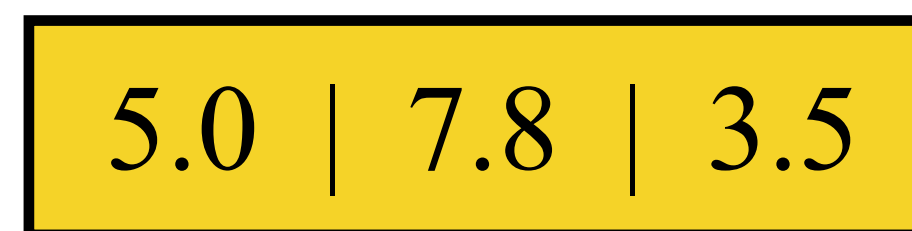
- $(:)$ and tail are $O(1)$,
- Indexing is $O(n)$



- Vectors are stored as arrays, with pointers to the values:



- Unboxed vectors store values instead of pointers in arrays:



vector

- (Un)boxed (im)mutable int-indexed arrays
 - Provides arrays in several flavours (i.e. underlying representation), but all with the same API
 - Data.Vector.Mutable
 - Boxed vectors (i.e. array of pointers) that can hold any structure
 - Data.Vector.Storable.Mutable
 - Unboxed vectors (i.e. array of values) that can hold only Storable (i.e. primitive) values
 - You can get a pointer directly to the array elements: useful for low-level atomic instructions
 - You can convert between different representations

- Provides the functionality of `atomicModifyIORef` on vectors
- For boxed vectors:

`atomicModifyIOVector`

`:: V.IOVector a → Int → (a → (a, b)) → IO b`

- For unboxed vectors:

`atomicModifyIOVectorFloat`

`:: M.IOVector Float → Int → (Float → (Float, b)) → IO b`

Conclusion

- The long, sequential, critical path was a problem.
- Trade-off between work overhead and parallelism:
 - We do some redundant work,
 - but when done properly, we will get a faster algorithm!
- Separation in light and heavy edges reduces work overhead.
- It is up to you to determine where the parallelism in the algorithm is (easy) and how to exploit this (hard)
- You are free to use IORefs, MVars, STM, mutable vectors, ...