# A Clean-Slate Look at Disk Scrubbing

## Alina Oprea and Ari Juels
RSA Laboratories

Presented by

## Arkady Kanevsky

# Risk of Data Loss in Hard Drives
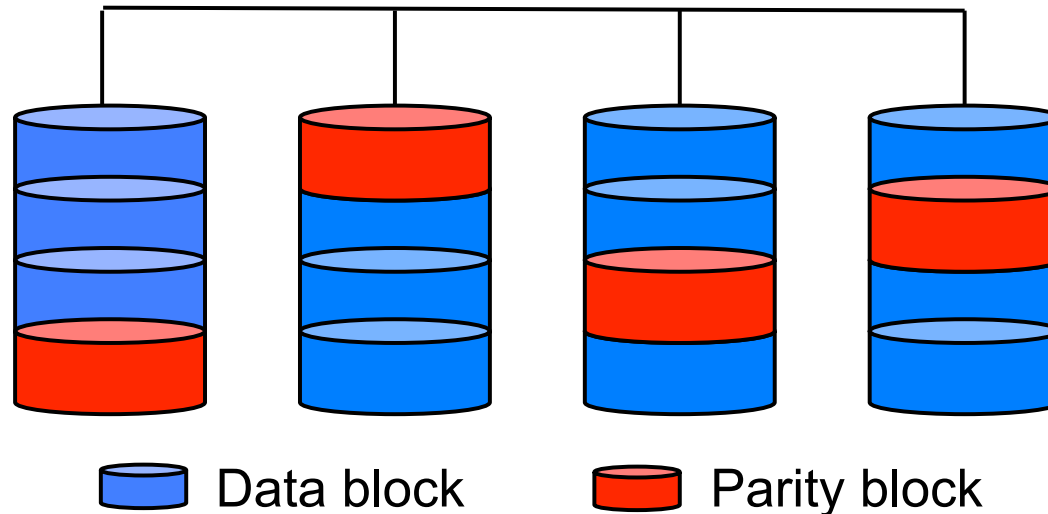
Drives fail!

Total
Crash

Bad sector!

- Latent sector errors (LSEs)
- Discovered only when sector is read

# RAID and LSEs

RAID-5

Data block    Parity block

- RAID-5 protects against one disk failure

- But…one disk failure + one LSE result in data loss

- Impact of LSEs on RAID reliability
  - [Elerath and Pecht 2007], [Baker et al. 2007]

# Mitigations Within a Single Drive

**EMC²**
where information lives®

- **Disk scrubbing [Schwartz et al. 2004]**
  - Background process that reads disk sectors during disk idle time to proactively discover LSEs

- **Intra-disk redundancy [Dholakia et al. 2008]**
  - Erasure code over consecutive disk sectors
  - Parity blocks stored on the same drive
  - Incurs write overhead

- Comparison of scrubbing and intra-disk redundancy
  - [Iliadis et al. 2008], [Mi et al. 2008], [Schroeder et al. 2010]

4

# Disk Scrubbing

**EMC²**
*where information lives®*

- Today: sequential reading of disk sectors, usually with a fixed pre-determined rate

- But LSEs do not occur with a fixed rate and uniformly across disk sectors (Sigmetrics 2007 study)

  – Temporal decay: subsequent errors develop after first LSE

  – Temporal locality: most errors occur within a short interval of previous error

  – Spatial locality: 50% of LSEs are at a logical distance of 10MB

Idea: enlarged design space of scrubbing strategies to account for distribution of LSEs and disk history

  ▪ Adaptively change scrubbing rate following error event

  ▪ Sample across disk regions for discovering errors faster than by sequential reading ("staggering")

# Outline

**EMC²**
where information lives®

- Motivation for more intelligent disk scrubbing techniques

- Our LSE model
  - Use known facts about LSE distribution (spatial and temporal locality)
  - New assumptions for usage error development

- Enlarged design space of scrubbing strategies
  - Staggered strategies
  - Strategies with adaptive rates

- Simulation model and evaluation
  - New metric for single drive reliability (MLET)
  - Reliability dependence on various disk parameters, and disk workloads

# Methodology for our LSE model

**EMC**[2]
where information lives®

- Published facts on LSE distribution
  - [Bairavasundaram et al. 2007] study on 1.53 million drives from various models and manufacturers over 24 month period
  - Two disk categories: nearline and enterprise
  - Consider only enterprise disks in our work
  - Data is not published, only some statistics on it
  - Translate known facts into scrubbing principles

- Need new assumptions to generate LSE model
  - Parameterized model aimed at capturing disks with various characteristics
  - Actual disk parameters are currently not transparent

- Validate LSE model against data published by Bairavasundaram et al.

