

---

# Advanced disk scheduling

## *“Freeblock scheduling”*

Eno Thereska

(slide contributions by Chris Lumb and  
Brandon Salmon)

PARALLEL DATA LABORATORY  
Carnegie Mellon University

# Why do I care?

---

- Windows Vista
  - background disk defragmenter
  - background backup
- Freeblock scheduling
  - one way to do the above (for free!)
  - minimal impact on foreground workloads

# Outline

---

- Freeblock scheduling: some theory
- Freeblock scheduling: applied
- Implementation considerations
- Q & A

# Some theory: preview

---

- Next few slides will review & show that:
  - disks are slow
    - mechanical delays (seek + rotational latencies)
  - there is nothing we can do during a seek
  - **there is a lot we can do during a rotation**
    - rotational latencies are very large
    - while rotation is happening go to nearby tracks and do useful work
  - “*freeblock scheduling*” = utilization of rotational latency gaps (+ any idle time)

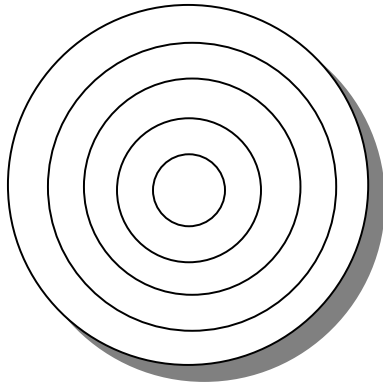
# Are disks slow?

---

- Are the xfer speeds that slow?
  - no, xfer speeds of 200MB/s are pretty good
- So what is slow?
  - workload often not sequential
  - disk head has to move from place to place
  - seek ( $\sim 4\text{ms}$ ) + rotation ( $\sim 3\text{ms}$ )
- Effective bandwidth can be very low
  - $\sim 10\text{-}30\text{MB/s}$ , even when SPTF is used
  - problem will exist for the next 10+ years

