# [537] Fast File System

Chapter 41
Tyler Harter
11/10/14

# File-System Case Studies

Local
- **FFS**: Fast File System
- **LFS**: Log-Structured File System

Network
- **NFS**: Network File System
- **AFS**: Andrew File System

# File-System Case Studies

Local
  - **FFS**: Fast File System [today]
  - **LFS**: Log-Structured File System


Network
  - **NFS**: Network File System
  - **AFS**: Andrew File System

# Review Basic FS

# Basic FS

Structures (on disk)

Operations

# Structure Overview

Core                                Performance

**Super Block**

# Structure Overview

Core                          Performance

**Super Block**

**Data Block**

# Structure Overview

Core                                    Performance

**Super Block**

**Data Block**

**Inode Table**

# Structure Overview

Core

Performance

**Super Block**

**Data Block**

directories

indirects

**Inode Table**

# Structure Overview

## Core

**Super Block**

**Data Block**
directories | indirects

**Inode Table**

## Performance

**Data Bitmap**

**Inode Bitmap**

# Layout

# Layout



super block | bit maps | inodes | Data Blocks

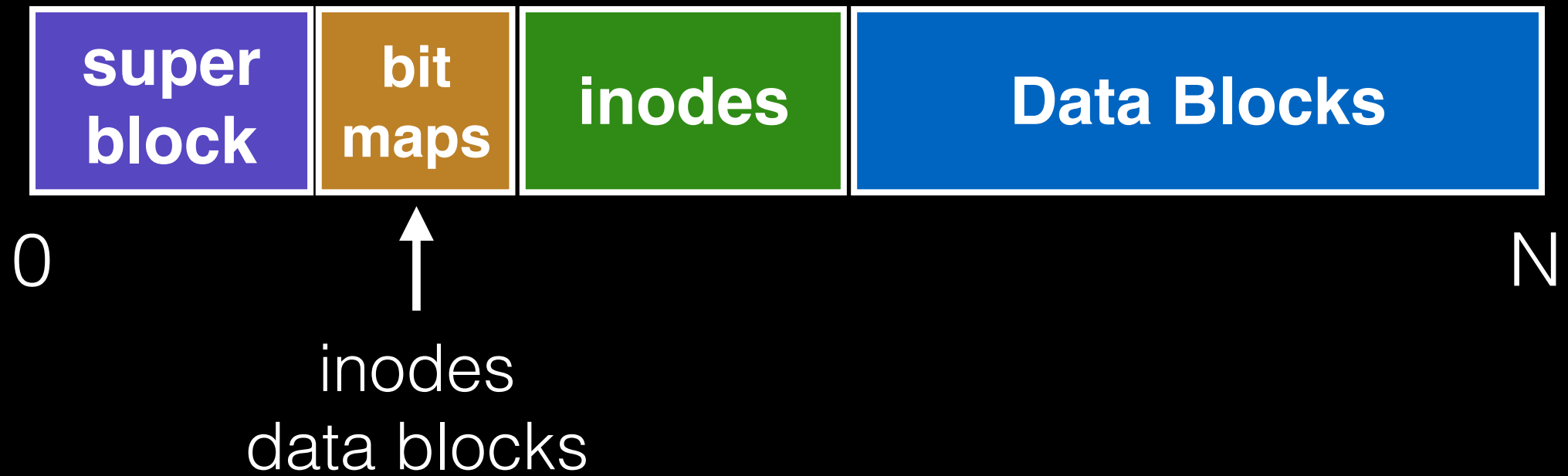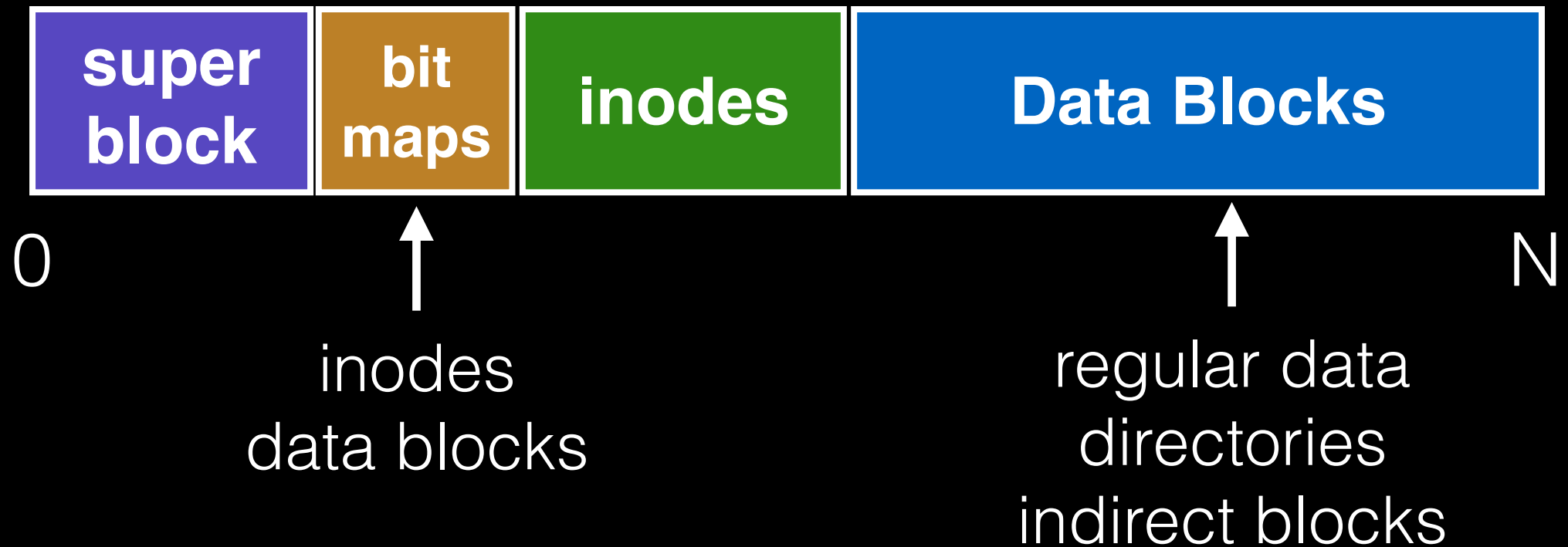0                                                    N

# Layout

# Layout

# Basic FS

Structures (on disk)

Operations

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
| --- | --- | --- | --- | --- | --- | --- |
| | | | | | | |

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | | | read | |

## create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | read | | | |
| | | | | | read | |
| | | | | | | read |

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | read | | | |
| | | | | | read | |
| | | | | | | read |

bar does not already exist

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | read | | | |
| | | | | | read | |
| | | | | | | read |

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | read | | | |
| | | | | | read | |
| | | | | | | read |
| | read write | | | | | |

# create /foo/bar   <span style="color:red">[populate inode]</span>

| data<br>bitmap | inode<br>bitmap | root<br>inode | foo<br>inode | bar<br>inode | root<br>data | foo<br>data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | | | read | |
| | | | read | | | |
| | | | | | | read |
| | read<br>write | | | | | |
| | | | | read<br>write | | |

# create /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data |
|---|---|---|---|---|---|---|
| | | read | | | | |
| | | | read | | | |
| | | | | | read | |
| | | | | | | read |
| | read write | | | | | |
| | | | | read write | | |
| | | | write | | | |
| | | | | | | write |

# append to /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data |
| --- | --- | --- | --- | --- | --- | --- | --- |

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data |
|---|---|---|---|---|---|---|---|
| | | | | read | | | |

# append to /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data |
|---|---|---|---|---|---|---|---|
|  |  |  |  | read |  |  |  |
| read |  |  |  |  |  |  |  |
| write |  |  |  |  |  |  |  |

append to /foo/bar

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data |
|---|---|---|---|---|---|---|---|
| | | | | read | | | |
| read | | | | | | | |
| write | | | | | | | |
| | | | | write | | | |

append to /foo/bar     <span style="color:red">[write to block]</span>

| data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data |
|---|---|---|---|---|---|---|---|
| read | | | | read | | | |
| write | | | | | | | |
| | | | | write | | | |
| | | | | | | | write |

# Review Locality

# Locality Types



address

time

address

time

# Locality Types



**Temporal Locality**

...

address

time

**Spatial Locality**

...

address

time

# Locality Usefulness

What types of locality are useful for a cache?