# Scrubbing principles

## Almost constant LSE rates

| [Bairavasundaram et al. 2007] | Scrubbing principles |
|---|---|
| • LSE rate is fairly low and constant in first 2 months of drive operation<br><br>• LSEs rate increases after 2 months, but is fairly constant before the first LSE develops | • Keep scrubbing rate low and constant during first 2 months<br><br>• Increase scrubbing rate after 2 months, and keep it constant before the first LSE develops |

# Scrubbing principles

## Temporal locality and decay of LSEs

### [Bairavasundaram et al. 2007]

- LSEs exhibit temporal locality – inter-arrival time distribution has very long tails

- LSEs exhibit decay – more LSEs develop shortly after a first LSE

### Scrubbing principles

- Use adaptive scrubbing rates
  - Increase scrubbing rate temporarily in a short interval after LSE detection
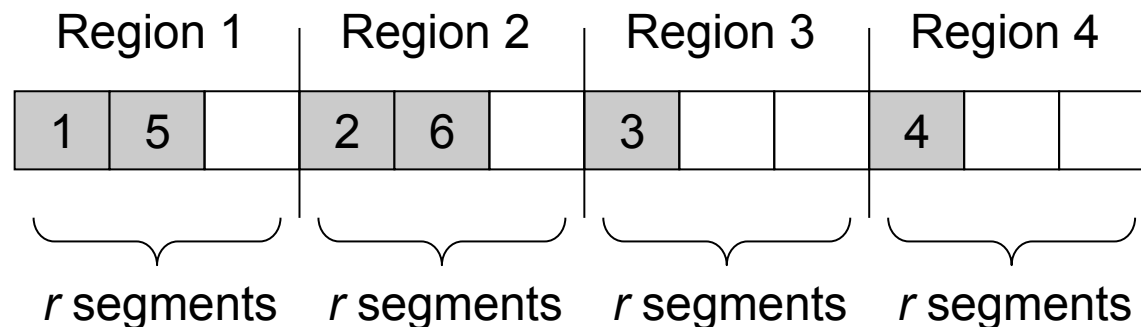
# Scrubbing principles
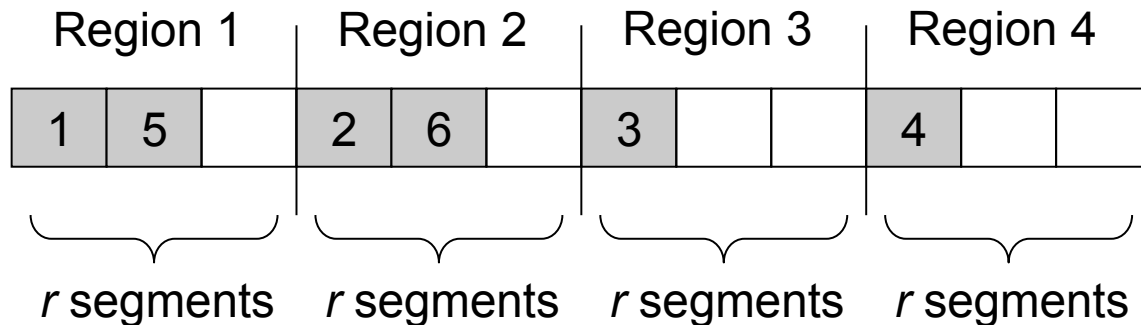
## Spatial locality of LSEs

| [Bairavasundaram et al. 2007] | Scrubbing principles |
|---|---|
| • LSEs develop clustered on disk at block logical level | • Staggering detects errors faster than sequential scrubbing |

