# Link Lifetimes and Randomized Neighbor Selection in DHTs

## Zhongmei Yao

Joint work with Dmitri Loguinov

Internet Research Lab
Department of Computer Science
Texas A&M University, College Station, TX 77843

April 15, 2008

1

# Agenda

- Background and Motivation
  - Terminology and related work

- Link Lifetime Model for switching systems
  - General DHT space, neighbor dynamics, A semi-Markov chain

- Lifetimes of Deterministic Links

- Lifetimes of Randomized Links

- Wrap-up

# Terminology

- User churn
  - User arrivals and departures are not synchronized

- Link creation in routing tables
  - Each user generates $k$ out-links pointing to its neighbors

- Non-switching systems (e.g., Kad and Gnutella)
  - The link points to the same neighbor until it fails

- Switching systems (classic DHTs)
  - Links switch to new neighbors before the current neighbor dies

- Link lifetimes
  - Time duration when the neighbor adjacent to the link is alive

# Terminology

- Repair of failed links
  - Detect failed links and replace with alive peers within $S$ time units

- Link churn
  - The dynamic behavior of links

# Background

**Analysis of link churn**

Unstructured P2P Networks

DHTs

Switching systems

Non-switching applied to randomized DHTs

Pandurangan 2001,
Leonard 2005a,
Leonard 2005b,
Yao 2006,
Yao 2007

Exponential lifetimes

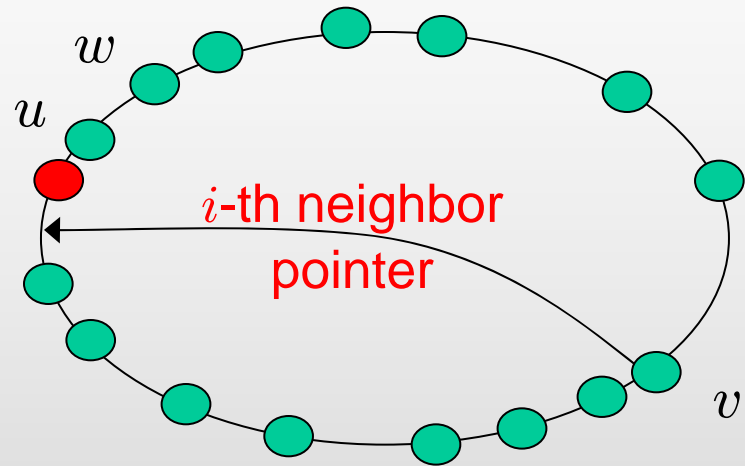Heavy-tailed lifetimes

Godfrey 2006,
Tan 2007

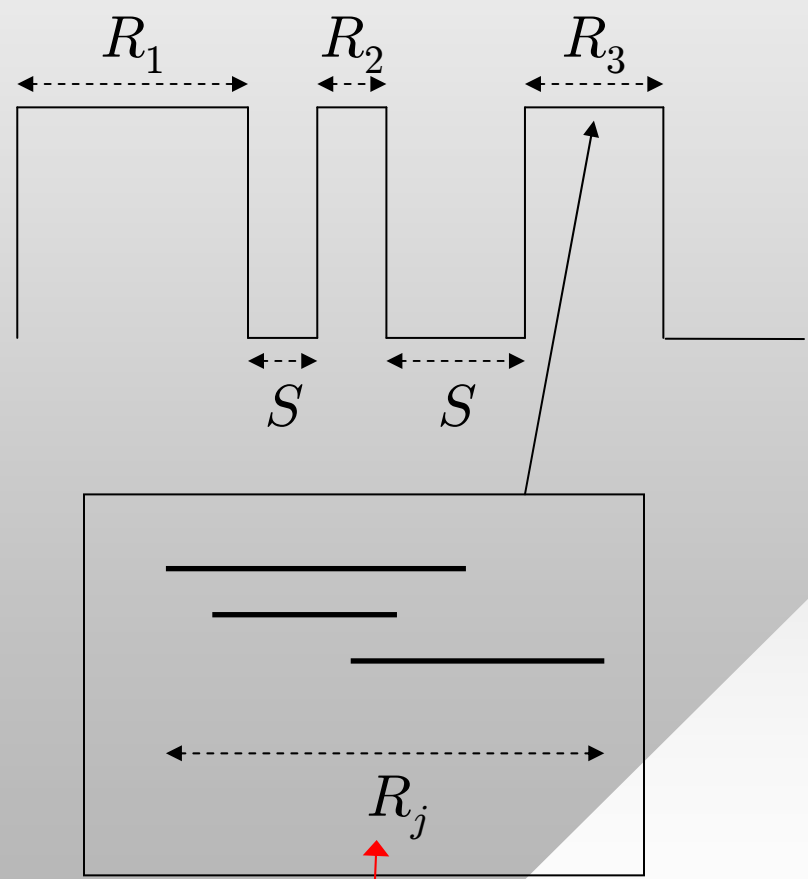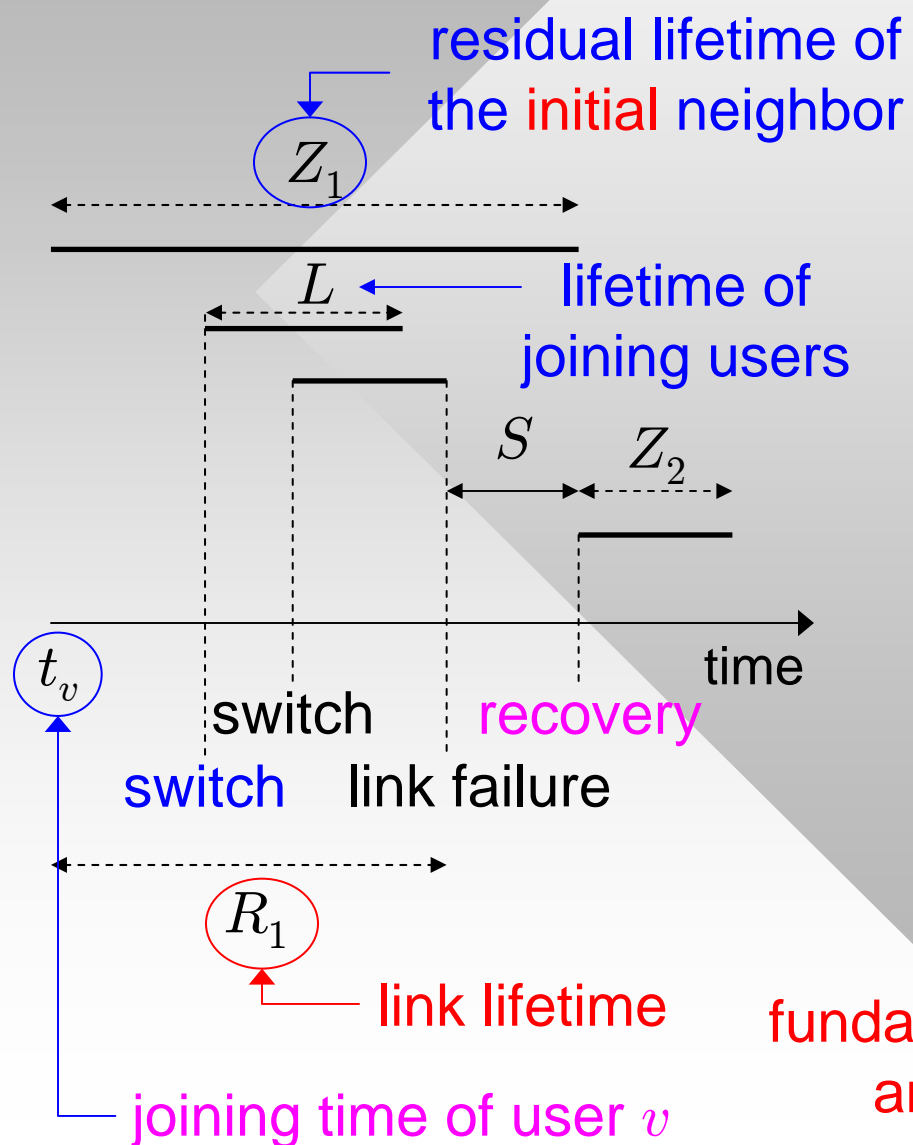Liben-Nowell 2002,
Krishnamurthy 2005

