# PageRank: The Complete Algorithm

- **Input:** **Graph $G$ and parameter $\beta$**
  - Directed graph $G$ with **spider traps** and **dead ends**
  - Parameter $\beta$

- **Output: PageRank vector $r$**
  - **Set:** $r_j^{(0)} = \frac{1}{N}, \quad t = 1$
  - **do:**

    - $\forall j:\ {r'}_j^{(t)} = \sum_{i \to j} \beta \, \frac{r_i^{(t-1)}}{d_i}$

      ${r'}_j^{(t)} = 0$ if in-deg. of $j$ is **0**

    - **Now re-insert the leaked PageRank:**

    $\forall j:\ r_j^{(t)} = {r'}_j^{(t)} + \frac{1-S}{N}$ where: $S = \sum_j {r'}_j^{(t)}$

    - $t = t + 1$
  - **while** $\sum_j \left| r_j^{(t)} - r_j^{(t-1)} \right| > \varepsilon$

# Sparse Matrix Encoding

- **Encode sparse matrix using only nonzero entries**
  - Space proportional roughly to number of links
  - Say 10N, or 4*10*1 billion = 40GB
  - **Still won't fit in memory, but will fit on disk**

| source node | degree | destination nodes |
|---|---|---|
| 0 | 3 | 1, 5, 7 |
| 1 | 5 | 17, 64, 113, 117, 245 |
| 2 | 2 | 13, 23 |

# Basic Algorithm: Update Step

- **Assume enough RAM to fit $r^{new}$ into memory**
  - Store $r^{old}$ and matrix **M** on disk
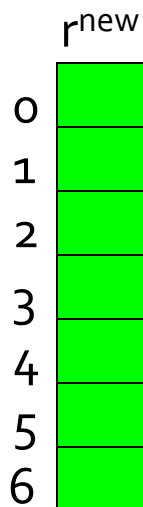- **Then 1 step of power-iteration is:**

  Initialize all entries of $r^{new}$ to $(1-\beta)/N$

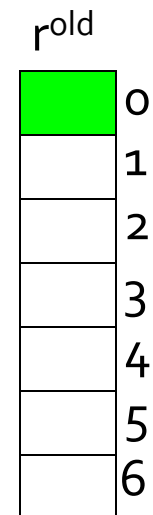  For each page $p$ (of out-degree $n$):

        Read into memory: $p, n, dest_1,…,dest_n, r^{old}(p)$

        for j = 1…n: $r^{new}(dest_j)$ += $\beta\ r^{old}(p) / n$

$r^{new}$

| 0 |
|---|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |

| src | degree | destination |
|-----|--------|-------------|
| 0 | 3 | 1, 5, 6 |
| 1 | 4 | 17, 64, 113, 117 |
| 2 | 2 | 13, 23 |

$r^{old}$

| 0 |
|---|
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |

# Analysis

- **Assume enough RAM to fit $r^{new}$ into memory**

  - Store $r^{old}$ and matrix **M** on disk

- **In each iteration, we have to:**

  - Read $r^{old}$ and **M**

  - Write $r^{new}$ back to disk

  - Input/Output cost = $2|r| + |M|$

- **Question:**

  - What if we could not even fit $r^{new}$ in memory?

# Block-based Update Algorithm

r^new

| | |
|---|---|
| 0 | 🟩 |
| 1 | 🟩 |

| 2 | |
| 3 | |

| 4 | |
| 5 | |

| src | degree | destination |
|-----|--------|-------------|
| 0 | 4 | 0, 1, 3, 5 |
| 1 | 2 | 0, 5 |
| 2 | 2 | 3, 4 |

r^old

| | |
|---|---|
| 🟩 | 0 |
| | 1 |
| | 2 |
| | 3 |
| | 4 |
| | 5 |

# Analysis of Block Update

- **Similar to nested-loop join in databases**
  - Break $r^{new}$ into $k$ blocks that fit in memory
  - Scan $M$ and $r^{old}$ once for each block
- **$k$ scans of $M$ and $r^{old}$**
  - $k(|M| + |r|) + |r| = k|M| + (k+1)|r|$

- **Can we do better?**
  - Hint: $M$ is much bigger than $r$ (approx 10-20x), so we must avoid reading it $k$ times per iteration

# Block-Stripe Update Algorithm

$r^{new}$

| | |
|---|---|
| 0 | |
| 1 | |

|  src | degree | destination |
|---|---|---|
| 0 | 4 | 0, 1 |
| 1 | 3 | 0 |
| 2 | 2 | 1 |

$r^{old}$

| | |
|---|---|
| | 0 |
| | 1 |
| | 2 |
| | 3 |
| | 4 |
| | 5 |

| | |
|---|---|
| 2 | |
| 3 | |

| | | |
|---|---|---|
| 0 | 4 | 3 |
| 2 | 2 | 3 |

| | |
|---|---|
| 4 | |
| 5 | |

| | | |
|---|---|---|
| 0 | 4 | 5 |
| 1 | 3 | 5 |
| 2 | 2 | 4 |

# Block-Stripe Analysis

- **Break _M_ into stripes**

  - Each stripe contains only destination nodes in the corresponding block of $r^{new}$

- Some additional overhead per stripe

  - But it is usually worth it

- Input/Output cost per iteration

  - $|M|(1+\varepsilon) + (k+1)|r|$

# Some Problems with Page Rank

- **Measures generic popularity of a page**
  - Biased against topic-specific authorities
  - **Solution:** Topic-Specific PageRank
- **Uses a single measure of importance**
  - Other models e.g., hubs-and-authorities
  - **Solution:** Hubs-and-Authorities
- **Susceptible to Link spam**
  - Artificial link topographies created in order to boost page rank
  - **Solution:** TrustRank