# Declarative languages for distributed systems

- A research group at UC Berkeley lead by Prof. Hellerstein:
  - claims that the problems with distributed software come from the usage of imperative sequential programming languages to describe systems that are inherently non-sequential
  - resulting systems tend to be much smaller: 20KLOC / 1KLOC for HDFS
- Related PhD theses we've studied in this class:
  - Peter Alvaro: Data-centric Programming for Distributed Systems, 2015
  - Peter Bailis: Coordination Avoidance in Distributed Databases, 2015. I-Confluence

UNIVERSITY OF
WATERLOO

# Project goals

- Decided to verify claims on applicability of declarative logic programming for development of distributed software systems

- Decided to build one of the distributed data processing models presented in class

- Decided to implement Google's Pregel, as a simple synchronous model for parallel computation based on Valiant's Bulk Synchronous Parallel BSP model

- To test correctness of our Pregel model - implemented PageRank on top of it

Overview

UNIVERSITY OF
WATERLOO

# Bloom Bud declarative framework

- All data is represented as collections of facts (or tables containing records)

- New facts can be derived by declaring transformational rules

```
workers_list <= connect{ |worker|
  [worker.worker_addr, worker.id, false]
}
```

- No shared state: nodes exchange data as network messages (Overlog)

```
channel <- message ["IP:port recipient", "IP:port sender", payload_object]
```

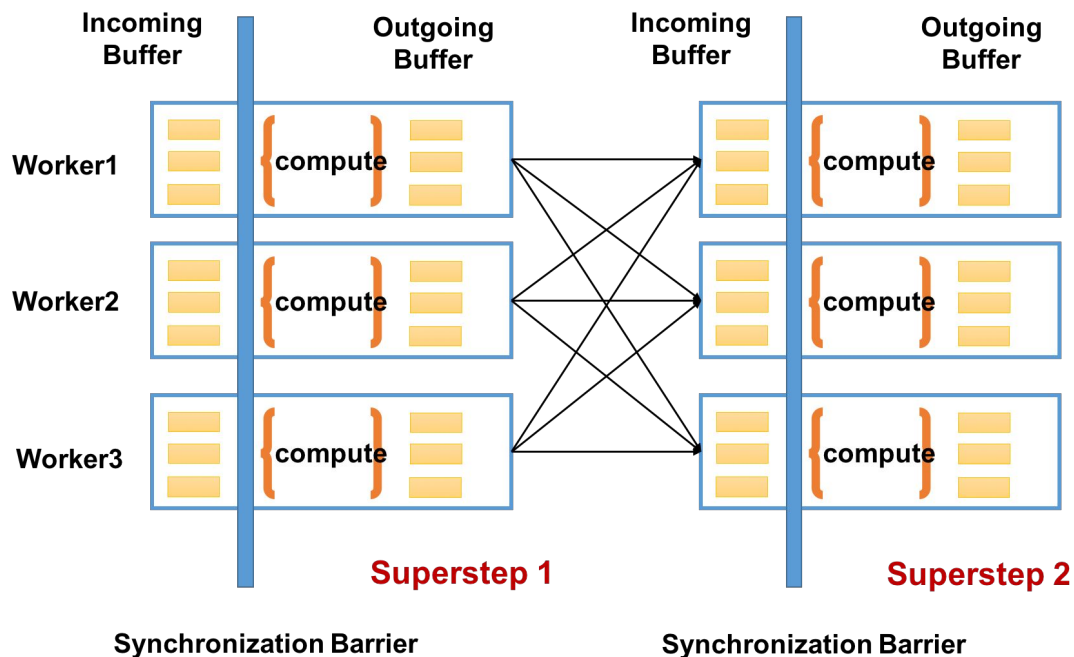- Introduction of notion of time - data collections evolve over time (Dedalus)

When is **counter** incremented?

```
counter(To,X+1) <= counter(To,X), request(To,_)
response(@From,X) <~ counter(@To,X), request(@To,From)
```

What does **response** contain?

Overview

UNIVERSITY OF
WATERLOO

# Building Pregel using Bud Bloom declarative framework

# Pregel distributed graph processing model

# Master node superstep coordination

```
1  #start iterating
2  supersteps <= master_stdio { |network_message|
3    [0, false, false] if(network_message.message=="start" and graph_loaded.reveal and supersteps.empty?)
4  }
5  #table :supersteps, [:id] => [:request_sent, :completed]
6
7
8  #if all workers completed the superstep, start a new one
9  supersteps <+ workers_list.group([], bool_and(:superstep_completed)) {|columns|
10   if columns.first == true and supersteps_count.reveal < MAX_SUPERSTEPS
11     [supersteps_count+1, false, false]
12   end
13 }
14
15 # for the latest superstep tuple in "supersteps", send a request to Workers to start the superstep
16 multicast <= supersteps.argmax([], :id) {|superstep|
17   if(!superstep.request_sent and !superstep.completed)
18     superstep.request_sent=true
19     ["start", {:superstep=>superstep.id}]
20   end
21 }
22 |
```