No prior work

# Agenda

- Background and Motivation
  - Terminology and related work

- Link Lifetime Model for switching systems

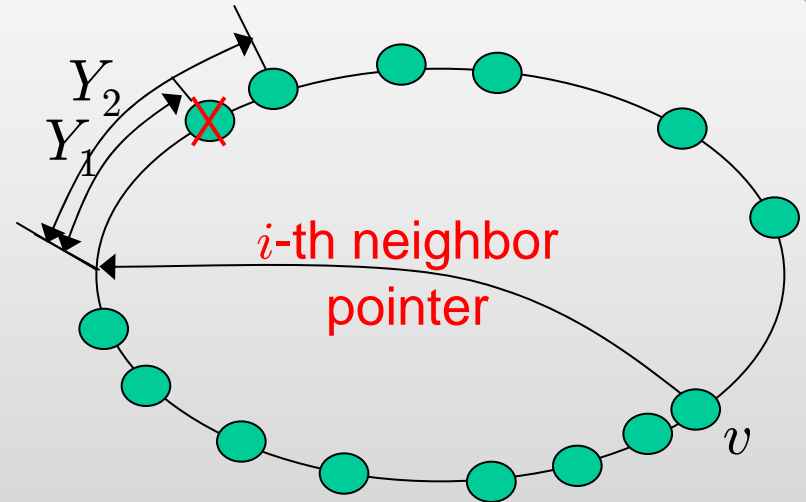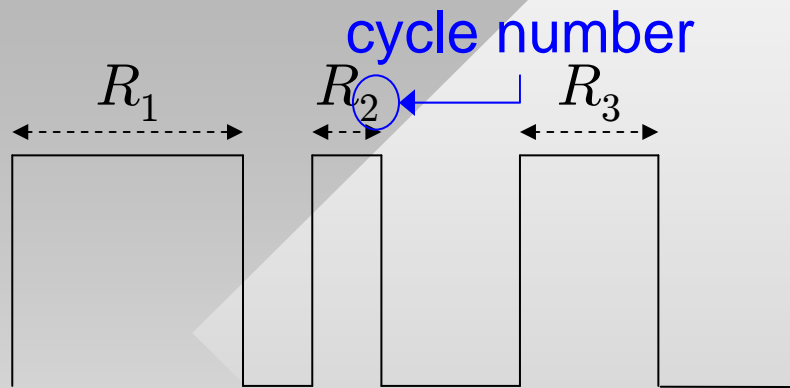- Deterministic Links

- Randomized Links

- Wrap-up

# Links in DHTs



$i$-th neighbor pointer

- The DHT space: consider a unit ring where hash indexes of users are uniform in $[0, 1)$
  - Random-split: zones are split at the hash indexes of joining users

- Fixed rules for selecting neighbors in routing tables
  - The successor of each neighbor pointer is $v$'s neighbor

- Link ownership changes under churn
  - Recovery: an existing neighbor dies and the ownership is assigned to the successor of the failed neighbor
  - Switch: a link switches to new users who arrive into the zone before the current successor fails

# Link ON/OFF Behavior

residual lifetime of the initial neighbor

$Z_1$

$L$

lifetime of joining users

$S$    $Z_2$

$t_v$

time

switch

recovery

switch    link failure

$R_1$

link lifetime

joining time of user $v$

$R_1$    $R_2$    $R_3$

$S$    $S$

$R_j$

fundamental to the studies of resilience and performance of the system

8

# Zone Size

cycle number

$R_1$    $R_2$    $R_3$

$Y_2$
$Y_1$

$i$-th neighbor pointer

$v$

- Denote by $Y_j$ is the zone size from the neighbor pointer to the initial neighbor who starts the $j$th cycle
  - Variable $Y_1$ is the zone size of the initial neighbor obtained when user $v$ joins the system

- It determines the arrival rate of new users that split the zone and become the owner of the neighbor pointer
  - Large $Y_j$ implies that more users arrive into the zone

- Other $Y_j$ correspond to link recovery: the initial neighbor is found in replacement of the failed neighbor for $j = 2, 3, \ldots$

# Link Lifetime Model Overview
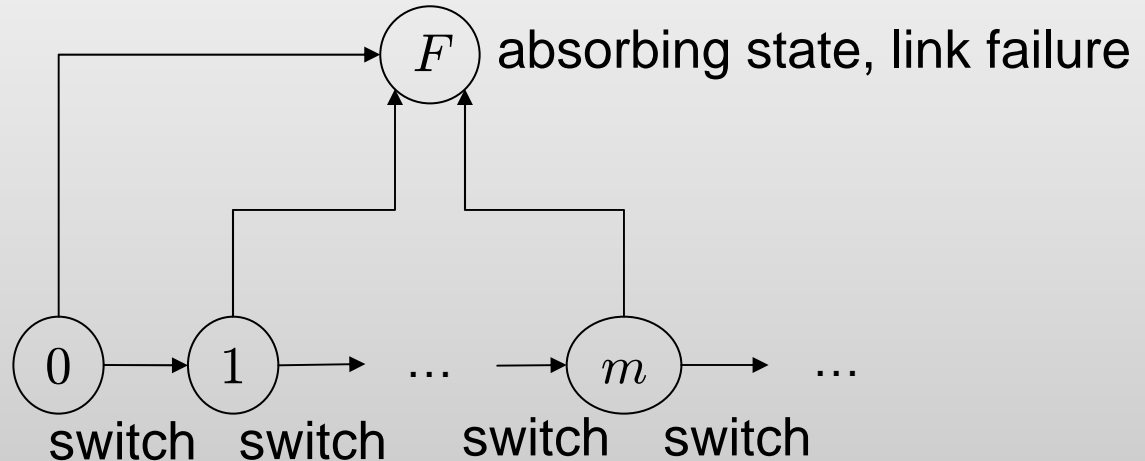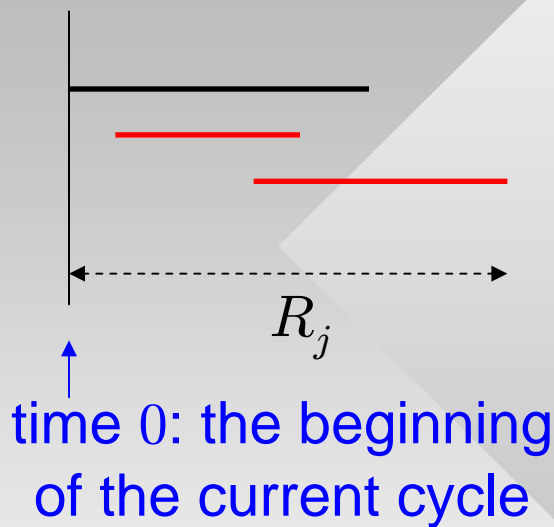
- Define the conditional link lifetime $R(y)$ as the link ON duration conditioned on the zone size $Y_j = y$
  - We use a semi-Markov chain to study $R(y)$