What types of locality are useful for a disk?

# Order Matters Now

# Order Matters Now

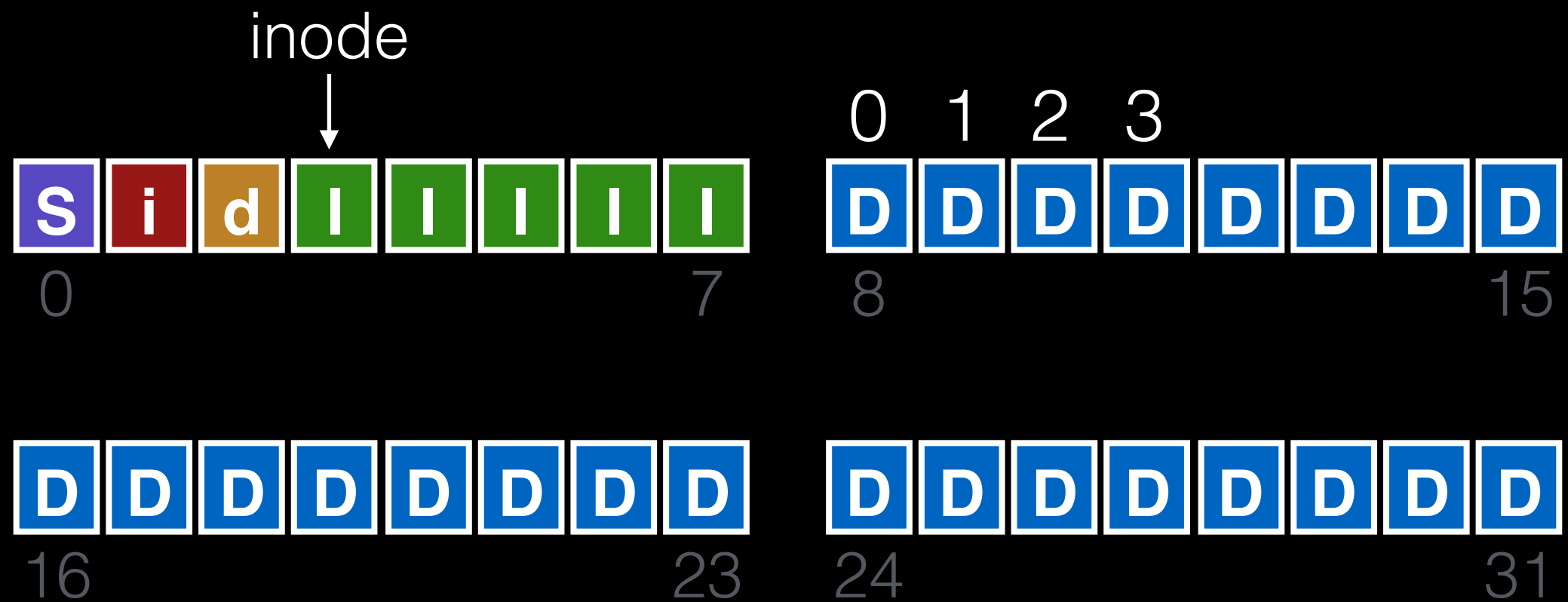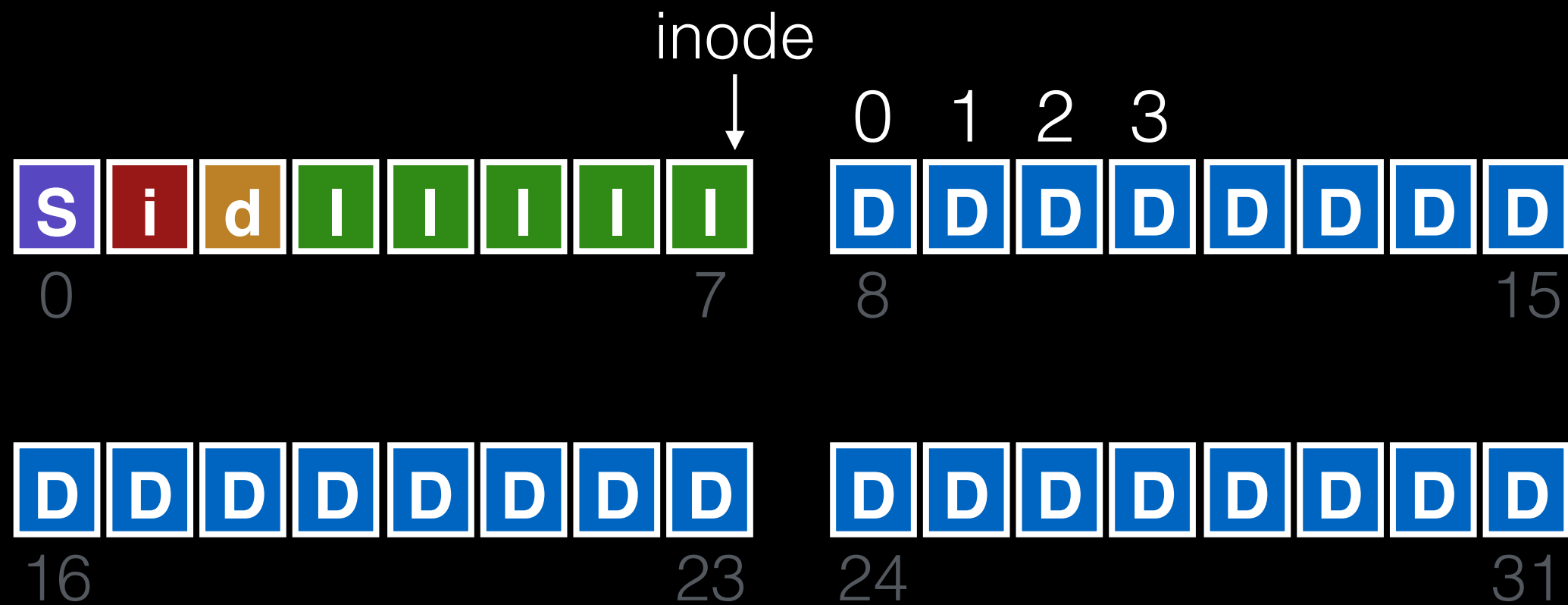# Policy: Choose Inode, Data Blocks

# Bad File Layout

inode

**S** **i** **d** **l** **l** **l** **l** **l**

0                            7

0

**D** **D** **D** **D** **D** **D** **D** **D**

8                           15

3                             2                           1

**D** **D** **D** **D** **D** **D** **D** **D**     **D** **D** **D** **D** **D** **D** **D** **D**

16                         23    24                         31

# Better File Layout

# Best File Layout

# Fast File System

# System Building

**noob approach**
1. get idea
2. build it!

# System Building

**noob approach**
1. get idea
2. build it!

**pro approach**
1. identify state of the art
2. measure it, identify problems
3. get idea
4. build it!

# System Building

**noob approach**
1. get idea
2. build it!

**pro approach**
1. identify state of the art
2. measure it, identify problems
3. get idea
4. build it!                    *measure then build*

# Old FS

State of the art: original UNIX file system.