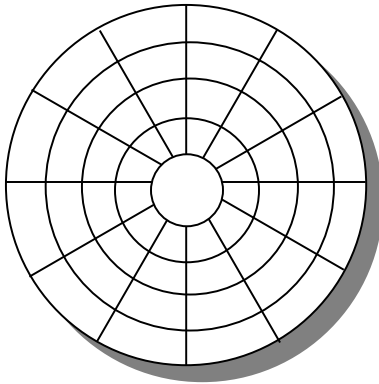
# Surface organized into tracks

---



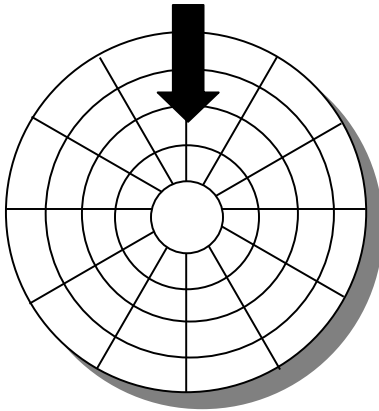
# Tracks broken up into sectors

---



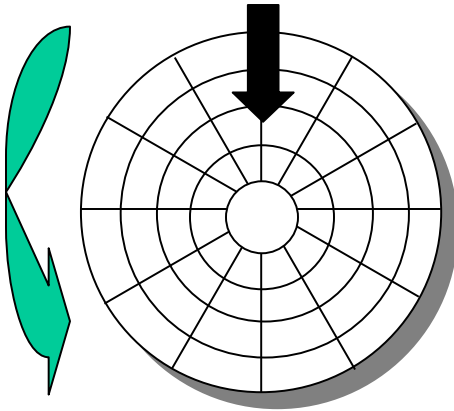
# Disk head position

---



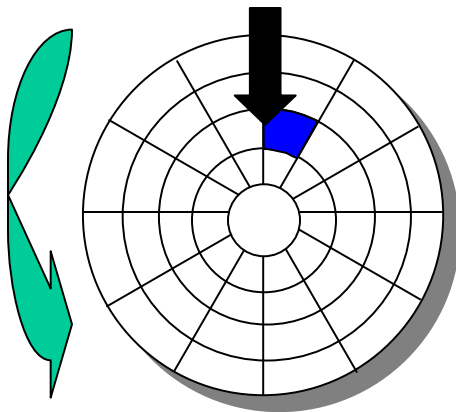
# Rotation is counter-clockwise

---



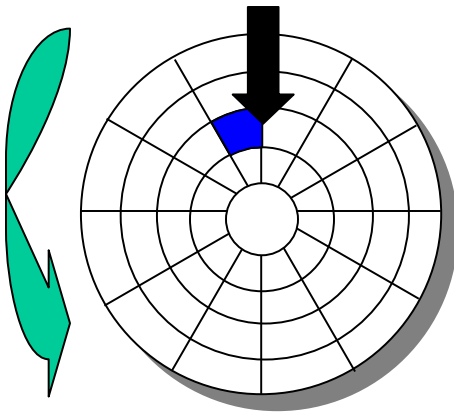
# About to read blue sector

---



# After reading blue sector

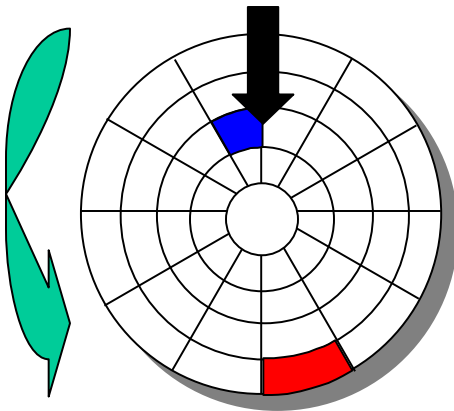
---



After **BLUE** read

# Red request scheduled next

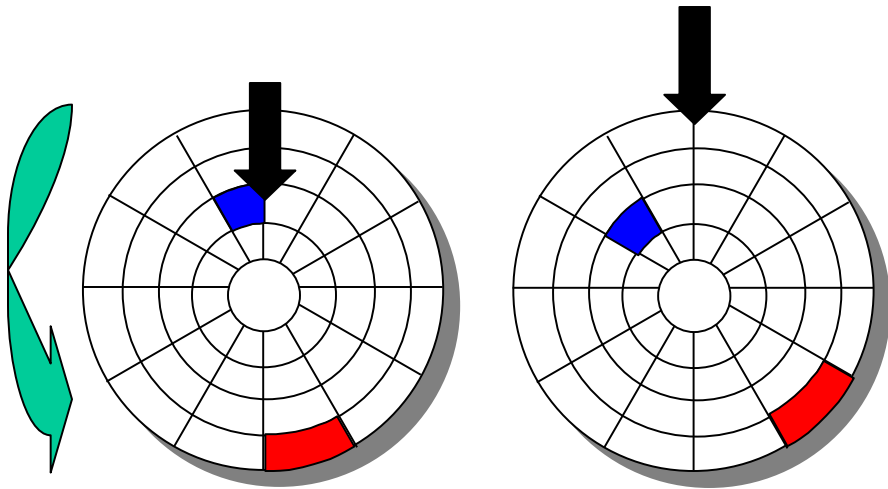
---



After **BLUE** read



# Seek to Red's track

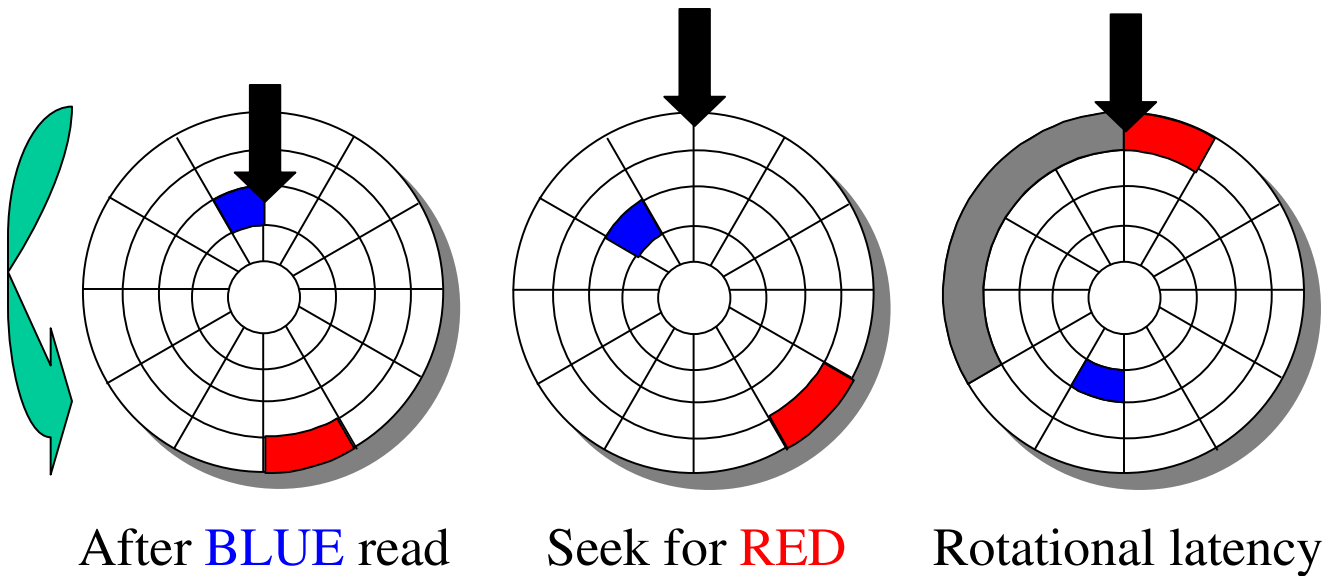


After **BLUE** read

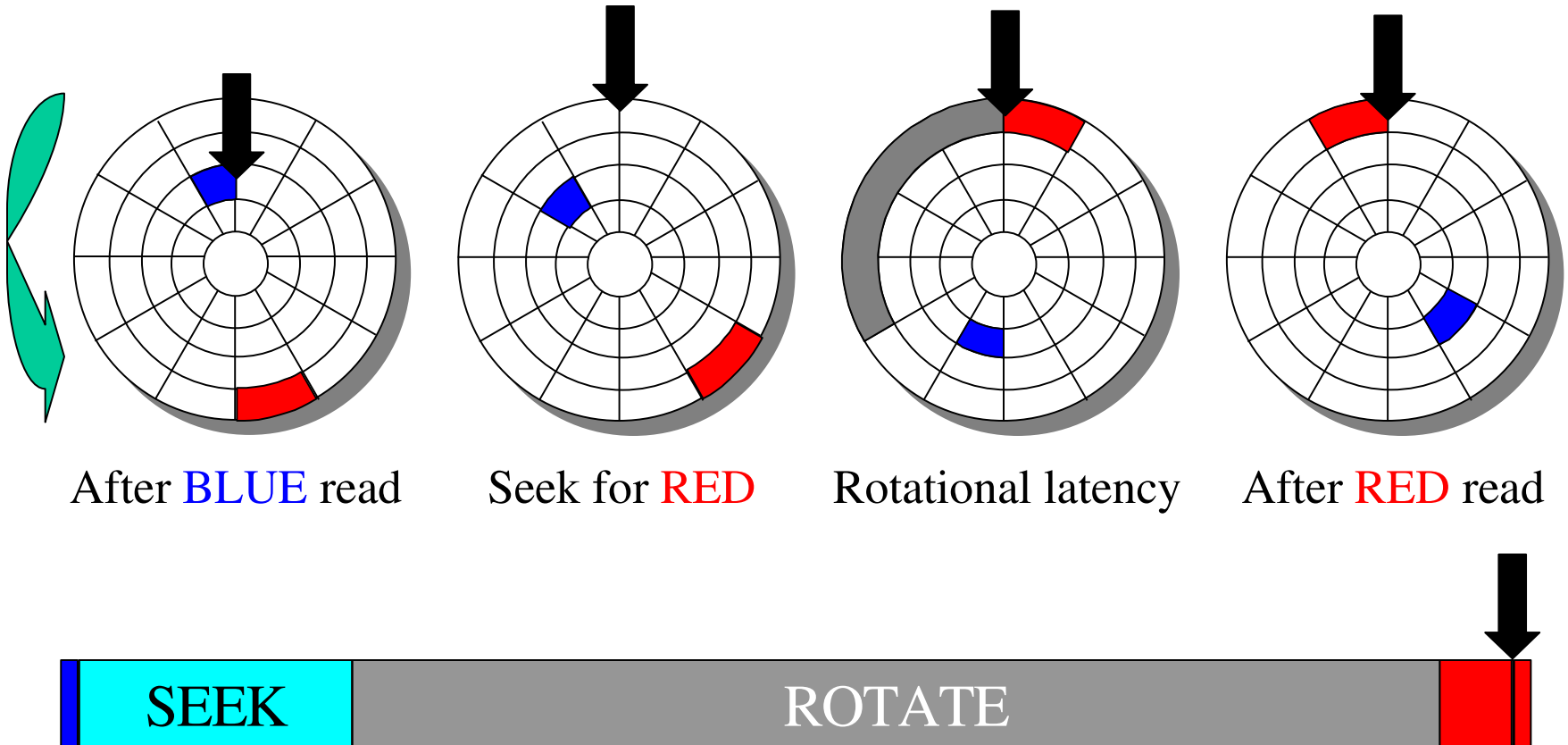
Seek for **RED**



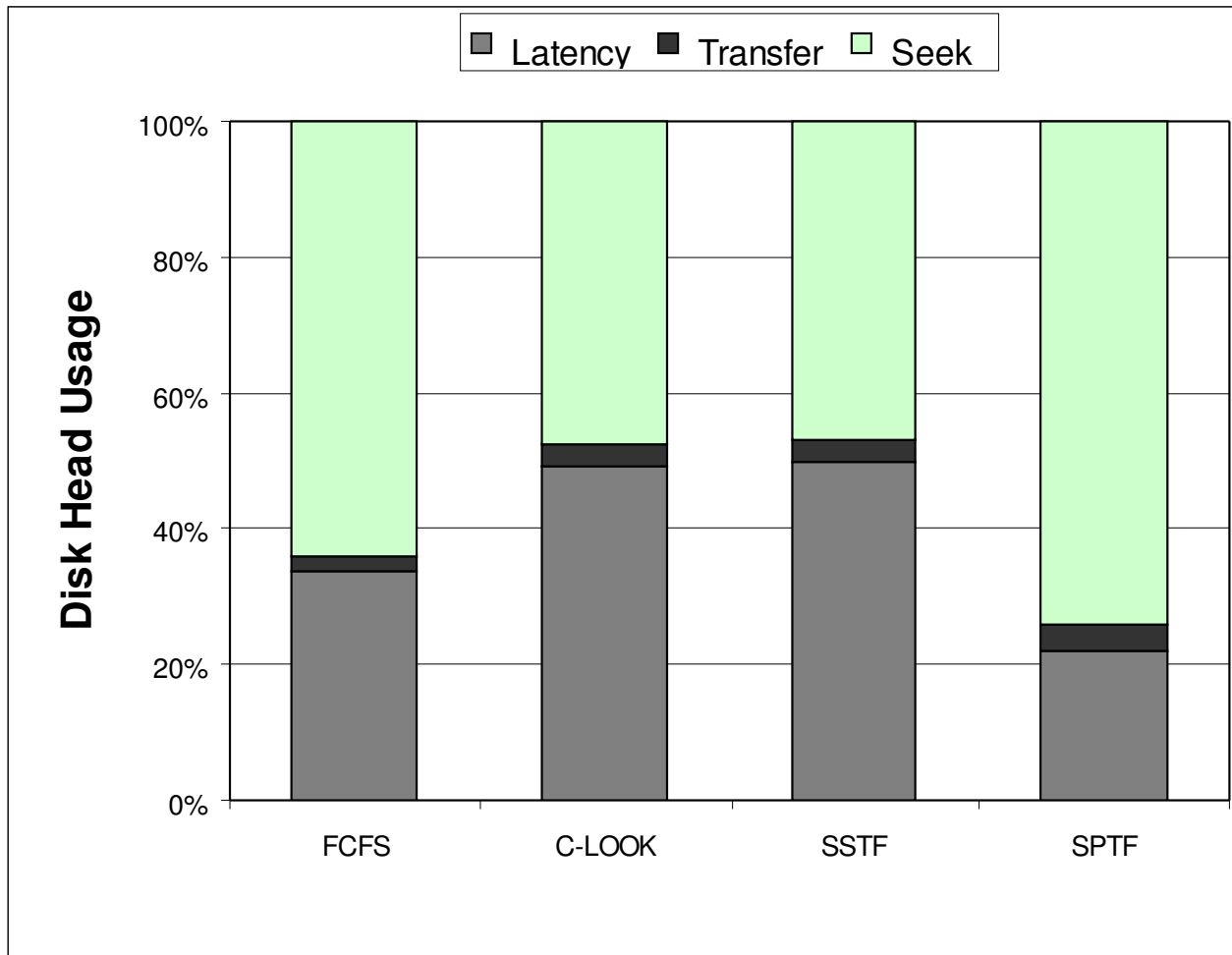
# Wait for Red sector to reach head



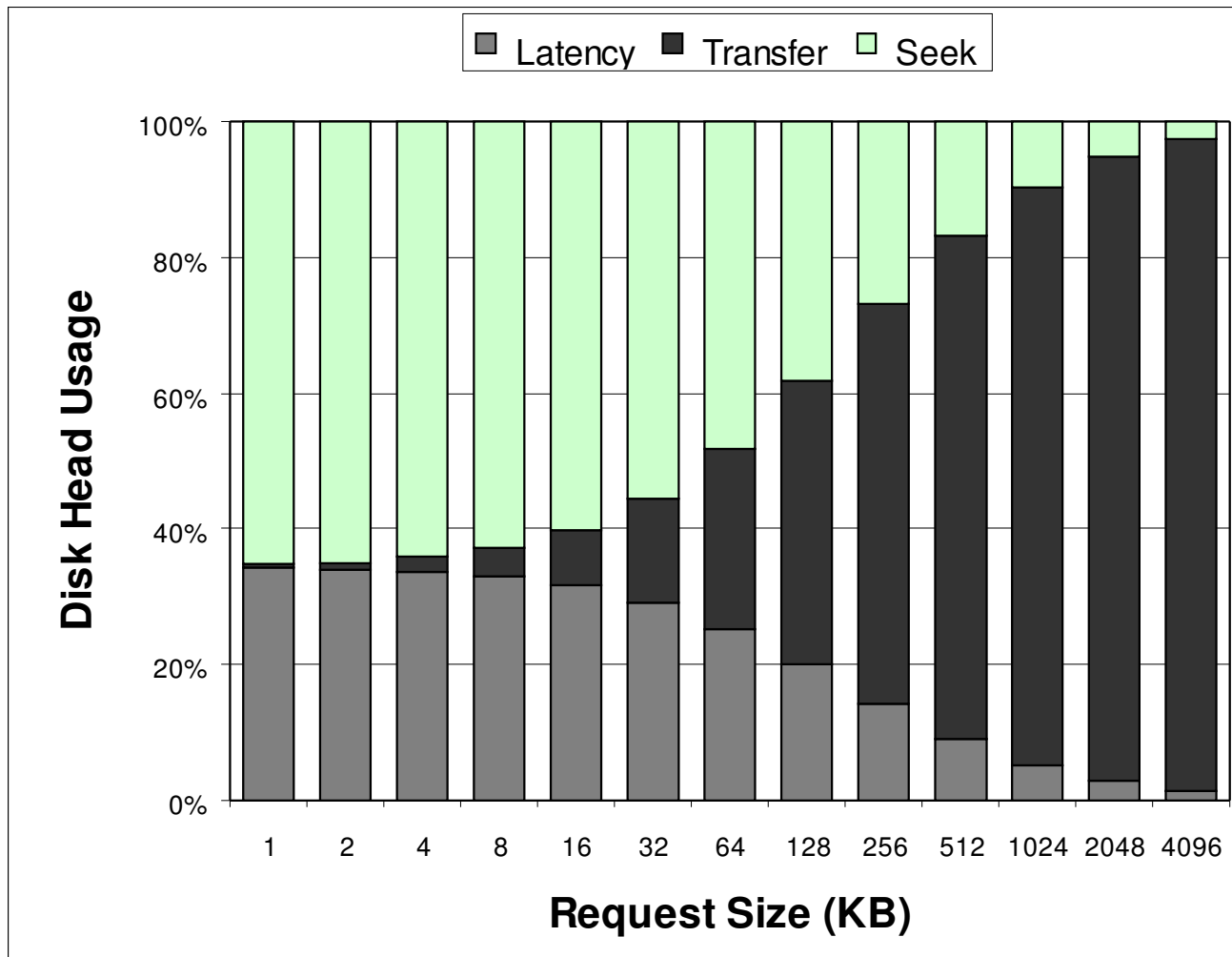