Region 1    Region 2    Region 3    Region 4

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

$r$ segments    $r$ segments    $r$ segments    $r$ segments

# Staggered strategy

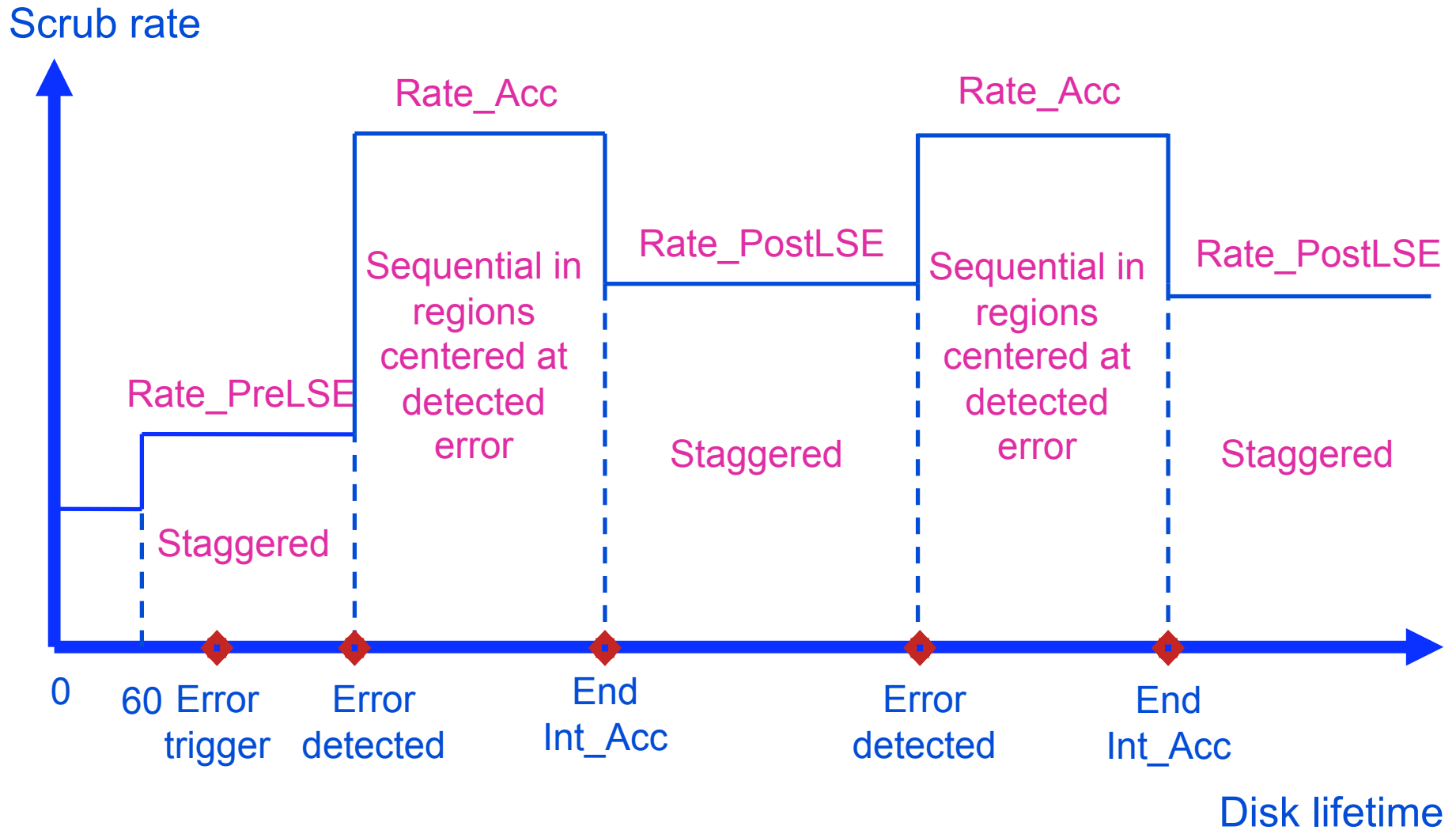| Region 1 | Region 2 | Region 3 | Region 4 |
|---|---|---|---|
| 1  5 | 2  6 | 3 | 4 |

*r* segments    *r* segments    *r* segments    *r* segments

- Performance overhead compared to sequential scrubbing
  - Small segment sizes (32-64KB): a factor of 5
  - Large segment sizes (1MB): only 2% overhead

- Parameter choices
  - Segment size 1MB
  - Region size 128MB: most LSEs are at distance lower than 128MB

# Staggered Adaptive Strategies

Scrub rate

Rate_Acc

Rate_Acc

Rate_PostLSE

Sequential in regions centered at detected error

Rate_PostLSE

Rate_PreLSE

Sequential in regions centered at detected error

Staggered

Staggered

Staggered

0    60 Error trigger    Error detected    End Int_Acc    Error detected    End Int_Acc

Disk lifetime

12

# Assumptions on LSE development

**EMC²**
where information lives®

- Usage error development
  - Bairavasundaram et al. study only characterizes age errors
  - Usage errors exhibit same spatial and temporal locality
  - Usage errors develop due to both reads and writes, albeit with different weights given by a parameter RW_Weight
  - Increase RW_Weight to minimize effect of reads on LSE development

- Usage errors are triggered when number of bytes accessed (weighted by RW_Weight) exceeds on average 1/BER
  - BER: byte-error rate, between $[10^{-15}, 10^{-13}]$

- Error distribution on disk
  - Errors are clustered on disk around a cluster centroid
  - Clusters of errors are uniformly distributed on disk

- Assumptions validated against Bairavasundaram et al. results

13

# Outline

**EMC²**
where information lives®

- Motivation for more intelligent disk scrubbing techniques

- Our LSE model
  – Use known facts about LSE distribution (spatial and temporal locality)
  – New assumptions for usage error development

- Enlarged design space of scrubbing strategies
  – Staggered strategies
  – Strategies with adaptive rates

- **Simulation model and evaluation**
  – New metric for single drive reliability (MLET)
  – Reliability dependence on various disk parameters, and disk workloads

# Simulation Model

**EMC²**
where information lives®

## Disk Model

- 24 months, interval of one hour

- 100,000 disks, 500GB each

- LSE model includes both age and usage errors

## Staggered adaptive space

- Scrubbing rates from 0 to one full disk scrub per day (in GB/hour)