# Layout



Free lists are embedded in inodes, data blocks.
Data blocks are 512 bytes.

# Old FS

State of the art: original UNIX file system.

# Old FS

State of the art: original UNIX file system.

Measure throughput for file reads/writes.

Compare to theoretical max, which is…

# Old FS

State of the art: original UNIX file system.

Measure throughput for file reads/writes.

Compare to theoretical max, which is…
disk bandwidth

# Old FS

State of the art: original UNIX file system.

Measure throughput for file reads/writes.

Compare to theoretical max, which is…
disk bandwidth

Old UNIX file system: only **2%** of potential.  Why?

# Measurement 1

What is performance before/after aging?

# Measurement 1

What is performance before/after aging?

New FS: **17.5%** of disk bandwidth
Few weeks old: **3%** of disk bandwidth

# Measurement 1

What is performance before/after aging?

New FS: **17.5%** of disk bandwidth
Few weeks old: **3%** of disk bandwidth

FS is probably becoming fragmented over time.

Free list makes contiguous chunks hard to find.

# Measurement 1

What is performance before/after aging?

New FS: **17.5%** of disk bandwidth

Few weeks old: **3%** of disk bandwidth

<span style="color:red">hacky solution: occasional defrag</span>

FS is probably becoming fragmented over time.

Free list makes contiguous chunks hard to find.

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

Performance **more** than doubled.

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

Performance **more** than <span style="color:red">doubled</span>.

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

Performance **more** than <span style="color:red">doubled</span>.

<span style="color:orange">Logically adjacent</span> blocks are probably not
<span style="color:blue">physically adjacent</span>.

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

Performance **more** than doubled.

Logically adjacent blocks are probably not physically adjacent.

# Measurement 2

How does <u>block size</u> affect performance?
Try doubling it!

Performance **more** than doubled.

Logically adjacent blocks are probably not
physically adjacent.

Smaller blocks cause more indirect I/O.

# Old FS Summary

Observations:
 - long distance between inodes/data
 - inodes in single dir not close to one another
 - small blocks (512 bytes)
 - blocks laid out poorly
 - free list becomes scrambled, causes random alloc

Result: **2%** of potential performance!
(and worse over time)

# Problem: old FS treats disk like RAM!

# Solution: a disk-aware FS

# Design Questions

How to use big blocks without wasting space.

How to place data on disk.

# Technique 1: Bitmaps

# Technique 1: Bitmaps

| super block | bitmaps | inodes | Data Blocks |
|:---:|:---:|:---:|:---:|

0                                                                    N

Use bitmaps instead of free list.
Provides more flexibility, with more global view.

# Techniques

Bitmaps

# Technique 2: Groups



fast

| super block | bitmaps | inodes | Data Blocks |

0                                                    N

before: whole disk

# Technique 2: Groups

# Technique 2: Groups

# Technique 2: Groups

slowest

| super block | bitmaps | inodes | Data Blocks |

0                                        N

before: whole disk

# Technique 2: Groups



super block | bitmaps | inodes | Data Blocks

0      N

before: whole disk

# Technique 2: Groups



| super block | bitmaps | inodes | Data Blocks |
|---|---|---|---|

0                                      G

now: one (smallish) group

# Technique 2: Groups



S B I D S B I D S B I D ...

0 · · · group 1 · · · G · · · group 2 · · · 2G · · · group 3 · · · 3G

zoom out

# Technique 2: Groups



strategy: allocate inodes and data blocks in same group.

# Groups

In FFS, groups were ranges of cylinders
 - called <u>cylinder group</u>


In ext2-4, groups are ranges of blocks
 - called <u>block group</u>

# Groups

In FFS, groups were ranges of cylinders
 - called <u>cylinder group</u>

In ext2-4, groups are ranges of blocks
 - called <u>block group</u>

# Techniques

Bitmaps
Locality groups

# Technique 3: Super Rotation

# Technique 3: Super Rotation



Is it useful to have multiple super blocks?

# Technique 3: Super Rotation



Is it useful to have multiple super blocks?
Yes, if some (but not all) fail.

# Problem

# Problem



All super-block copies are on the top platter. What if it dies?

# Problem

All super-block copies are on the top platter. What if it dies?

solution: for each group, store super-block at different offset.

# Techniques

Bitmaps
Locality groups
Rotated super

# Block Size

Doubling the block size for the old FS over doubled performance.

Strategy: choose block size so we never have to read more than two indirect blocks to find a data block (2 levels of indirection max).  Want 4GB files.

How large is this?

# Techniques

Bitmaps
Locality groups
Rotated super
Large blocks

# Large Blocks

Why not make blocks huge?

Most file are very small.

# Large Blocks

Why not make blocks huge?
Lots of waste in remainder of blocks.

# Large Blocks

Why not make blocks huge?
Lots of waste in remainder of blocks.

# Large Blocks

Why not make blocks huge?
Lots of waste in remainder of blocks.

Time vs. Space
Tradeoffs…

# Solution: Fragments

Hybrid!

Introduce "fragment" for files that use parts of blocks.

Only tail of file uses fragments.

# Fragment Example

Block size = 4096
Fragment size = 1024

bits: 0000   0000 1111 0010
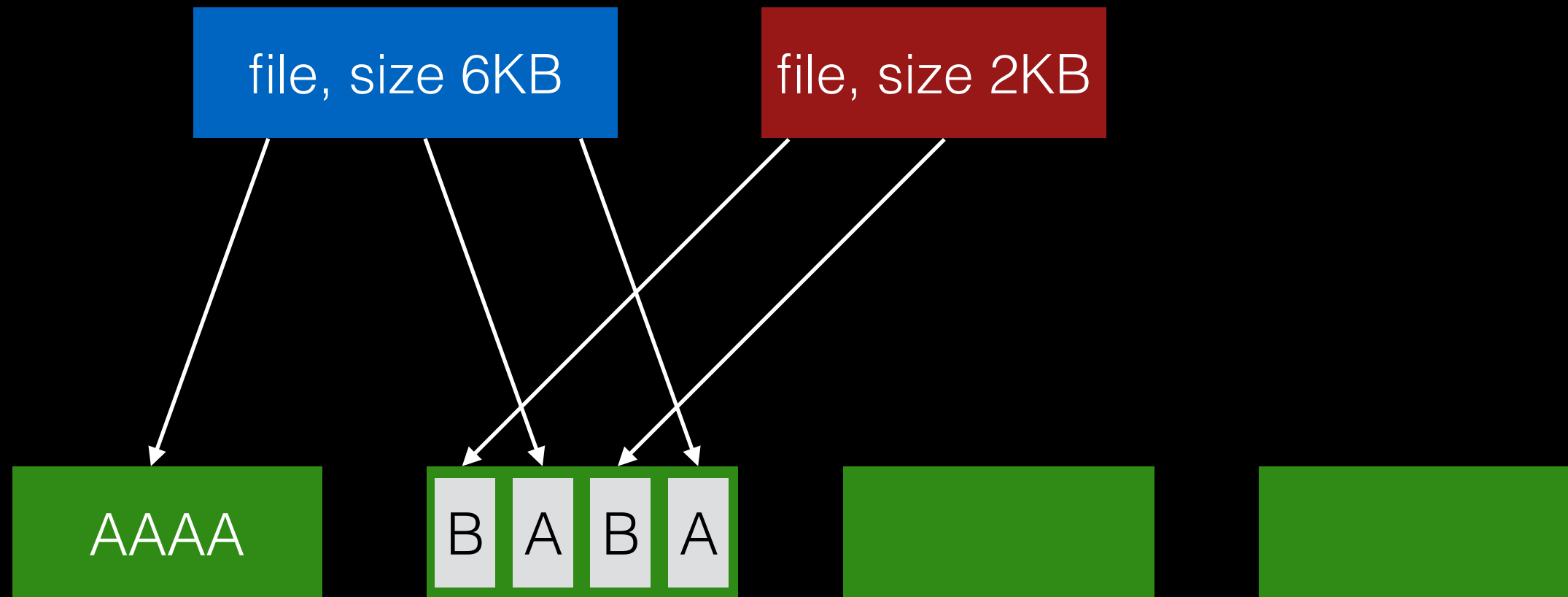      blk1   blk2  blk3  blk4