# Read Red sector



# Scheduling algorithm impact



# Impact of request sizes



# What can we do?

---

- Nothing we can do during a seek
  - disk head has to move to the right track
- Rotational latency is fully wasted
  - let's use this latency
- During a rotational latency
  - go to nearby tracks and do useful work
  - then, just-in-time, seek back to the original request

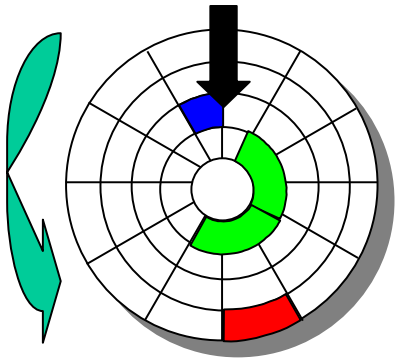
# A quick glance ahead...

---

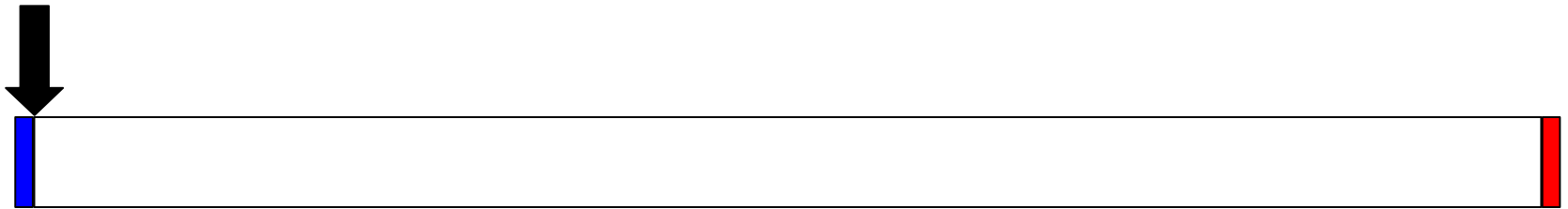
- What kind of “useful work” are we doing?
  - work that belongs to a “background” app
  - things like backup, defrag, virus scanning
- What do we really gain?
  - background apps don't interfere with fore. apps
  - background apps still complete
- What's in it for me?
  - can run defrag + virus scanner + backup in the background while working on your homework and you won't notice they are running ☺