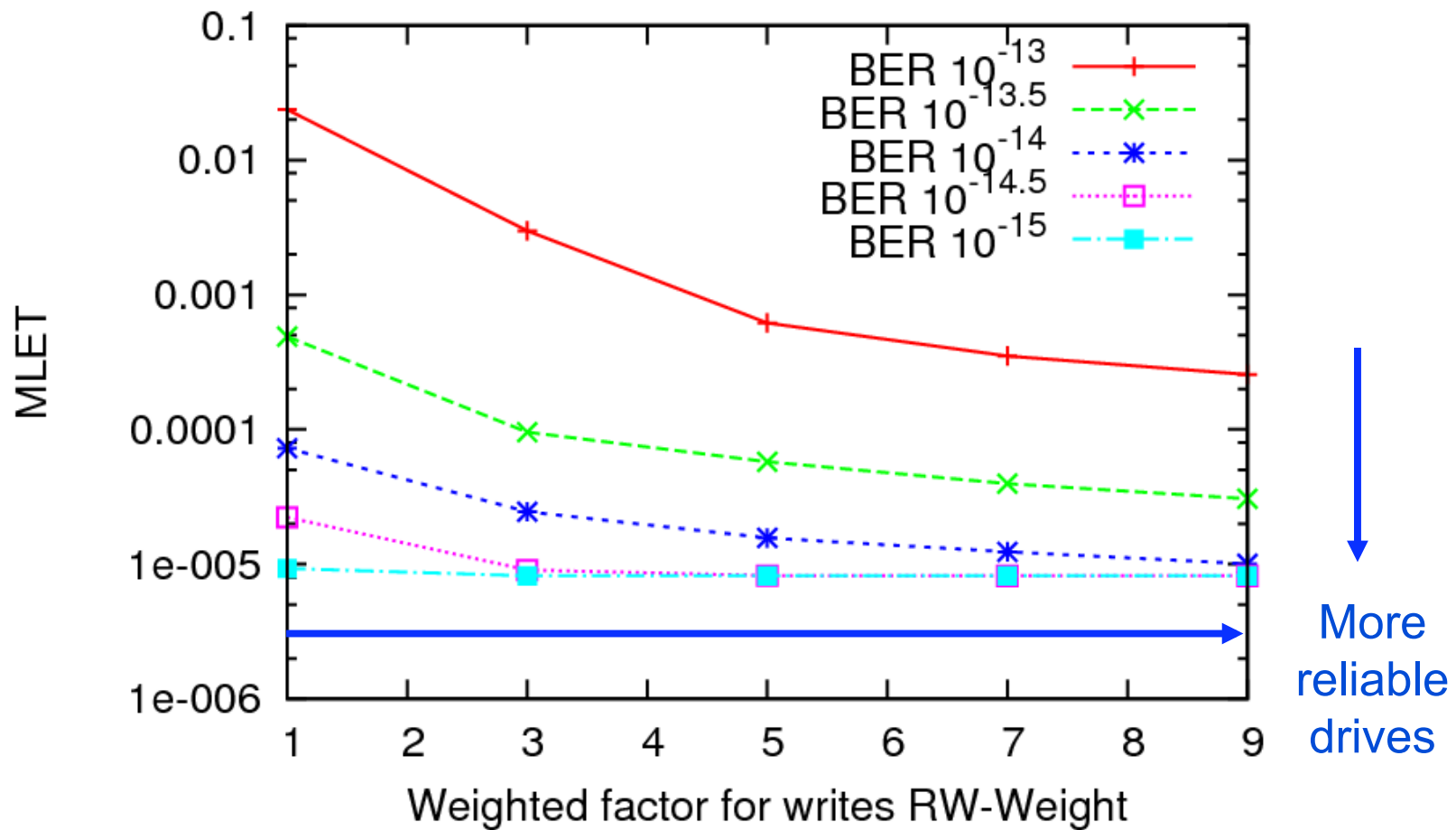- Length of accelerated interval from 3 hours to time to scrub the full disk