# How to Decide

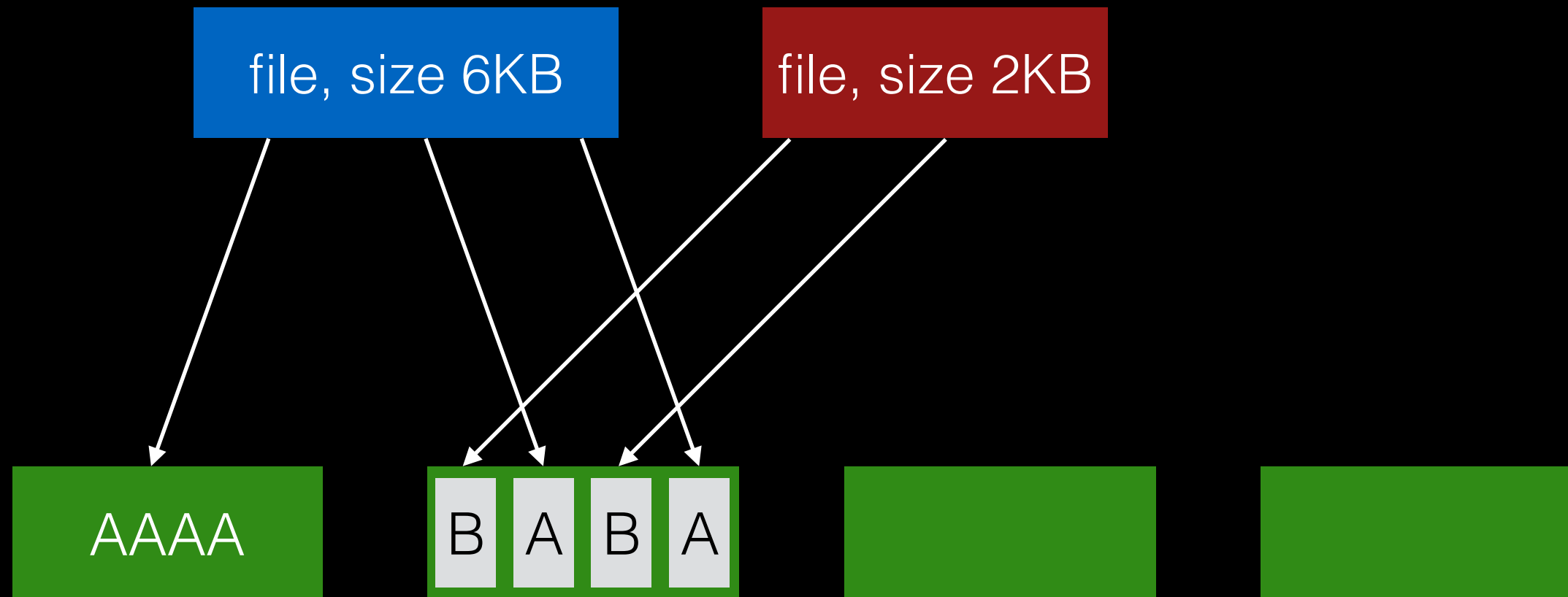Whether addr refers to block or fragment is inferred by the file size.
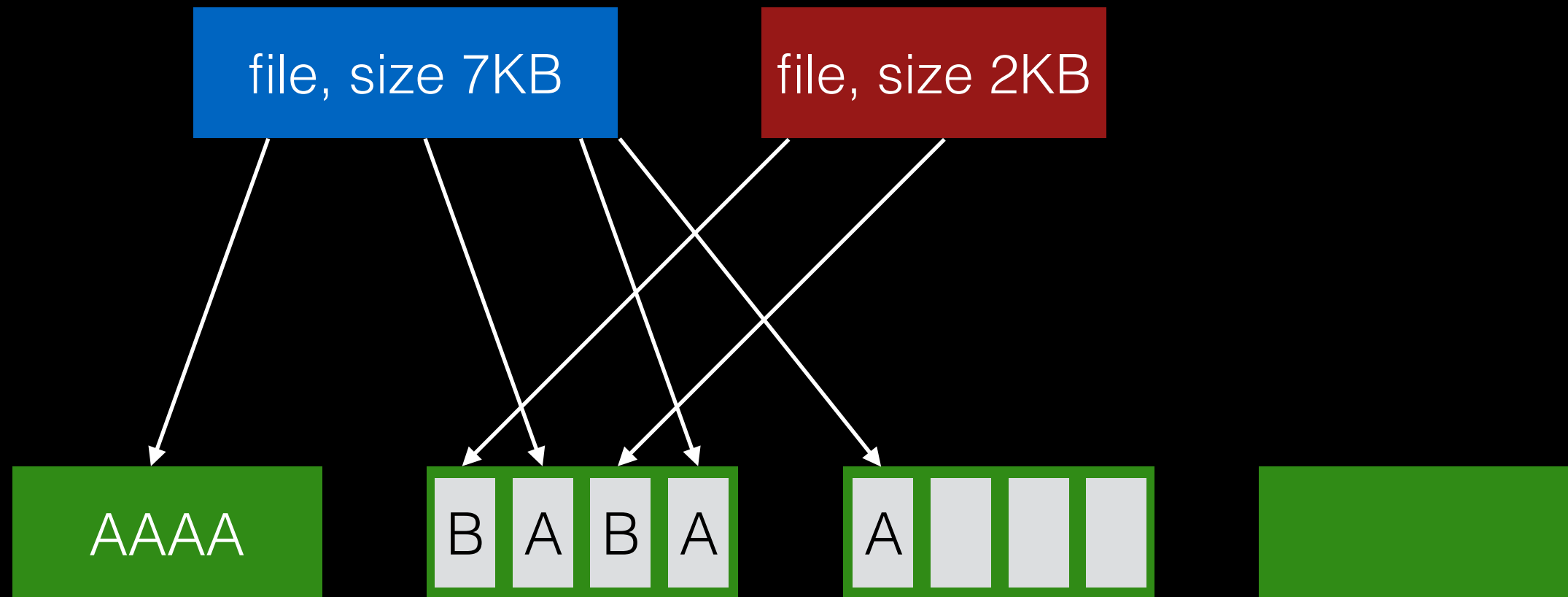
What about when files grow?