- Properties of link lifetimes can then be examined:

$$P(R_j < x) = \int_0^1 P(R(y) < x) f_{Y_j}(y) dy$$

the PDF of $Y_j$

  - Compute the distribution of $Y_j$ for deterministic DHTs and randomized DHTs accordingly

10

# Conditional Link Lifetimes



$R_j$

time 0: the beginning of the current cycle

$F$ absorbing state, link failure

0 → 1 → … → $m$ → …

switch   switch   switch   switch

- Denote by $A_t^y$ the number of switches (to new users) that have occurred in $[0, t]$ for given zone size $Y_j = y$

- Using notation $A_t^y$, we describe:

$$R(y) = \inf\{t > 0 : A_t^y = F | A_0^y = 0, Y_j = y\}$$
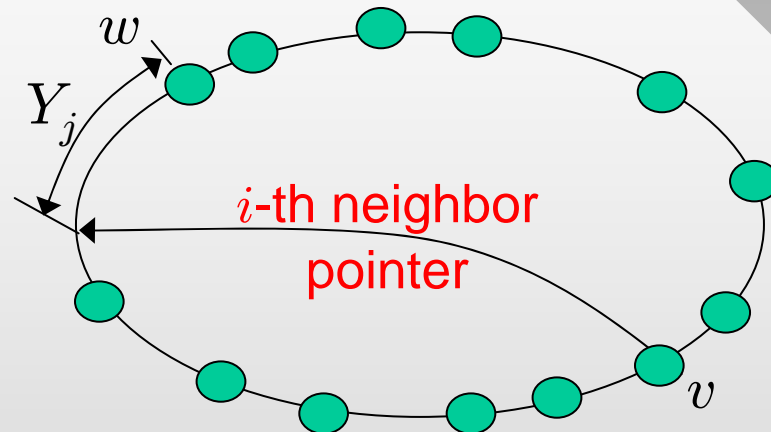
Conditional link lifetime

Conditional on zone size

11

# Conditional Link Lifetimes

- <u>Theorem 1</u>: Process $\{A_t{}^y\}$ is a <span style="color:red">semi-Markov</span> process where the sojourn time $\tau_i$ in each state $A_t{}^y = i$ follows a <span style="color:red">general</span> distribution

  – This process is fully determined by the distribution of residual lifetime $Z$ of the initial neighbor, the distribution of user lifetimes $L$, and the arrival process of new users splitting the given zone

  – Based on this semi-Markov chain, one can obtain the distribution $P(R(y)<x)$ and expectation $E[R(y)]$

- Next, we focus on the distribution of zone size $Y_j$ to get final results on link lifetimes