# Rotational latency gap utilization

---

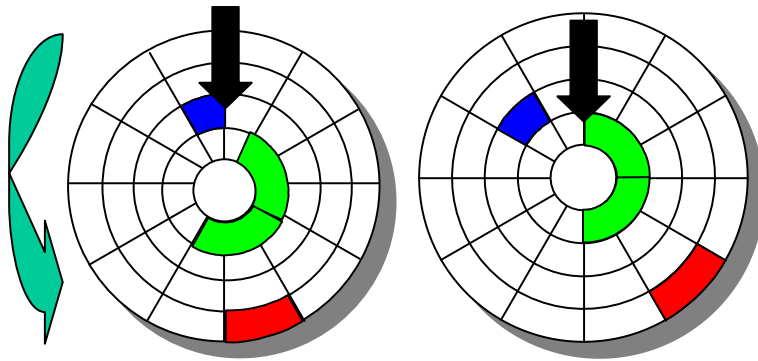


After BLUE read

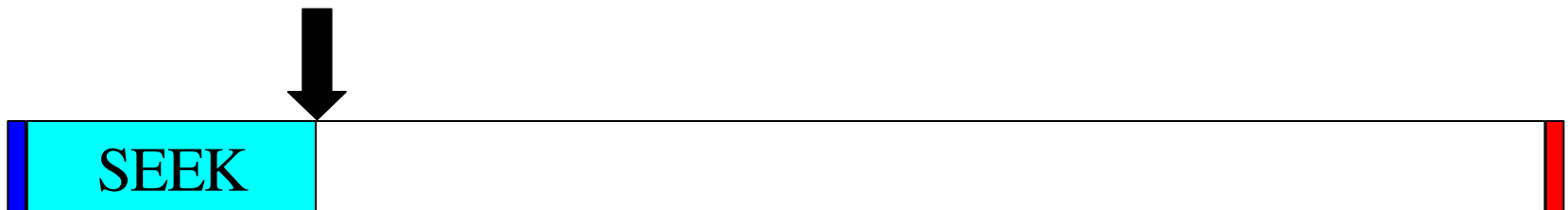


# Seek to Third track

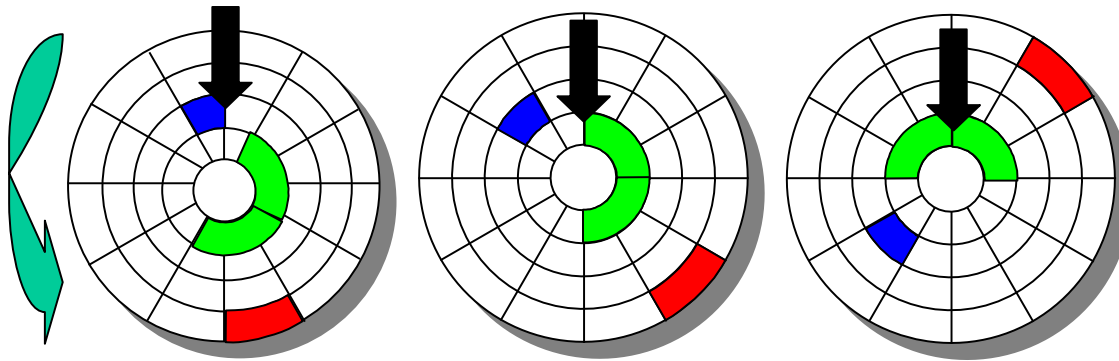
---



After BLUE read      Seek to Third



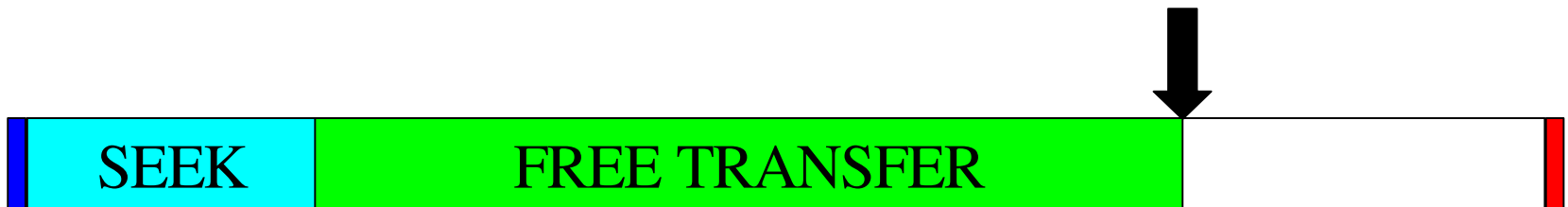
# Free transfer



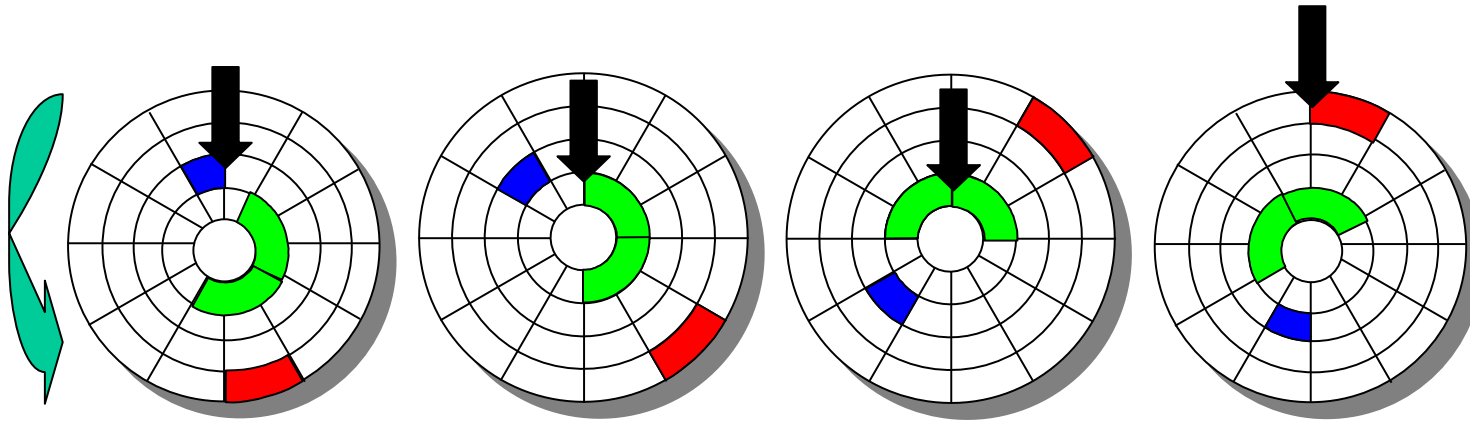
After **BLUE** read

Seek to Third

Free transfer



# Seek to Red's track



After **BLUE** read

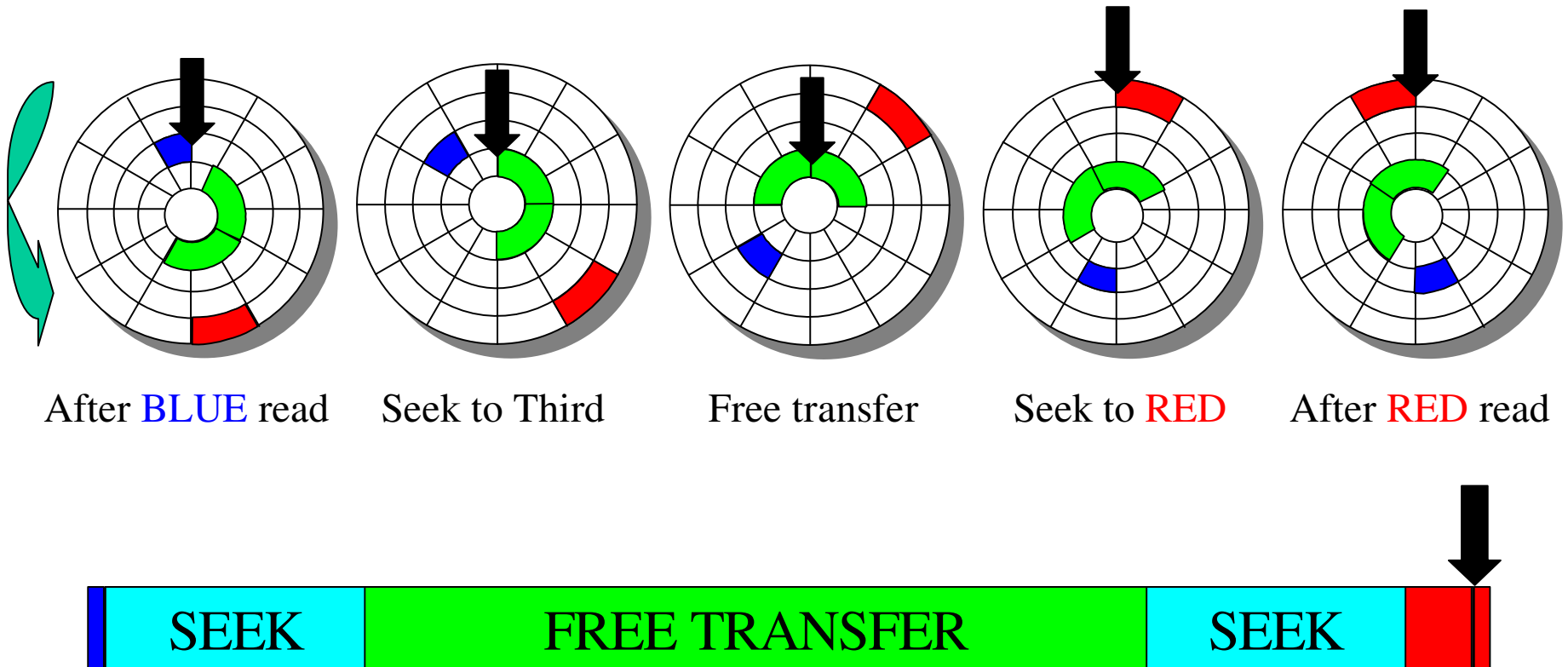
Seek to Third

Free transfer

Seek to **RED**



# Read Red sector



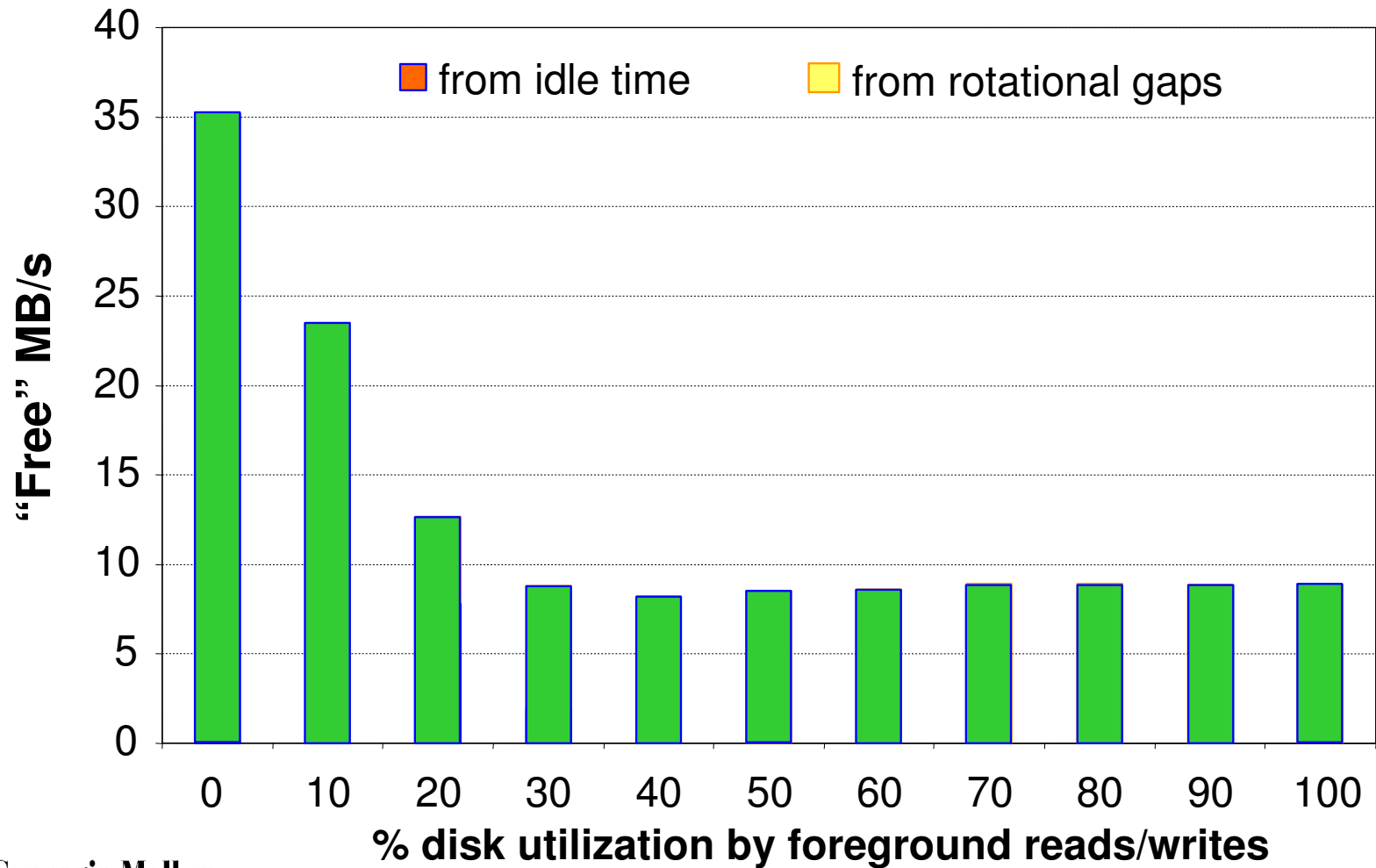
# Final theory details

---

- Scheduler also uses disk idle time
  - high end servers have little idle time
