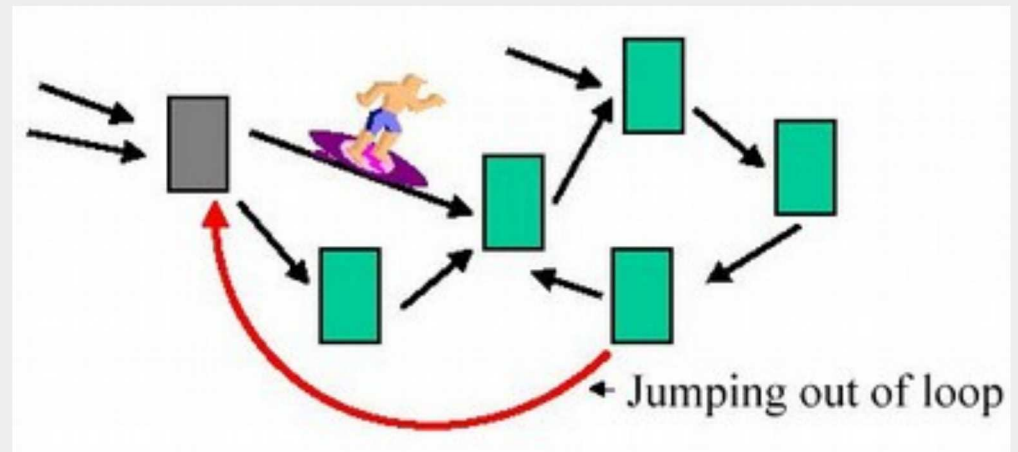# Example 2: PageRank

- Input: pages with output links
- Output: ranking of page
- Based on the random surfer model
  - The page rank is the probability to reach a page by:
    - Randomly clicking on links
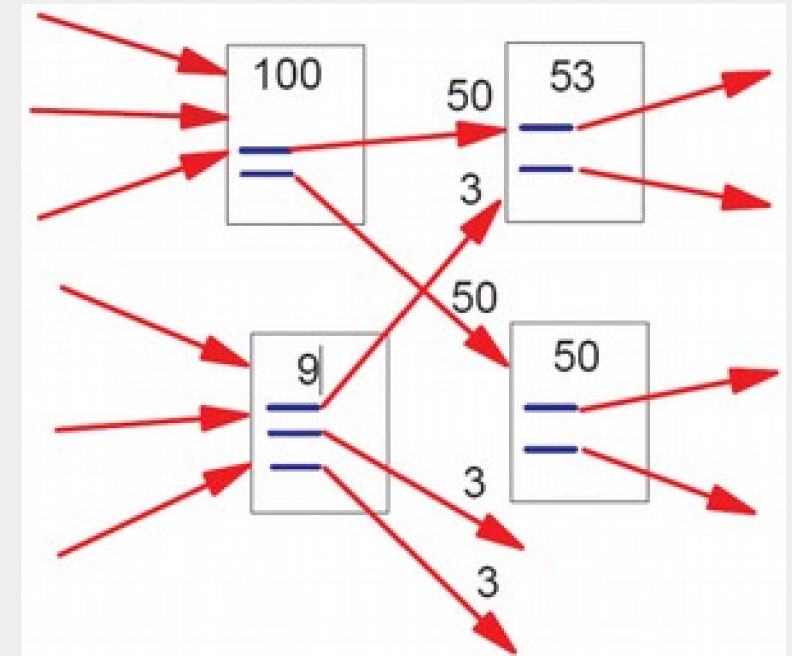    - Sometimes the user input an URL directly: put a small probability on this


← Jumping out of loop

# PageRank Formula

- $Rk(P) = (1-d)/n + d [ Rk(T_1) / out(T_1) + \cdots + Rk(T_n) / out(T_n) ]$

  with :

  - n: number of pages

  - P: ranked page

  - d: damping factor (usually 0.85)

  - $T_1,...,T_n$: pages pointing to P

  - $out(T_i)$: number of output links in $T_i$

- Iterate: $Rk^{i+1}(P)$ computed from $Rk^i(T_n)$

  - $Rk^0(P) = 1$

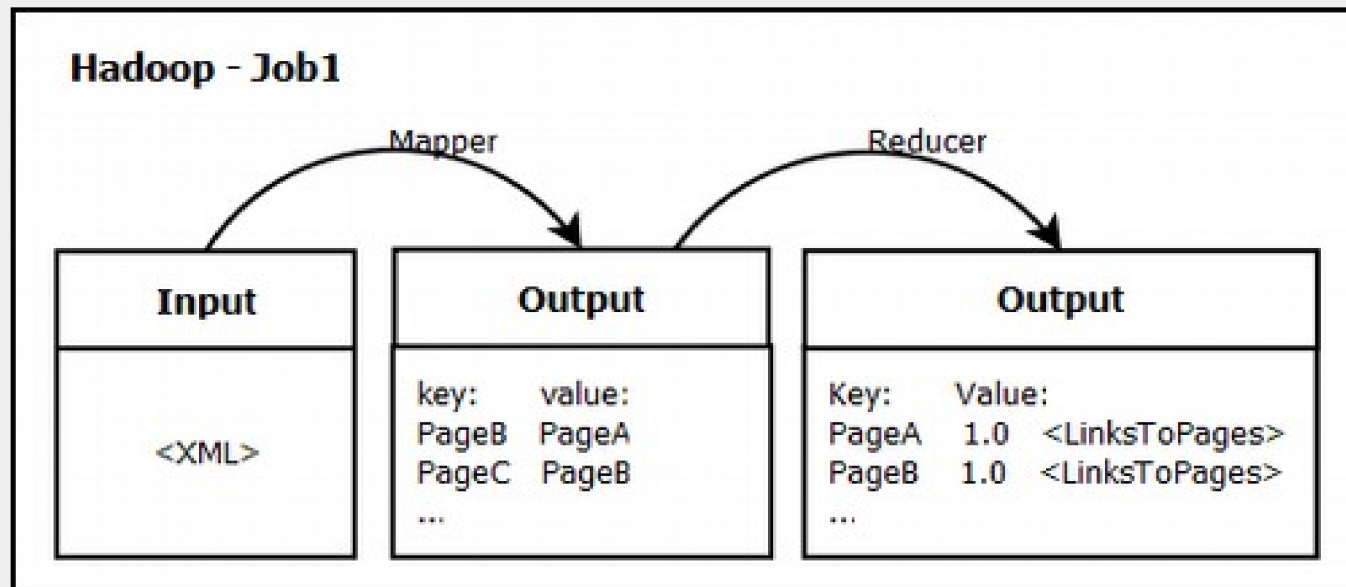- Stops when all the pageranks are stable during two iterations

# PageRank stage 1

- Extract links from XML pages

map: (pageid, body) → (pageid, (Rk$^0$, list(pageid)))

reduce: identity

# PageRank stage 2

- Iterates on page ranks computation

$$\text{map}^i: (pid, (Rk^i, list(pid))) \rightarrow (pid, (Rk^i, list(pid))) \text{ and}$$
$$\text{listof } (pid, (out^i, list()))$$

$$\text{reduce}^i: (pid, (Rk^i, list(pid))) \text{ and listof } (pid, (out^i, list()))$$
$$\rightarrow (pid, (Rk^{i+1}, list(pid)))$$