UNIVERSITY OF
WATERLOO

# Worker node superstep processing

```
1  state do
2    table :vertices, [:id] => [:value, :total_adjacent_vertices, :vertices_to, :messages_inbox]
3    table :queue_in_next, [:vertex_id, :vertex_from] => [:message_value]
4    table :queue_out, [:adjacent_vertex_worker_id, :vertex_from, :vertex_to] => [:message, :sent, :delivered]
5  end
6
7  # Generating vertex messages on superstep start
8  queue_out <+ (vertices * worker_input).pairs.flat_map do |vertex, worker_input_command|
9    if(worker_input_command.message.command=="start")
10     vertex_messages = @pregel_vertex_processor.compute(vertex)
11   end
12 end
13
14 # delivery of the vertex messages to adjacent vertices for the next Pregel superstep
15 vertex_pipe <~ (queue_out * workers_list).pairs(:adjacent_vertex_worker_id => :id) do |vertex_message, worker|
16     if(vertex_message.sent == false)
17       vertex_message.sent = true
18       [worker.worker_addr, ip_port(), vertex_message]
19     end
20 end
21
22 # remove all outgoing vertex messages from "queue_out" in next timestep
23 # They are sent to recipients in the current timestep.
24 queue_out <- (queue_out * queue_out.group([], bool_and(:sent))).lefts
25
26 # send back a confirmation to Master that the superstep is complete
27 # This  message is sent after all vertex_messages were *sent*, not *delivered*
28 control_pipe <~ queue_out.group([], bool_and(:sent)) {|vertex_messages_sent|
29   [@master_address, ip_port, "success"] if vertex_messages_sent==true
30 }
```

# Pregel implementation

UNIVERSITY OF
WATERLOO

# PageRank implementation

```ruby
1  class PageRankVertexProcessor
2    def compute(vertex)
3      messages = []
4      if(!vertex.messages_inbox.nil? and !vertex.messages_inbox.empty?)
5        new_vertex_value=0
6        vertex.messages_inbox.each {|message|
7          new_vertex_value+=message[1]
8        }
9        vertex.value = 0.15/@graph_loader.vertices_all.size + 0.85*new_vertex_value
10     end
11
12     vertex.vertices_to.each { |adjacent_vertex|
13       adjacent_vertex_worker_id = @graph_loader.graph_partition_for_vertex(adjacent_vertex)
14       messages << [adjacent_vertex_worker_id, vertex.id, adjacent_vertex,
15         vertex.value.to_f / vertex.total_adjacent_vertices]
16     }
17     messages
18   end
19 end
```

# Pregel implementation

UNIVERSITY OF
WATERLOO

# Comparing declarative and imperative programming

# Advantages - less code

```
1   #send commands to all workers
2   control_pipe <~ (workers_list * multicast).combos do |worker, message|
3     [worker.worker_addr, ip_port, message]
4   end
5
6   # update workers list on job-completion messages
7   workers_list <+- (workers_list * control_pipe)
8     .pairs(workers_list.worker_addr => control_pipe.from) do |worker, command|
9       if(command.message.command == "load" and command.message.params[:status]=="success")
10        [worker.worker_addr, worker.id, true, worker.superstep_completed]
11      elsif(command.message.command == "start" and command.message.params[:status]=="success")
12        #worker completed the current superstep
13        [worker.worker_addr, worker.id, worker.graph_loaded, true]
14      end
15  end
```

UNIVERSITY OF WATERLOO

# Troubles, limitations

```
bloom :superstep_initialization do
  # table :queue_in_next, [:vertex_id, :vertex_from] => [:message_value]
  vertices <+- (vertices * control_pipe).pairs do |vertex, payload|
    if payload.message.command=="start"
      messages = []
      queue_in_next.each {|message|
        if vertex.id == message.vertex_id
          messages << [message[1], message[2]]
        end
      }
      [vertex.id, vertex.value, vertex.total_adjacent_vertices, vertex.vertices_to, messages]
    end
  end
end
```

Iterating over all network messages, imperative code

UNIVERSITY OF
WATERLOO

# Demo

|   | y | a | m |
|---|---|---|---|
| y | 0.5 | 0.5 | 0 |
| a | 0.5 | 0 | 1 |
| m | 0 | 0.5 | 0 |

| r1 | r2 | r3 | r4 | r5 | r6 | r7 | r8 | r9 | r10 | r11 | r12 | r13 | r14 | r15 |
|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| 0.333 | 0.333 | 0.417 | 0.375 | 0.417 | 0.385 | 0.411 | 0.391 | 0.408 | 0.394 | 0.405 | 0.396 | 0.403 | 0.397 | 0.402 |
| 0.333 | 0.500 | 0.333 | 0.458 | 0.354 | 0.438 | 0.370 | 0.424 | 0.380 | 0.416 | 0.387 | 0.410 | 0.392 | 0.407 | 0.394 |
| 0.333 | 0.167 | 0.250 | 0.167 | 0.229 | 0.177 | 0.219 | 0.185 | 0.212 | 0.190 | 0.208 | 0.194 | 0.205 | 0.196 | 0.203 |

PageRank calculation by matrix multiplication:   PageRank weights assignment using Pregel vertex-centric model:

Y_next_iteration_value: r2[0]=[0,0]*r1[0]+[0,1]*r1[1]+[0,2]*r1[2]
A_next_iteration_value: r2[1]=[1,0]*r1[0]+[1,1]*r1[1]+[1,2]*r1[2]
M_next_iteration_value: r2[2]=[2,0]*r1[0]+[2,1]*r1[1]+[2,2]*r1[2]

UNIVERSITY OF WATERLOO

# TCP network communication (instead of UDP)