- Idle time + rotational latency usage = *“freeblock scheduling”*

(it means we are getting things for free)

# Steady background I/O progress



# Applied freeblocks: preview

---

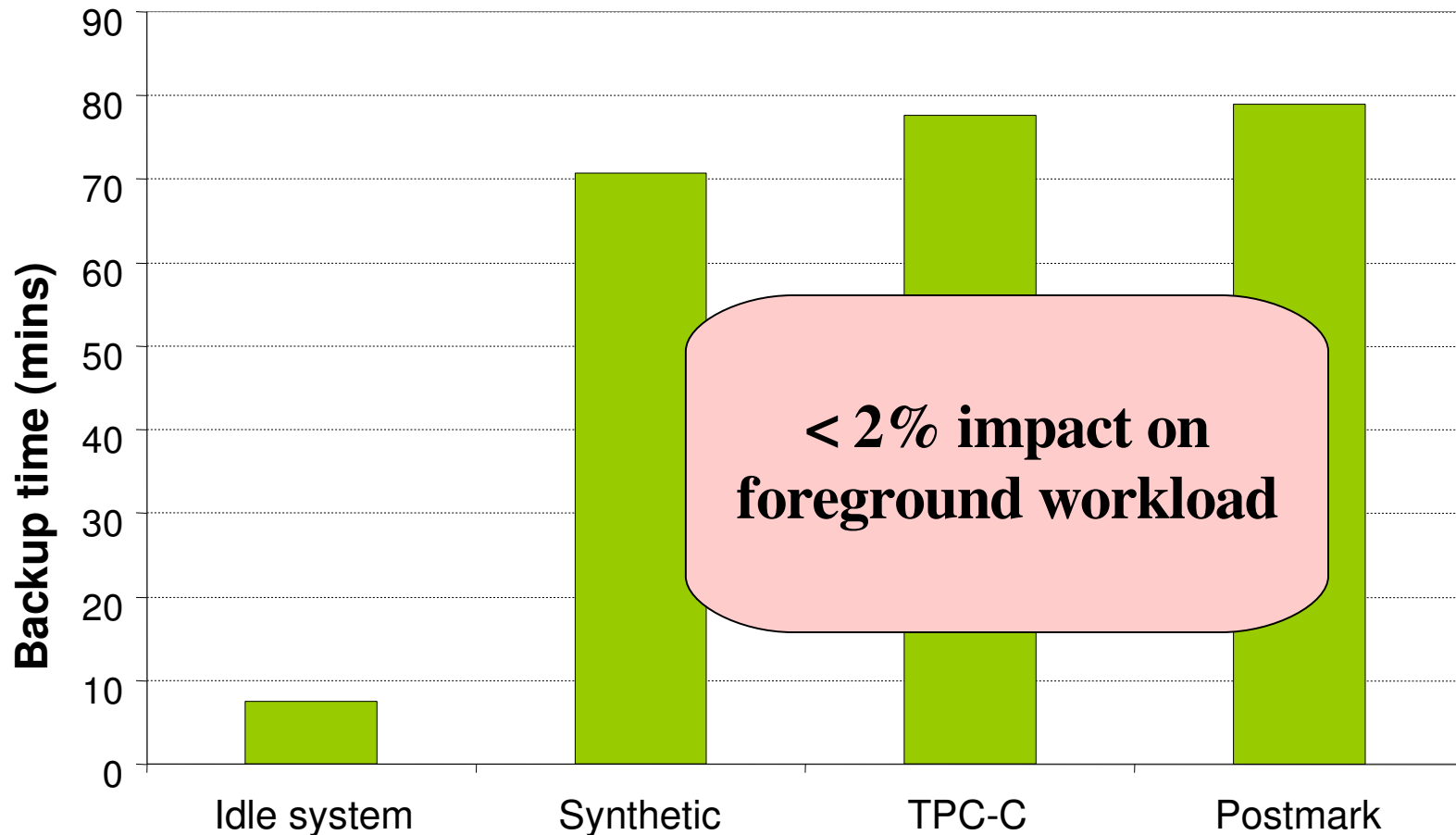
- Next few slides will show that:
  - we can build background apps
    - that do not interfere with foreground apps
    - that complete eventually
    - things like *backup*, *defrag*, *virus scanners*, etc
    - imagine the possibilities...

# App 1: Backup

---

- Frequent backup improves data reliability and availability
  - companies take very frequent backups
  - a backup every 30 mins is not uncommon
- Our experiment:
  - disk used is 18GB
  - we want to back up 12GB of data
  - goal: back it up for free