Exhaustive search for optimized scrubbing

**Optimized**: Min **MLET** (**M**ean **L**atent **E**rror **T**ime)
Fraction of time disk has latent sector errors
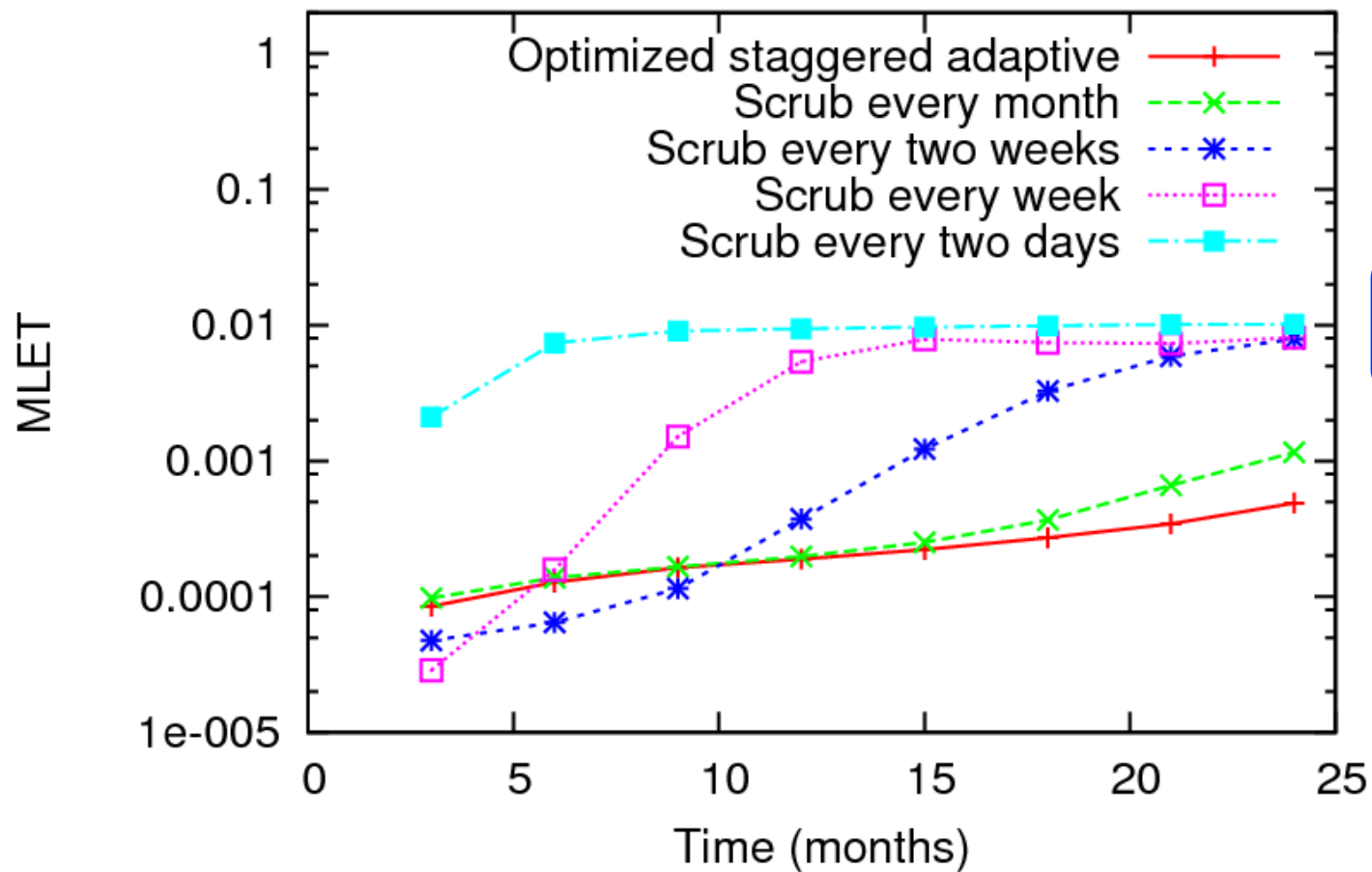
# Dependence on BER

Optimized staggered adaptive strategy

Legend:
- BER $10^{-13}$ — red, +
- BER $10^{-13.5}$ — green, ×
- BER $10^{-14}$ — blue, *
- BER $10^{-14.5}$ — magenta, □
- BER $10^{-15}$ — cyan, ■

Y-axis: MLET
X-axis: Weighted factor for writes RW-Weight

More reliable drives

16

# Staggered adaptive vs sequential



MLET for BER=$10^{-13.5}$ and RW-Weight=1 → High number usage errors

Legend:
- Optimized staggered adaptive
- Scrub every month
- Scrub every two weeks
- Scrub every week
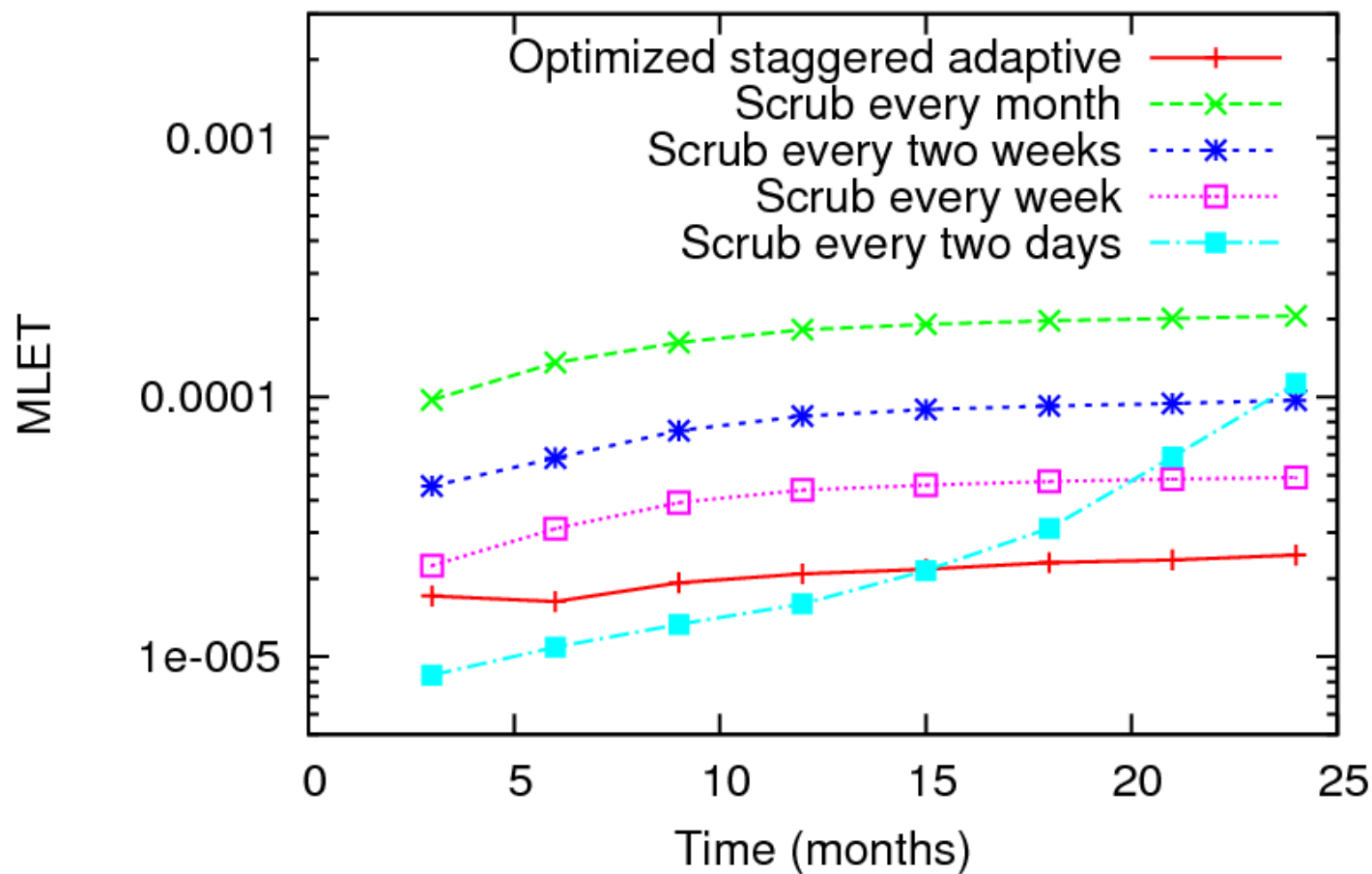- Scrub every two days