# Agenda

- Background and Motivation
  - Terminology and related work

- Link Lifetime Model for switching systems

- Deterministic Links

- Randomized Links

- Wrap-up

13

# Deterministic DHTs

- <u>Theorem 2</u>: In deterministic DHTs, the limiting distribution of $Y_1$ is <span style="color:blue">exponential</span> with mean $1/E[N]$ and that of $Y_j$ for $j = 2, 3, \ldots$ is <span style="color:blue">Erlang-$2$</span> with mean $2/E[N]$ as system population $N$ becomes sufficiently large
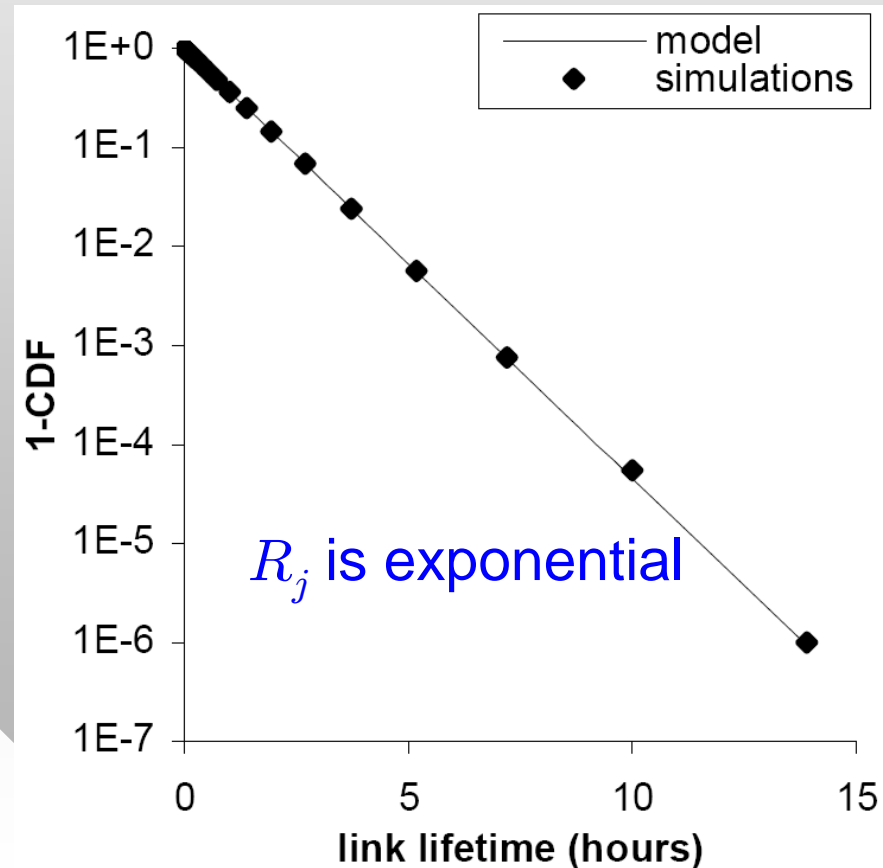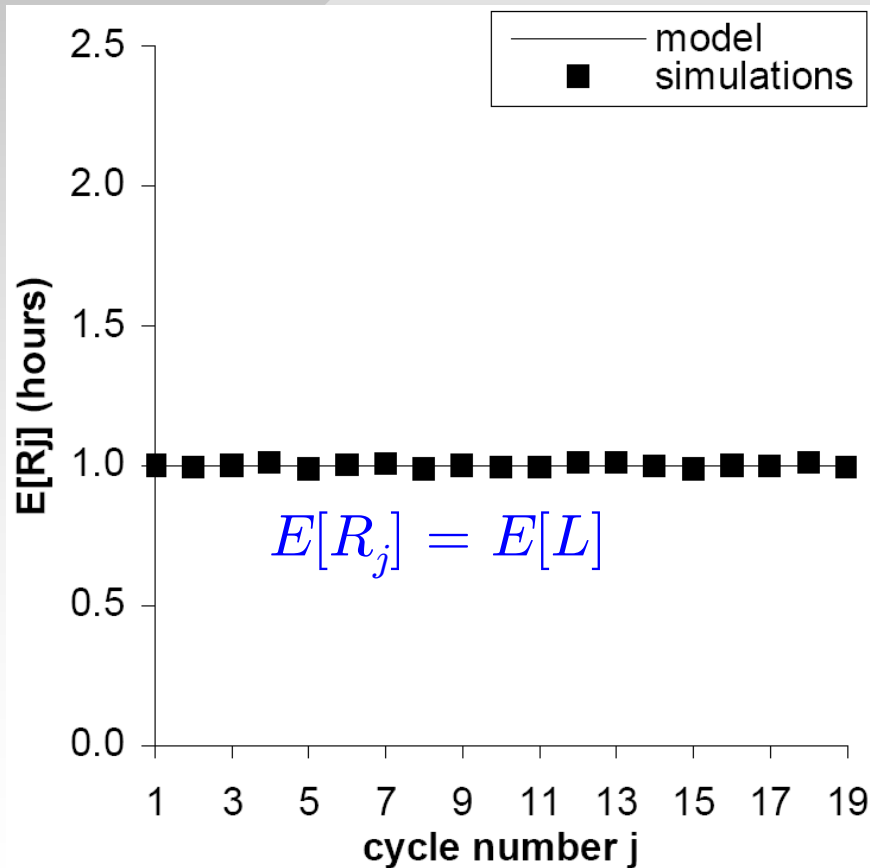
- The mean link lifetime is given by:

$$E[R_j] = \int_0^\infty E[R(y)] f_{Y_j}(y) dy$$

the mean conditional link lifetime

the PDF of $Y_j$
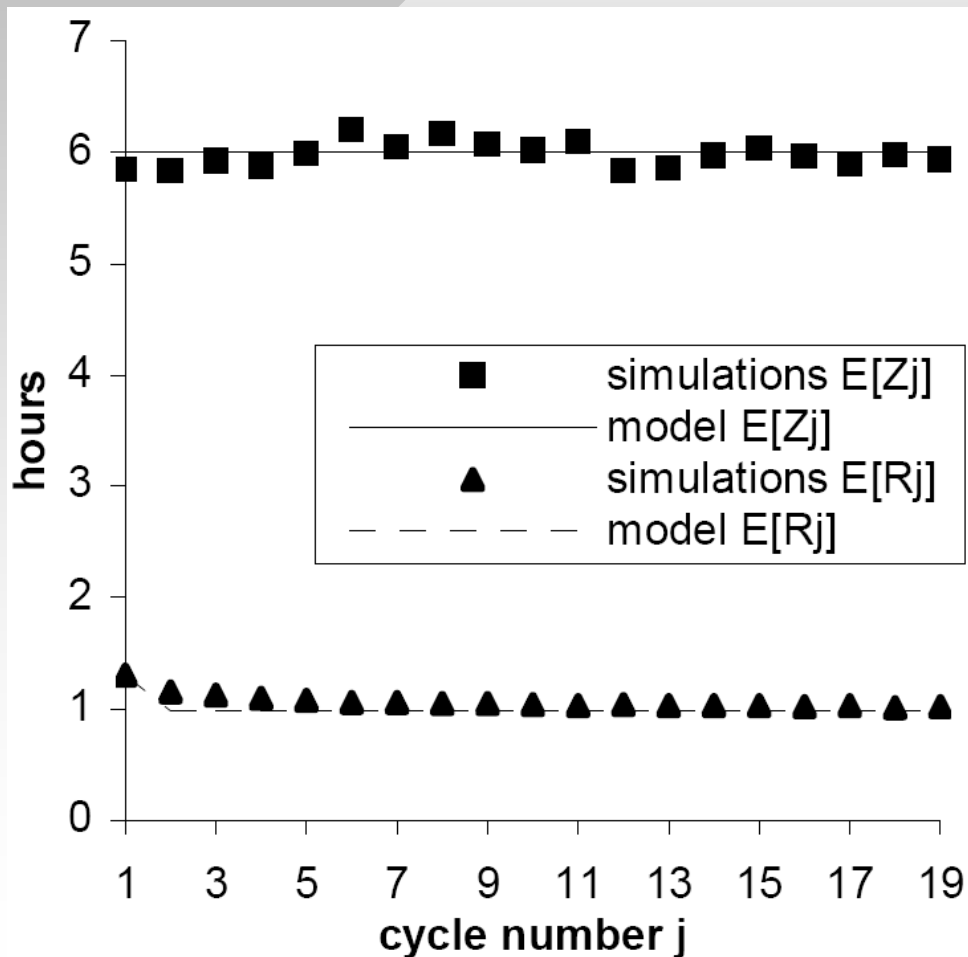
14

# Exponential User Lifetimes

- Properties of link lifetimes $R_j$ for exponential user lifetimes with $E[L] = 1$ hour



$$E[R_j] = E[L]$$

$R_j$ is exponential

15

# Pareto User Lifetimes

Pareto $L$ with $\alpha = 2.2$ and $E[L] = 1$ hour



- The initial neighbor is reliable since $E[Z_j] > E[L]$

- $E[R_j]$ is very close to $E[L]$ for $j = 2, 3, \ldots$

- $E[R_1] > E[R_2]$ since $Y_1$ is stochastically smaller than $Y_2$
  – A smaller zone size leads to a larger mean link lifetime

16

# Discussion

- Our model shows that link lifetime $R$ in deterministic DHTs is stochastically <span style="color:red">smaller</span> than residual lifetime $Z$ of the initial neighbor holding the link
  - Switching to newly arriving users makes $R$ smaller
  - Unlike non-switching systems, classic DHTs do not obtain benefits from heavy-tailed $L$

- Abandon switching systems?
  - Non-switching DHTs create inconsistence in routing tables and may expect longer routing delay

- We propose a new method that not only retains the advantage of switching systems, but increases link lifetimes