# Backup completed for free

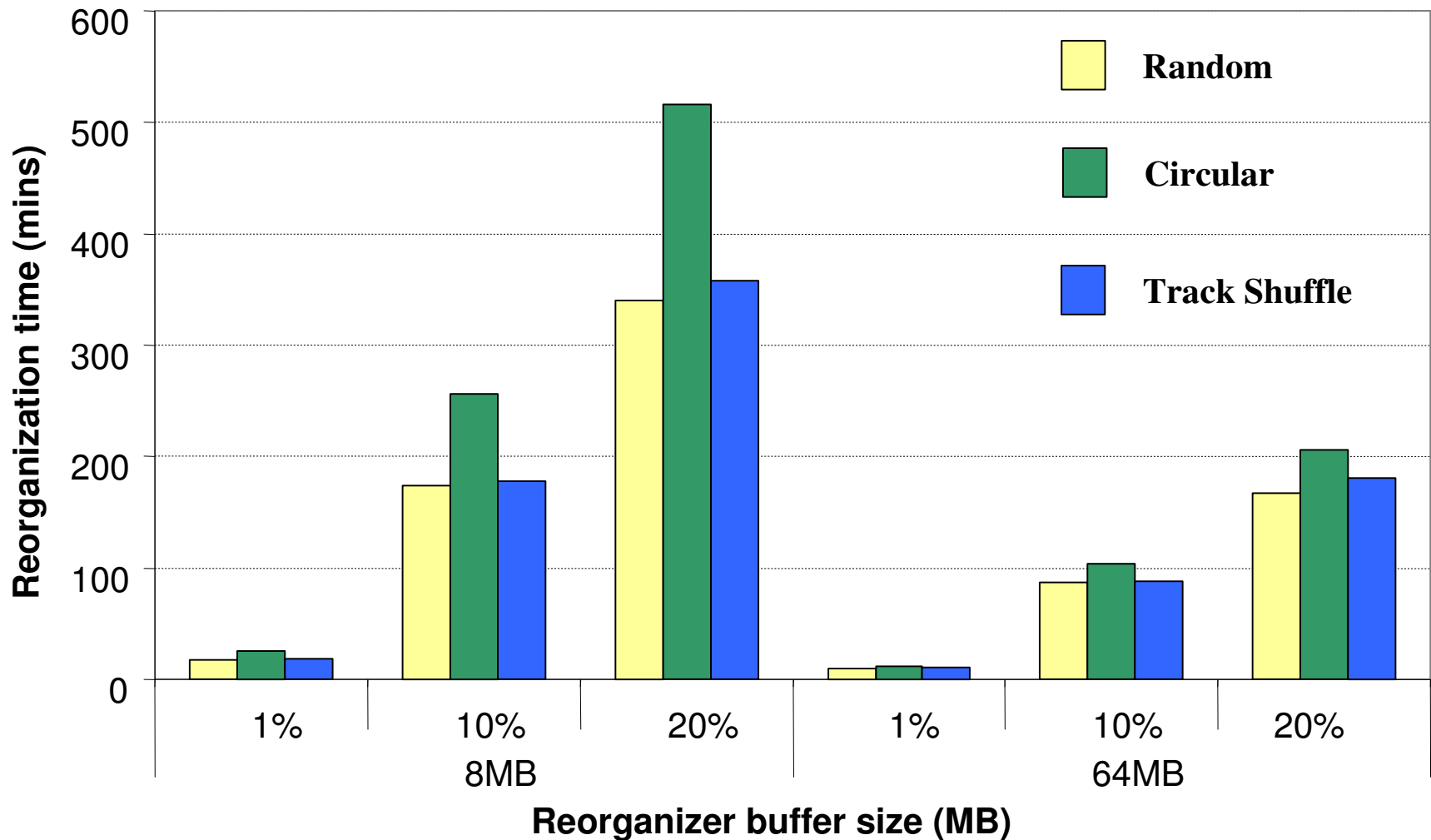


## App 2: Layout reorganization

---

- Layout reorganization improves access latencies
  - defragmentation is a type of reorganization
  - typical example of background activity
- Our experiment:
  - disk used is 18GB
  - we want to defrag up to 20% of it
  - goal: defrag for free

# Disk Layout Reorganized for Free!



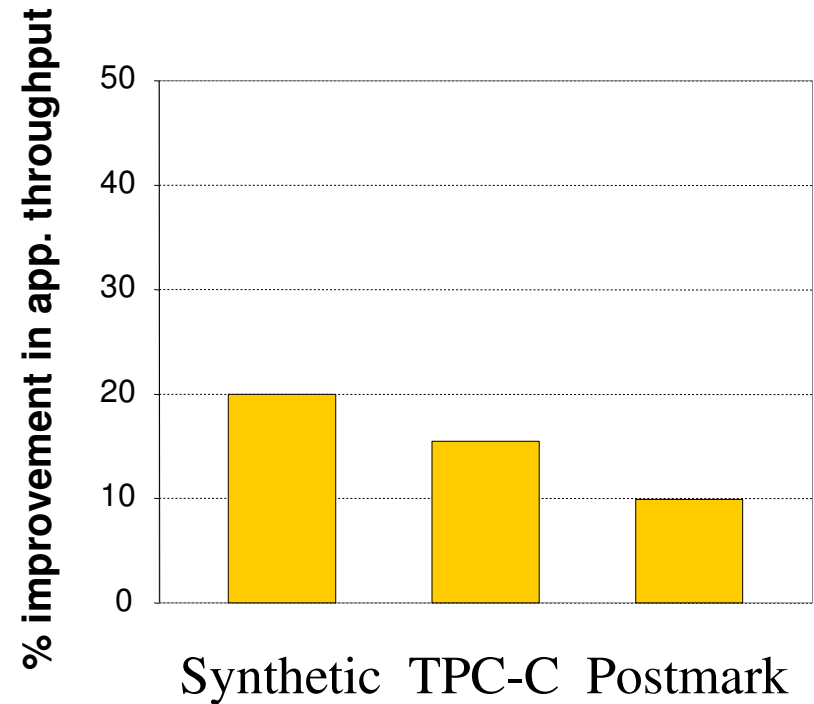
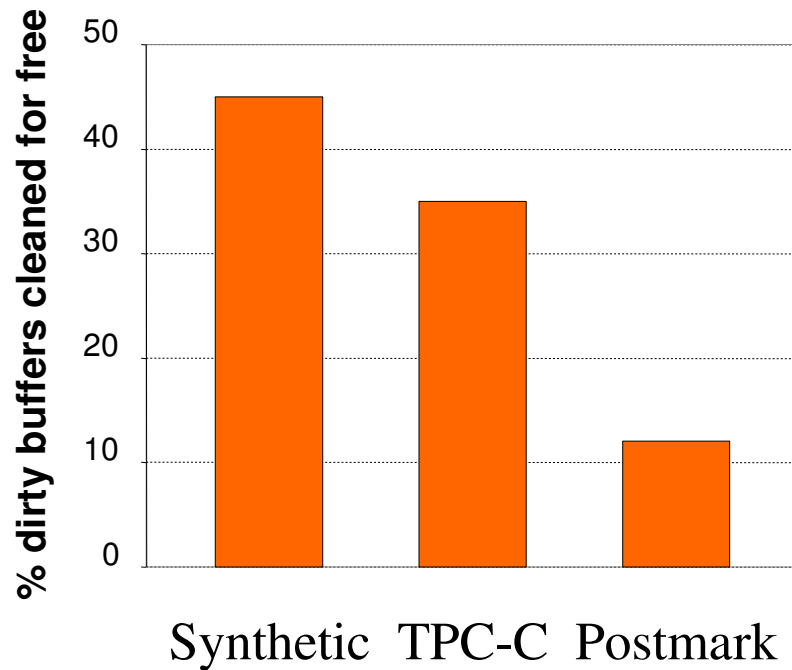
# App 3: Cache write-backs

---

- Must flush dirty buffers
  - for space reclamation
  - for persistence (if memory is not NVRAM)
- Our experiment
  - PIII with 384MB of RAM
  - controlled experiments with synthetic workload
  - benchmarks (same as used before) in FreeBSD
  - syncer daemon wakes up every 1 sec and flushes entries that have been dirty > 30secs
  - goal: write back dirty buffers for free

# 10-20% improvement in overall perf.

---



# Other maintenance applications

---

- Virus scanner
- LFS cleaner
- Disk scrubber
- Data mining
- Data migration

# Summary I

---

- Disks are slow
  - but we can squeeze extra bw out of them
- Use freeblock scheduling to extract free bandwidth
- Utilize free bandwidth for background applications
  - they still complete eventually
  - with no impact on foreground workload

# Implementation considerations: preview

---

- Next few slides will show that:
  - it's hard to do fine grained scheduling at the device driver
  - background apps need new interfaces to express their desires to the background scheduler
  - file consistency issues