Scrub infrequently

# Staggered adaptive vs sequential
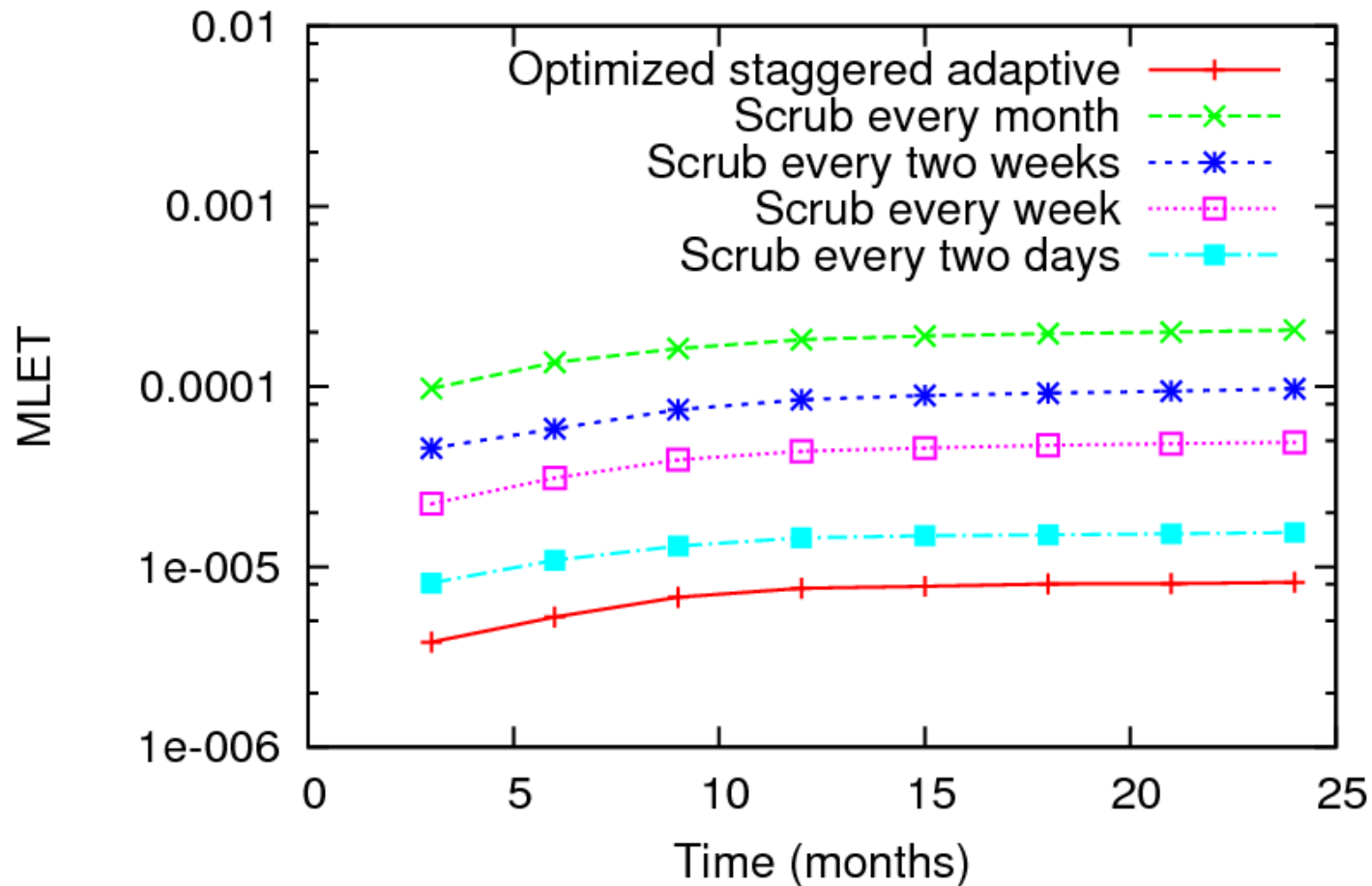
MLET for BER=$10^{-14}$ and RW-Weight=3 → Medium number usage errors