17

# Agenda

- Background and Motivation
  - Terminology and related work

- Link Lifetime Model for switching systems
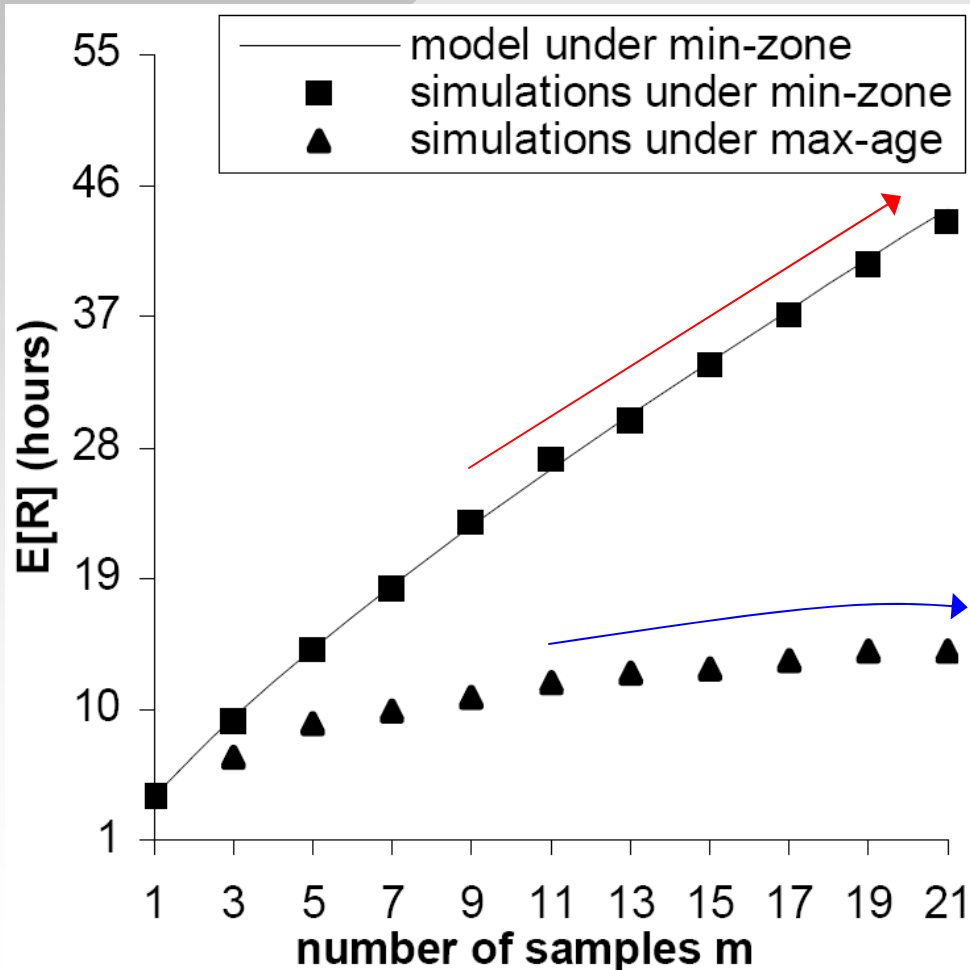
- Deterministic Links

- Randomized Links

- Wrap-up

# Improvement: Randomized Links

- We utilize the freedom of selecting links in randomized DHTs to propose min-zone selection
  - User $v$ uniformly samples $m$ points in $[\text{id}(v) + 2^i/2^{64}, \text{id}(v) + 2^{i+1}/2^{64}]$, and then selects the point whose successor has the minimum zone size
  - Upon link failure, user $v$ uses the same strategy to find a replacement
  - Zone size $Y_j$ is exponential but has a smaller mean $E[Y_j] = 1/(mE[N])$, where $N$ is system population, for all $j$

- For comparison purpose, we also examine max-age selection
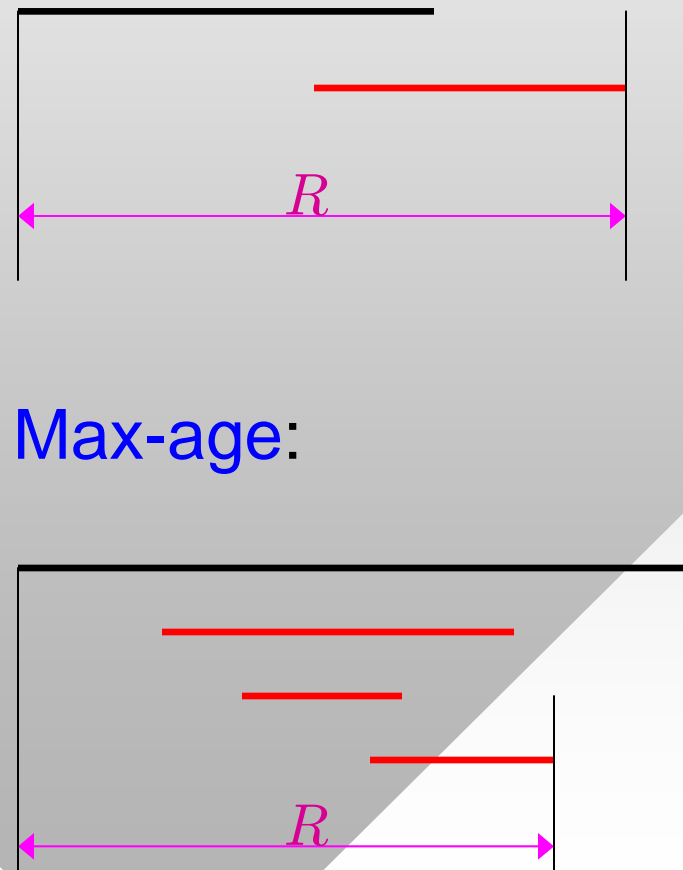  - The only difference is that age is used as selection criteria