# Implementing freeblock scheduling

---

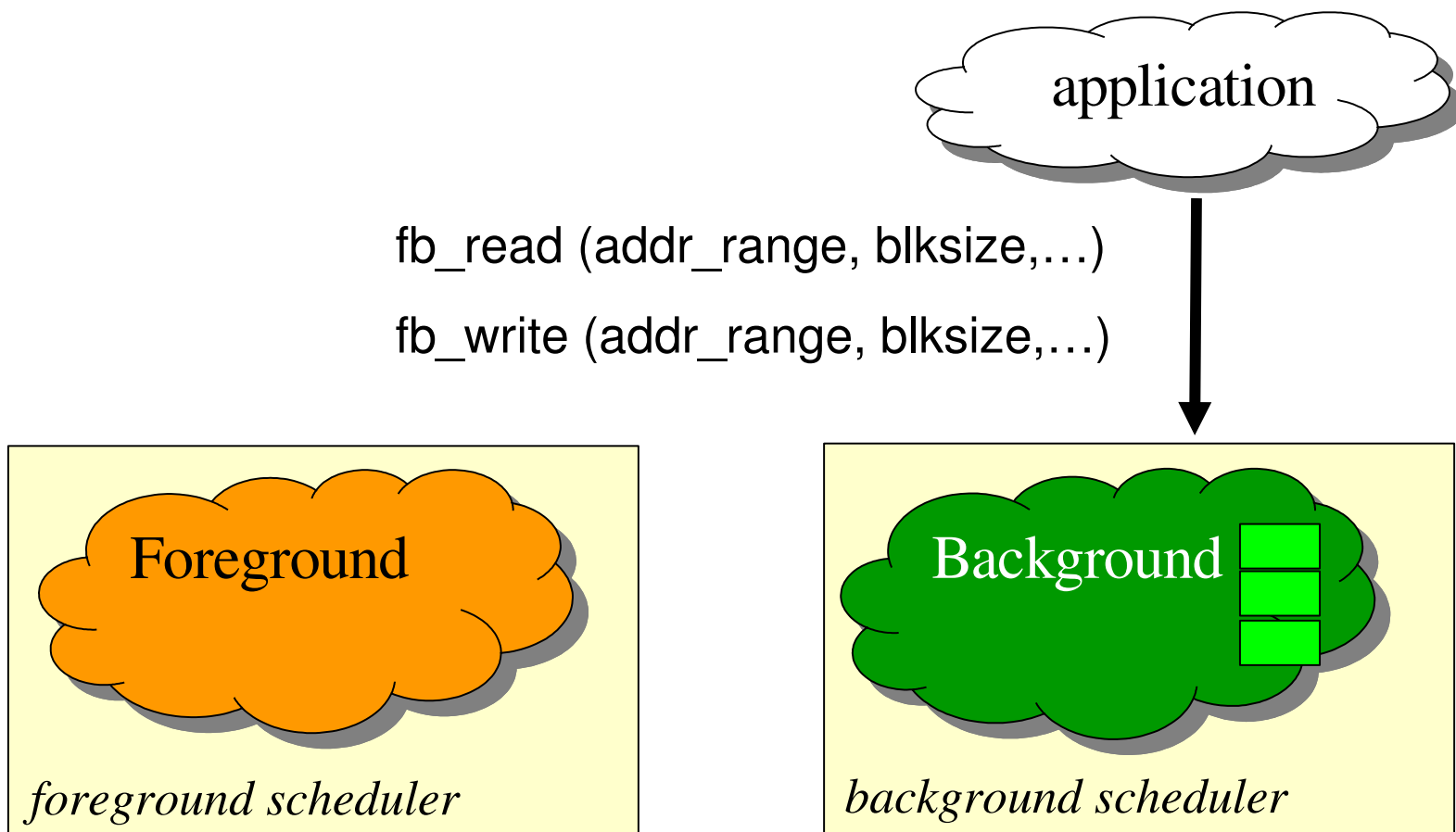
- Hard to do at the device driver
  - need to know the position of the disk head
  - however, we have done it!
  - it's more efficient inside the disk drive
    - try to convince your disk vendor to put it in
- Efficient algorithms
  - SPTF for foreground (0.5% of 1GHz PIII)
  - freeblock scheduling for background (<<5% of 1GHz PIII)
  - small memory footprint

# Application programming interface (API) goals

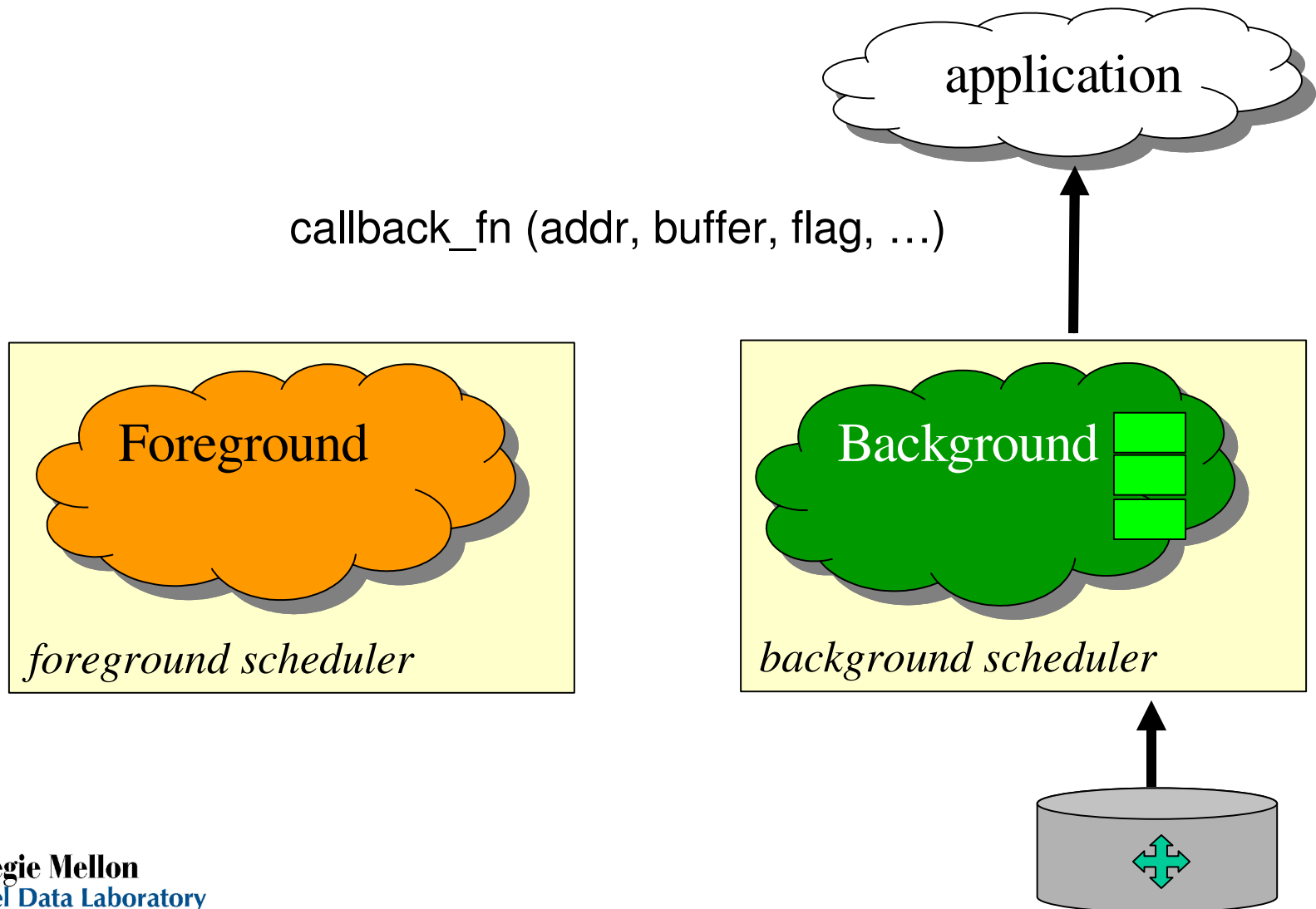
---

- Work exposed but done opportunistically
  - all disk accesses are asynchronous
- Minimized memory-induced constraints
  - late binding of memory buffers
  - late locking of memory buffers
- “Block size” can be application-specific
- Support for speculative tasks
- Support for rate control