Scrub every two weeks

# Staggered adaptive vs sequential



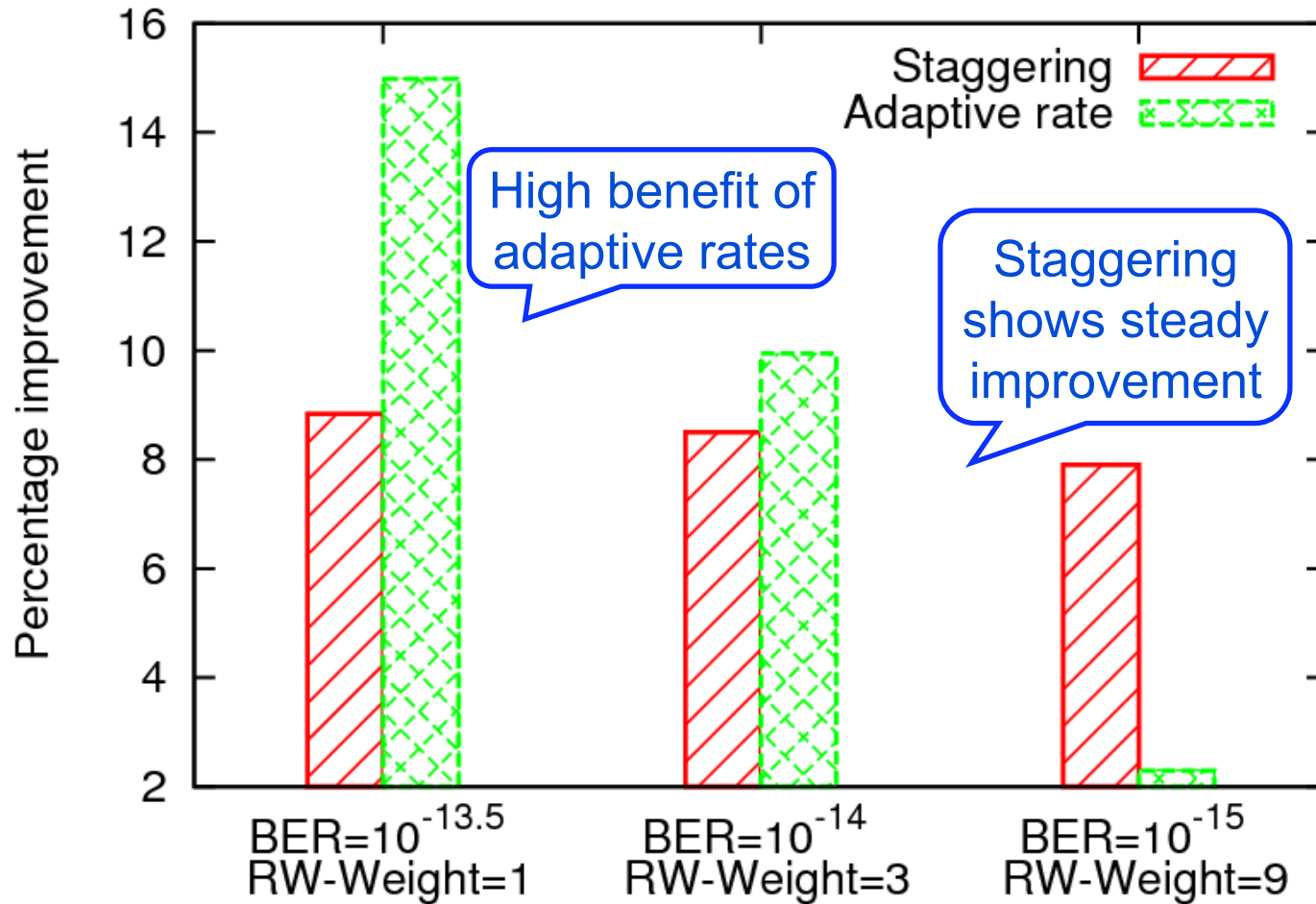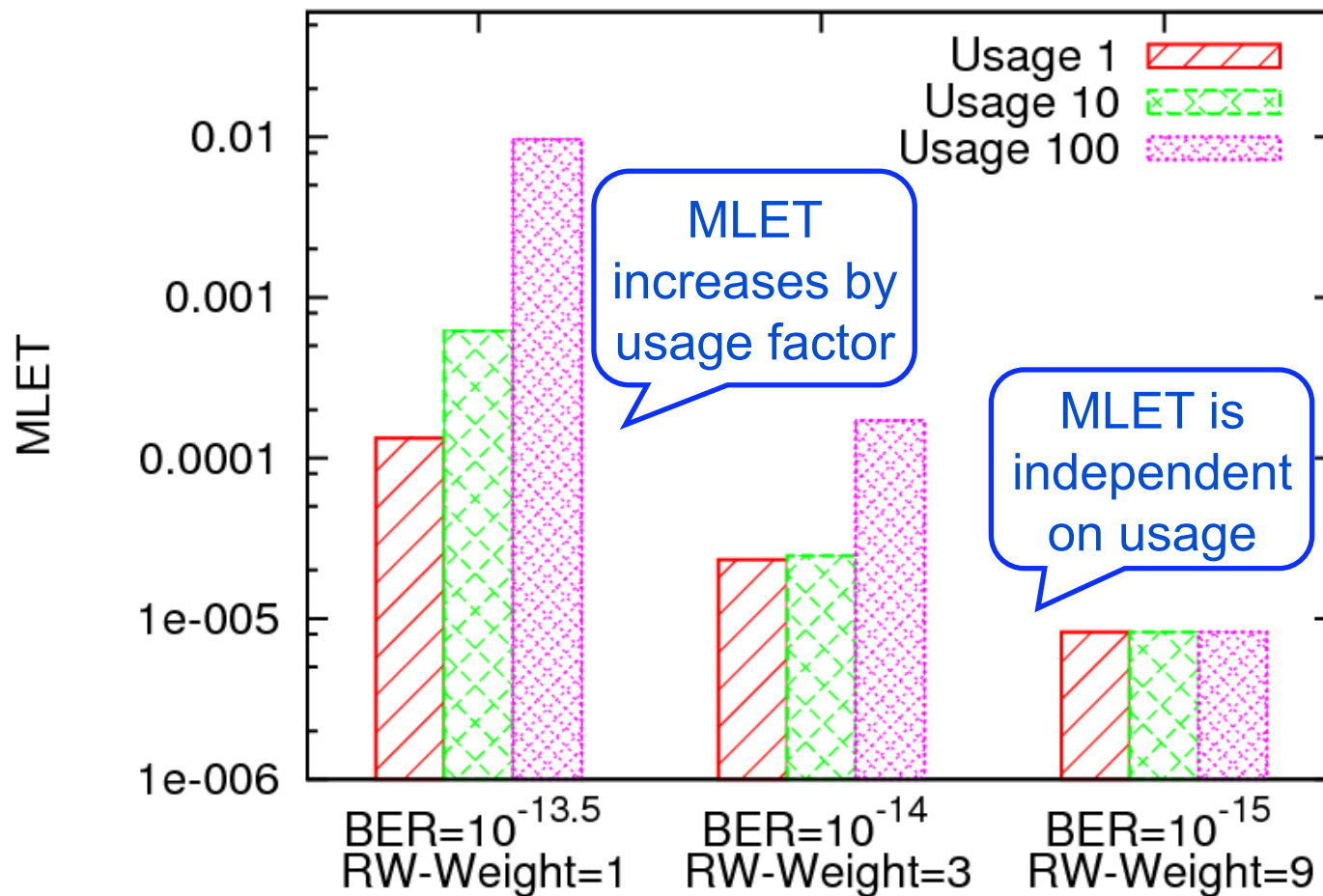MLET for BER=$10^{-15}$ and RW-Weight=9 → Low number usage errors