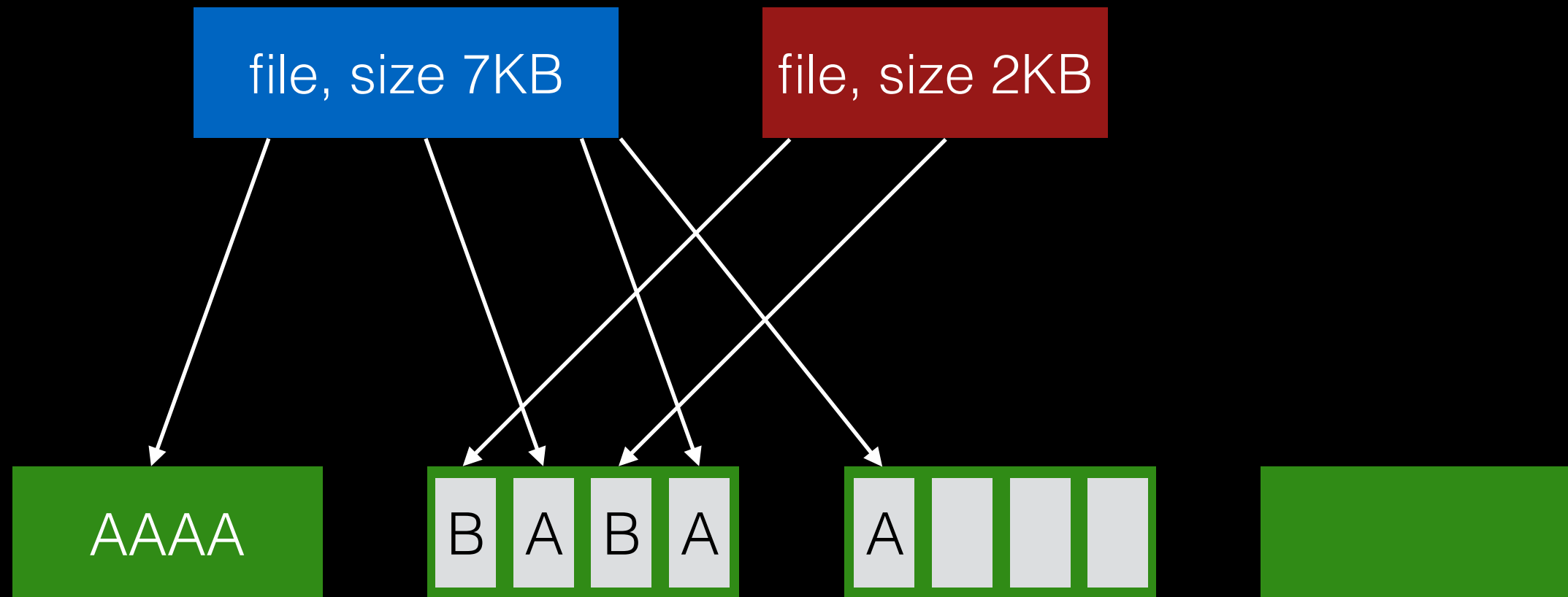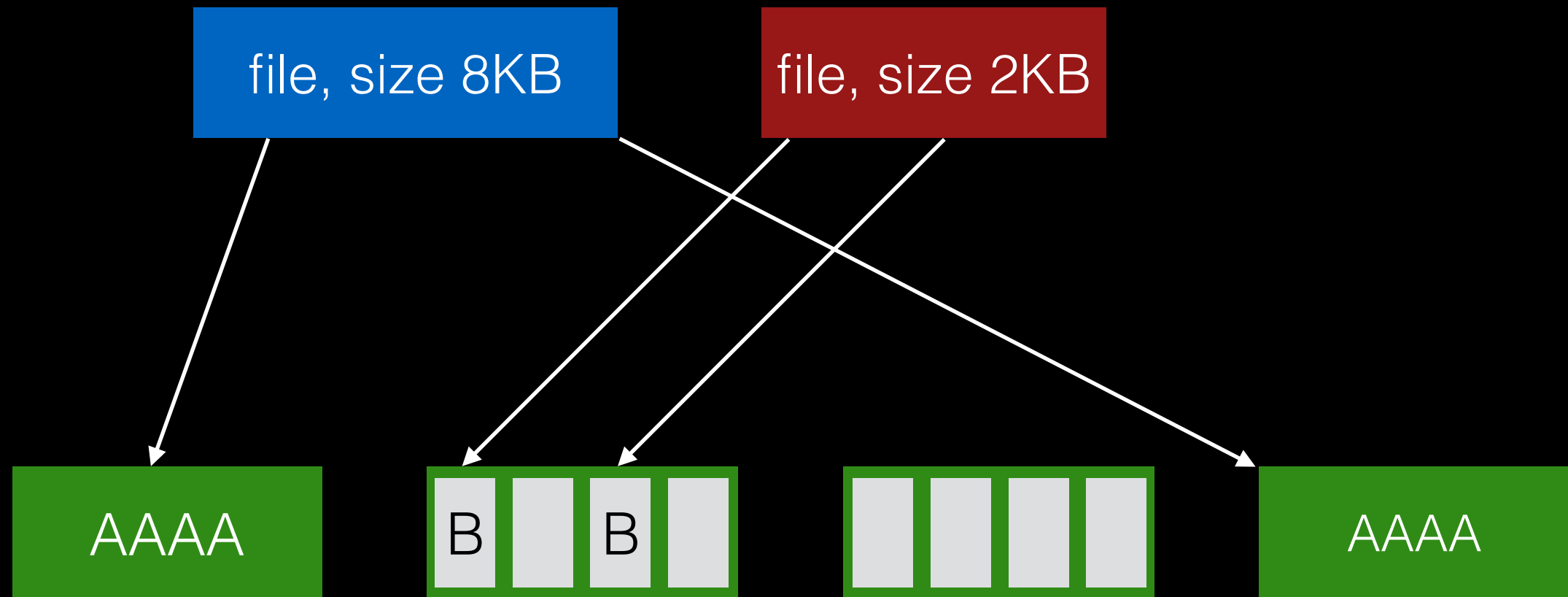Must copy fragments to new block if there's not room to grow.

file, size 6KB    file, size 2KB

AAAA    B A B A

append A to first file

file, size 6KB    file, size 2KB

AAAA    B A B A

append A to first file, copy to fragments to new block.

# Optimal Write Size

Writing less than a block is inefficient.

Solution: new API exposes optimal write size.

For pipes and sockets, the new call returns the buffer size.

The stdio library uses this call.

# Techniques

Bitmaps
Locality groups
Rotated super
Large blocks
Fragments

# Smart Policy



Where should new inodes and data blocks go?

# Strategy

Put related pieces of data near each other.

# Strategy

Put related pieces of data near each other.

Rules:
1. Put directory entries near directory inodes.
2. Put inodes near directory entries.
3. Put data blocks near inodes.

# Strategy

Put related pieces of data near each other.

Rules:
1. Put directory entries near directory inodes.
2. Put inodes near directory entries.
3. Put data blocks near inodes.

Sound good?

# Challenge

The file system is one big tree.

All directories and files have a common root.

In some sense, all data in the same FS is related.

# Challenge

The file system is one big tree.

All directories and files have a common root.

In some sense, all data in the same FS is related.

Trying to put everything near everything else will leave us with the same mess we started with.

# Revised Strategy

Put more-related pieces of data near each other.

Put less-related pieces of data far from each other.