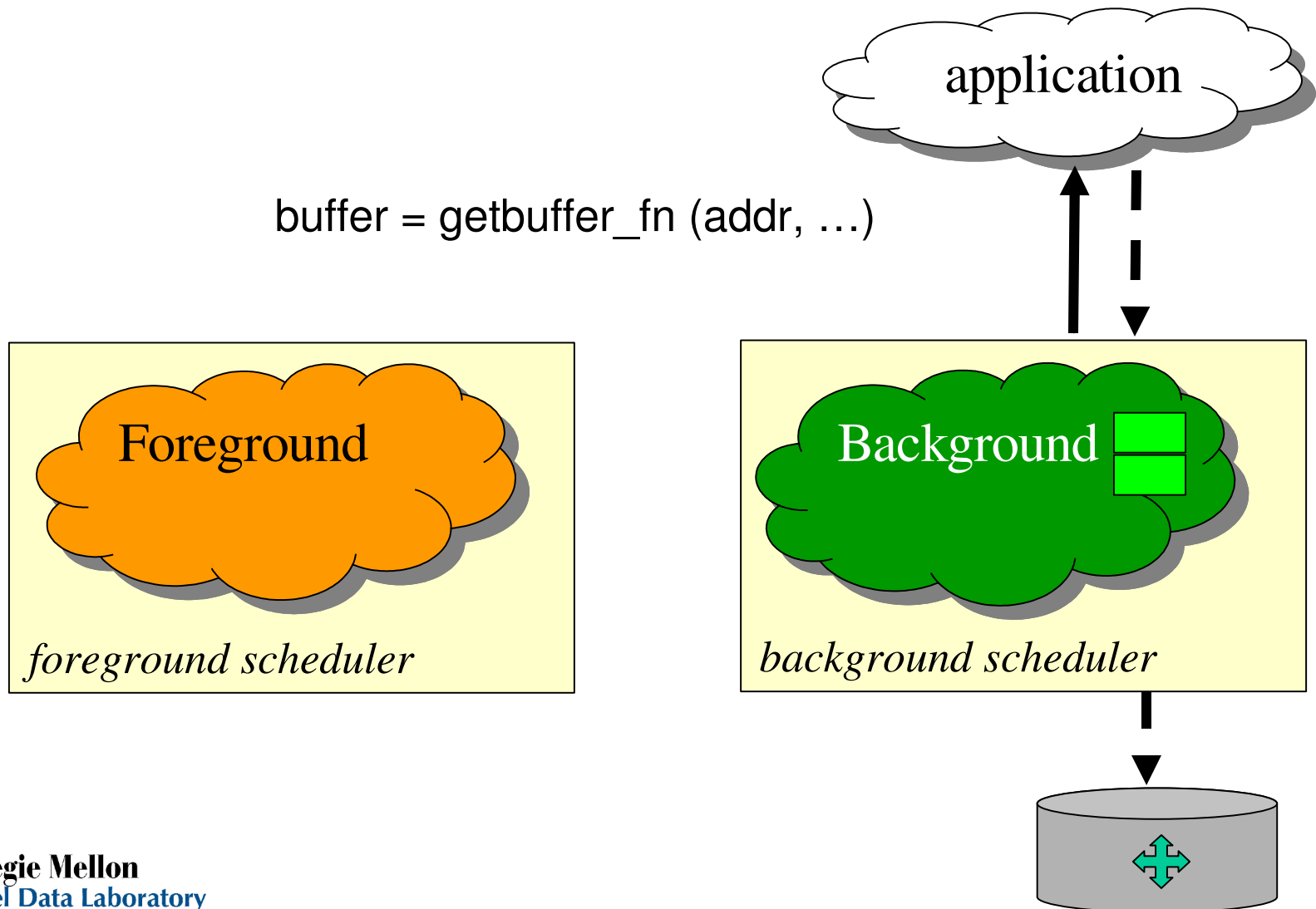
# API description: task registration



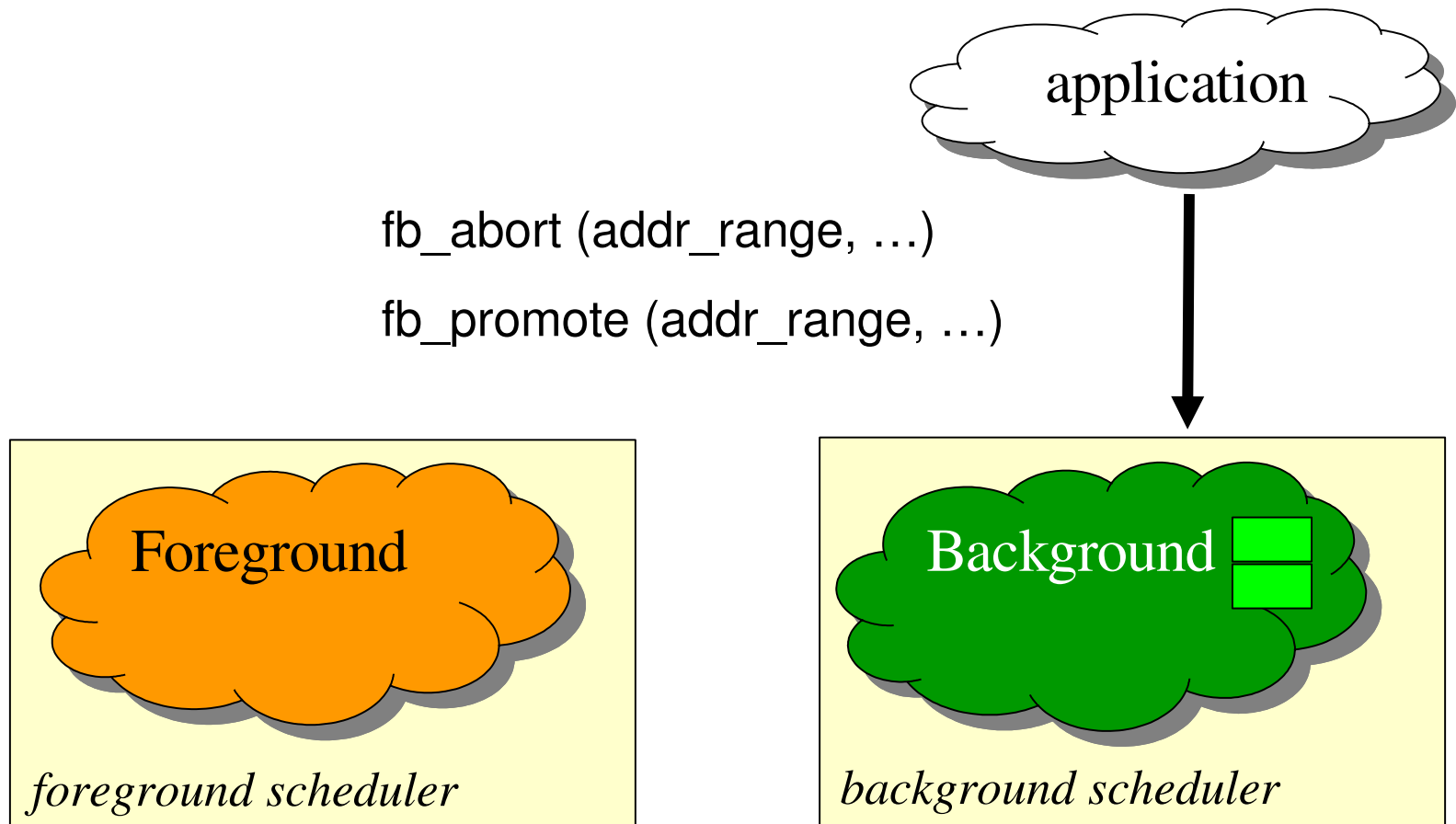
# API description: task completion



# API description: late locking of buffers



# API description: aborting/promoting tasks



# Complete API

---

| Function Name          | Arguments  | Description  |
|------------------------|--|--|
| <i>fb_open</i>         | <i>priority, callback_fn, getbuffer_fn</i>             | Open a freeblock session (ret: <i>session_id</i> )     |
| <i>fb_close</i>        | <i>session_id</i>                                      | Close a freeblock session                              |
| <i>fb_read</i>         | <i>session_id, addr_range, blksize, callback_param</i> | Register a freeblock read task                         |
| <i>fb_write</i>        | <i>session_id, addr_range, blksize, callback_param</i> | Register a freeblock write task (ret: <i>task_id</i> ) |
| <i>fb_abort</i>        | <i>session_id, addr_range</i>                          | Abort parts of registered task                         |
| <i>fb_promote</i>      | <i>session_id, addr_range</i>                          | Promote parts of registered task                       |
| <i>fb_suspend</i>      | <i>session_id</i>                                      | Suspend scheduling of a session's tasks                |
| <i>fb_resume</i>       | <i>session_id</i>                                      | Resume scheduling of a session's tasks                 |
| <i>*(callback_fn)</i>  | <i>session_id, addr, buffer, flags, callback_param</i> | Report that part of task completed                     |
| <i>*(getbuffer_fn)</i> | <i>session_id, addr, callback_param</i>                | Get memory address for selected write                  |

# Designing disk maintenance applications

---

- APIs talk in terms of logical blocks (LBNs)
- Some applications need structured version
  - as presented by file system or database
- Example consistency issues
  - application wants to read file “foo”
  - registers task for inode’s blocks
  - by time blocks read, file may not exist anymore!

# Designing disk maintenance applications

---

- Application does not care about structure
  - scrubbing, data migration, array reconstruction
- Coordinate with file system/database
  - cache write-backs, LFS cleaner, index generation
- Utilize snapshots
  - backup, background *fsck*

# Summary II

---

- Utilize free bandwidth for background applications
  - they still complete eventually
  - with no impact on foreground workload
- Scheduling is fine-grained
- Need new APIs

# Q & A

---

- See <http://www.pdl.cmu.edu/Freeblock/>
- Windows Vista!!!
- Talk to me if interested in summer internships or research with the PDL