Legend:
- Optimized staggered adaptive
- Scrub every month
- Scrub every two weeks
- Scrub every week
- Scrub every two days

Scrub every two days

# Relative improvement of staggering and adaptive rates

**EMC²**
where information lives®



Improvement of staggering and adaptive rates over fixed-rate sequential

High benefit of adaptive rates

Staggering shows steady improvement

# Dependence on disk workload

# Discussion

**EMC²**
*where information lives*®

- More intelligent scrubbing strategies by taking into account disk characteristics and the history of error development

- Optimal strategies are highly dependent on disk BER and disk workloads
  - High sensitivity to disk parameters that are not always public

- Staggering improves resilience to LSEs for all disks

- Adaptively changing scrubbing rates in a short interval after detecting an LSE benefits most disks that develop a high number of usage errors

- Optimized adaptive staggered strategies can reduce MLET by several orders of magnitude compared to fixed-rate sequential strategies used today

# Future work

- Expansion of search space for scrubbing strategies

- Use more sophisticated search heuristics
  – E.g., hill-climbing or simulated annealing

- Performance overhead of real scrubbers in conjunction with typical workloads

- Translation of results to FLASH

- Extension of results to replication and RAID systems

- Questions?
  – Alina Oprea (aoprea@rsa.com)
  – Ari Juels (ajuels@rsa.com)