# Revised Strategy

Put more-related pieces of data near each other.
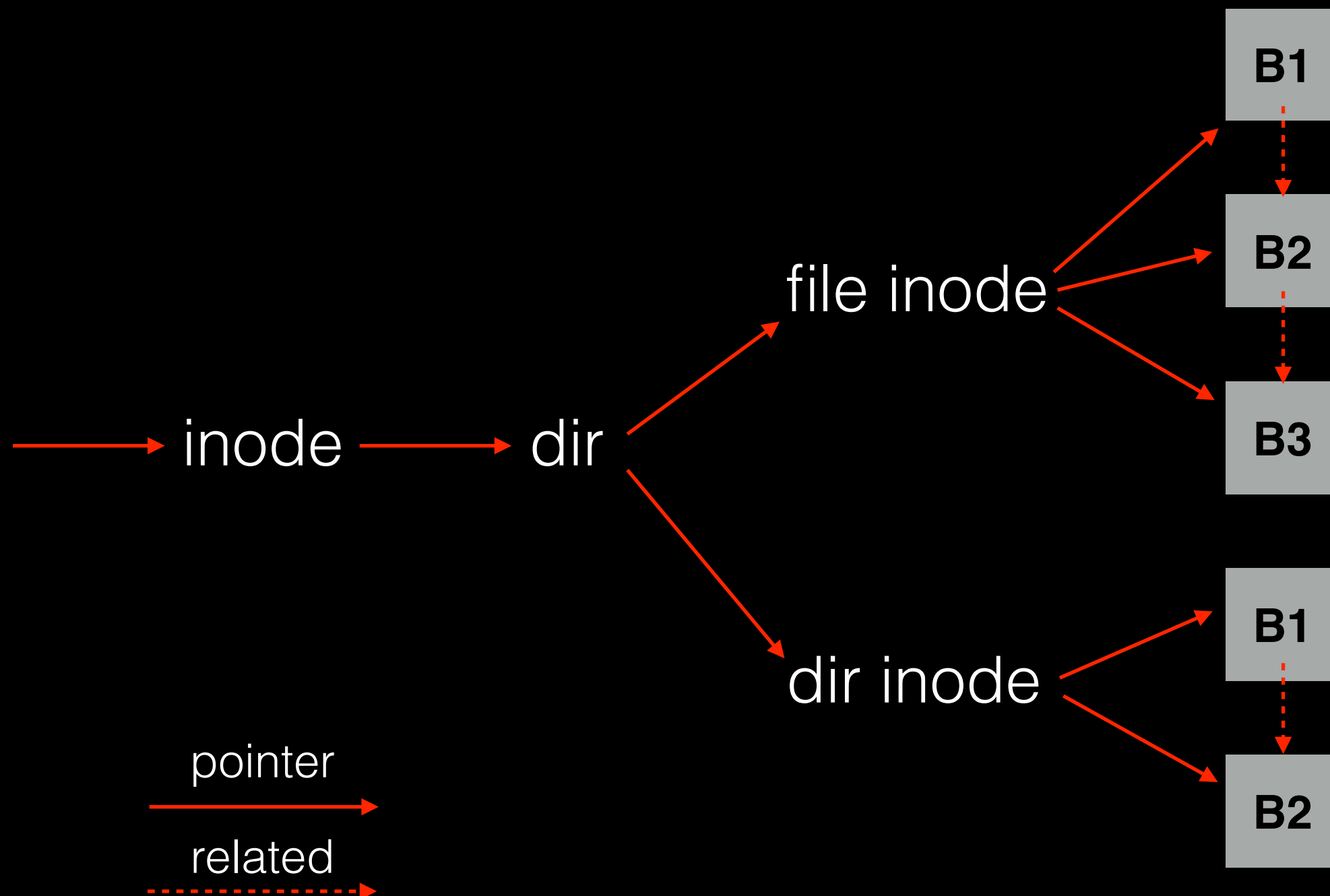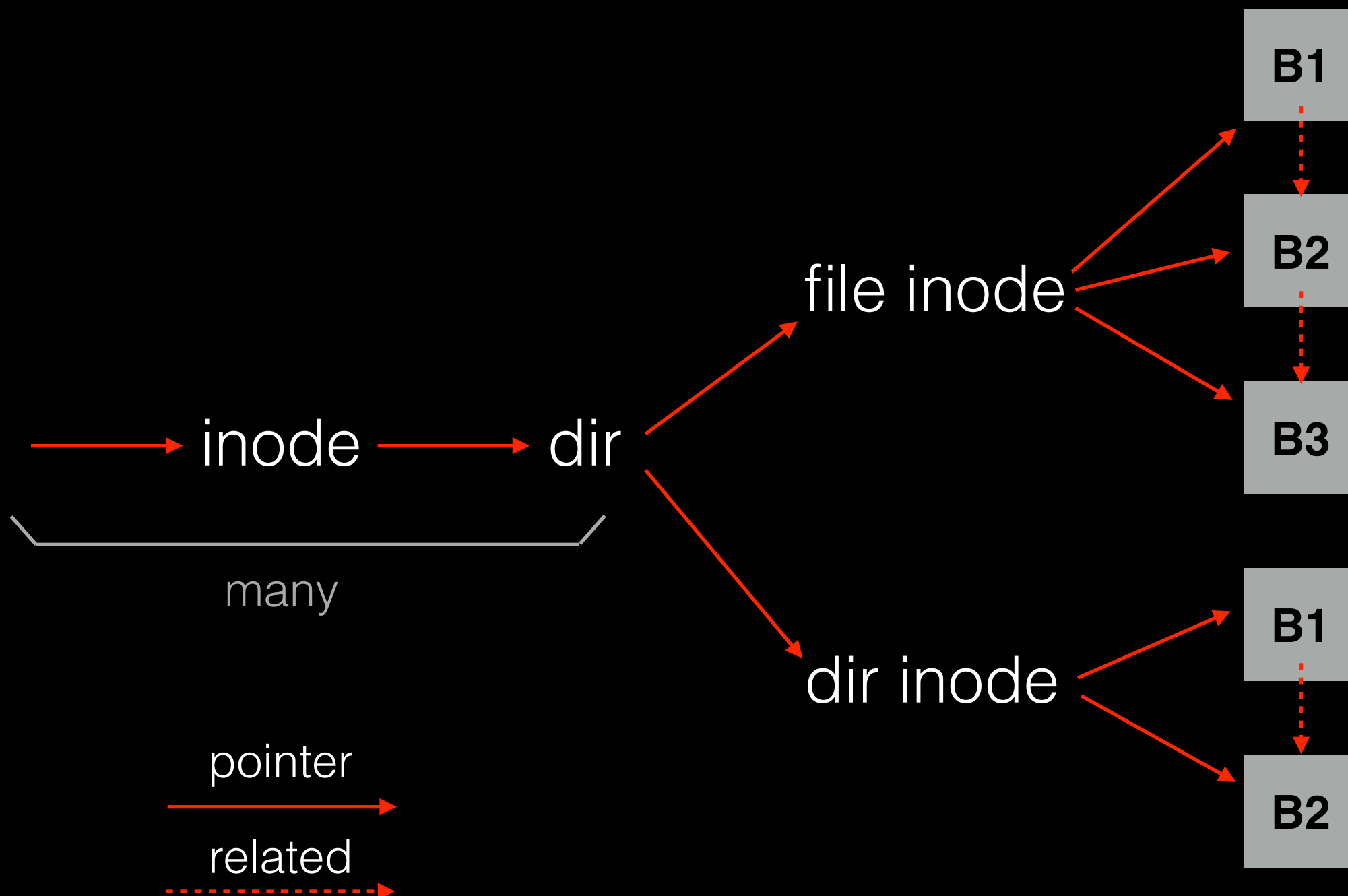
Put less-related pieces of data far from each other.