# Link Lifetimes under Min-Zone Selection



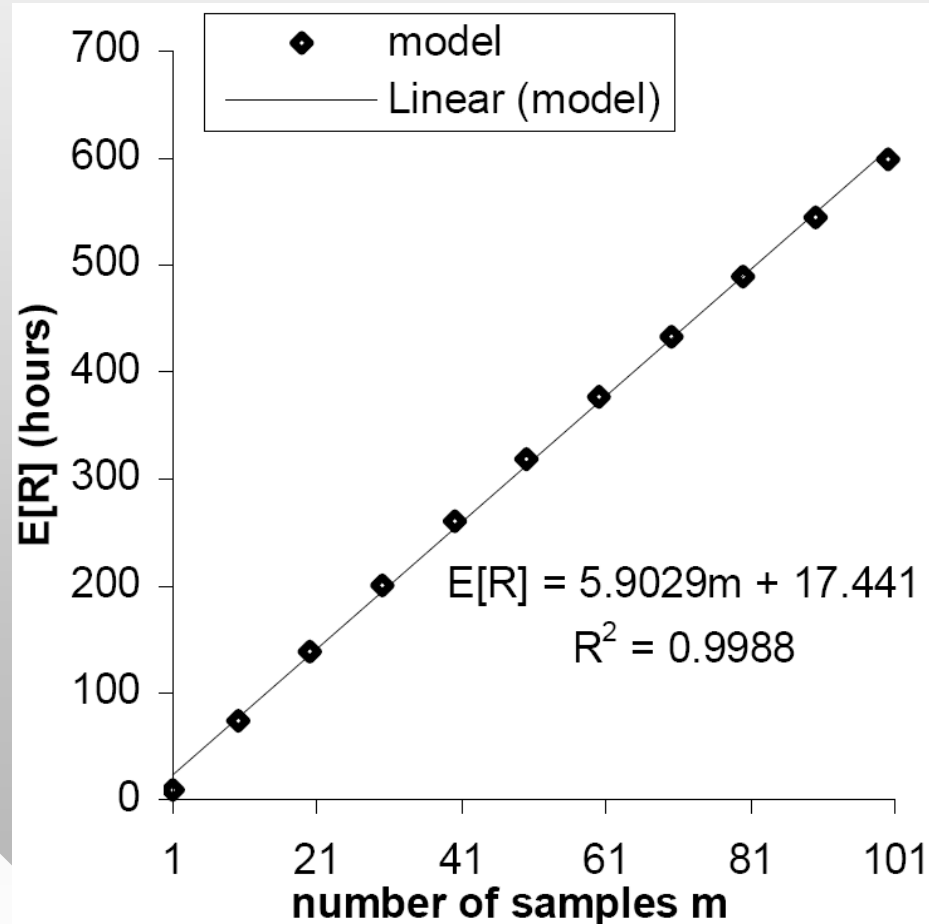Pareto $L$ with $\alpha = 1.2$ and $E[L] = 1$ hour

- **Min-zone:**

- **Max-age:**





20

# Link Lifetimes for $\alpha \leq 2$

- **Theorem 3**: For Pareto user lifetimes with shape $1 < \alpha \leq 2$ and <span style="color:red">min-zone</span> selection, $E[R] \to \infty$ as the <span style="color:blue">system size</span> and $m$ approach $\infty$. For <span style="color:blue">max-age</span> selection and any $\alpha > 1$, $E[R]$ is finite.



$E[R] = 5.9029m + 17.441$
$R^2 = 0.9988$

Pareto $L$ with $\alpha = 1.09$ and $E[L] = 1$ hour

21

# Wrap-up

- We developed a model for link lifetimes $R$ in DHTs
  - The mean link lifetime in deterministic DHTs is very close to the mean user lifetime
  - Switching leads to smaller link lifetimes

- We proposed min-zone selection which sufficiently increases $R$ for heavy-tailed user lifetimes
  - It allows us to achieve a spectrum of neighbor selection strategies while keeps routing tables consistent
  - For $m = 1$, it is the regular switching in DHTs
  - For $m = \infty$, the probability of switching is reduced to be $0$
  - Additionally, it benefits DHTs by balancing load such that users with smaller zone sizes are responsible for fewer keys while forwarding more queries