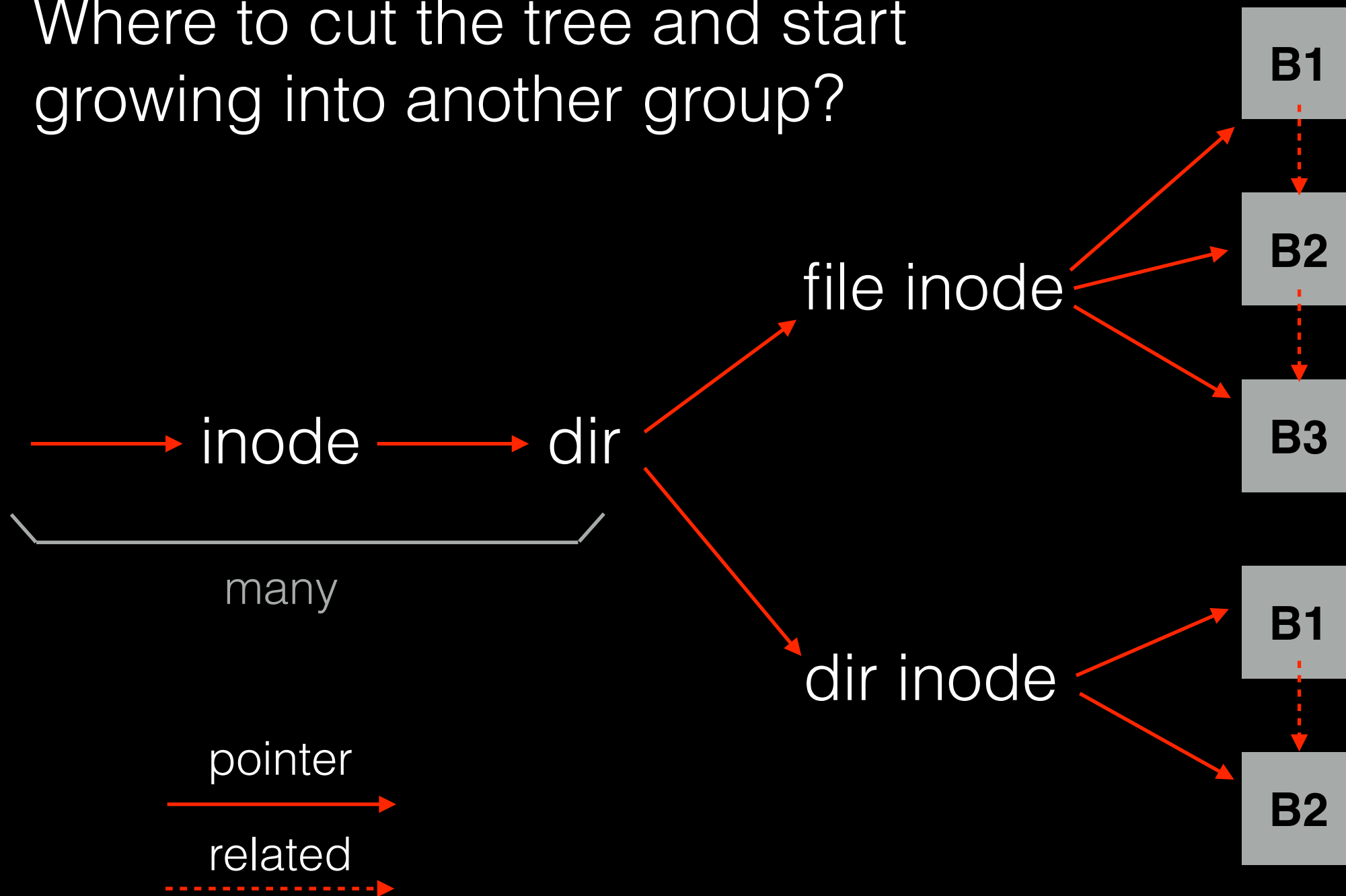FFS developers used their best judgement.

# FFS: Two-Level Allocator

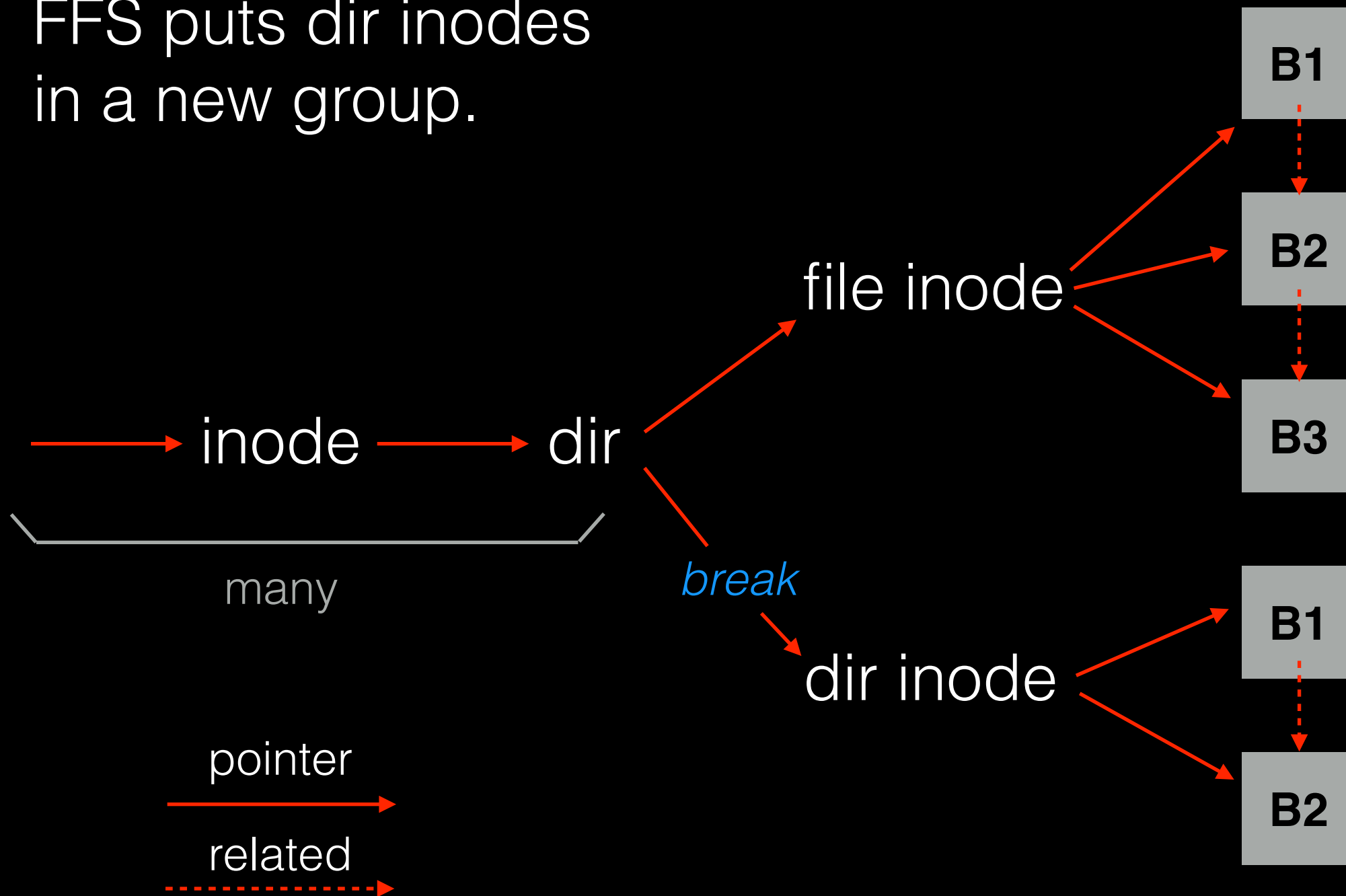Level 1: decide which group

Level 2: decide where in group

# Where to cut the tree and start growing into another group?

FFS puts dir inodes
in a new group.

inode → dir

many

file inode

B1
B2
B3

break

dir inode

B1
B2

pointer

related

"ls" is fast on directories with many files.

inode → dir

file inode

dir inode

*break*

many

B1
B2
B3

B1
B2

pointer

related

# Preferences

**File inodes**: allocate in <u>same</u> group with dir

**Dir inodes**: allocate in <u>new</u> group with fewer inodes than the average group
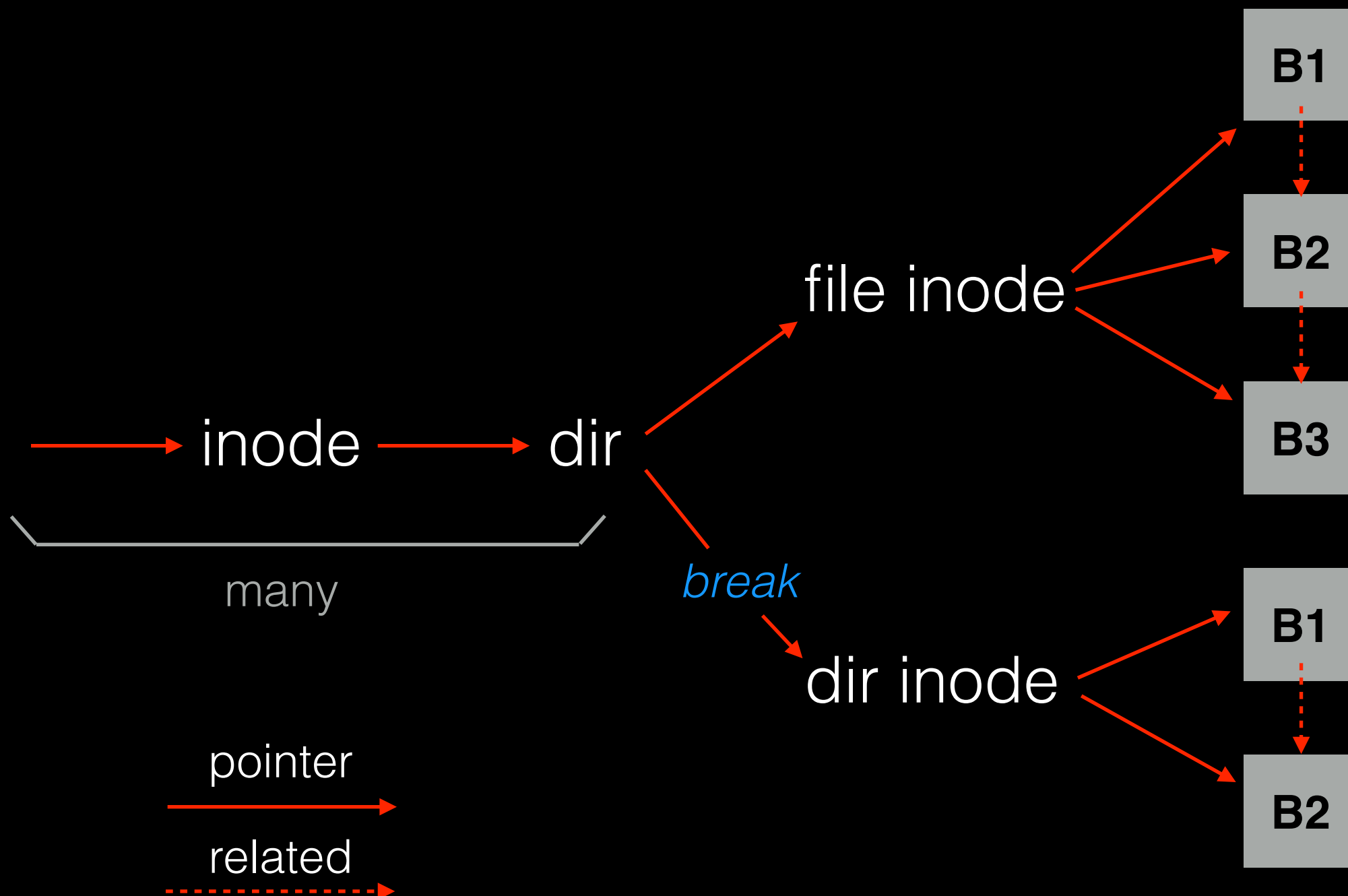
**First data block**: allocate near inode

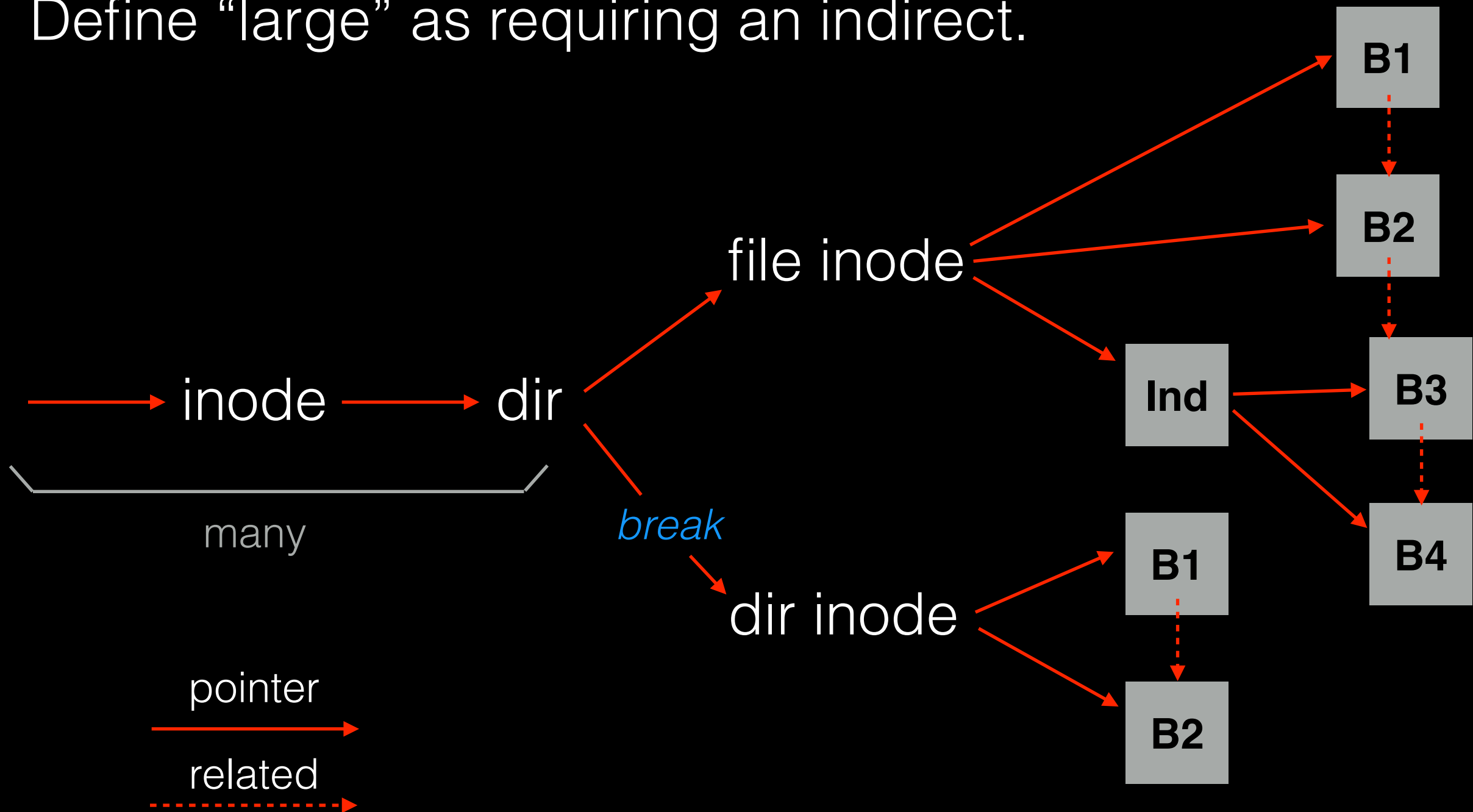**Other data blocks**: allocate near previous block

# Problem: Large Files

A single large file can use nearly all of a group.

This displaces data for many small files.

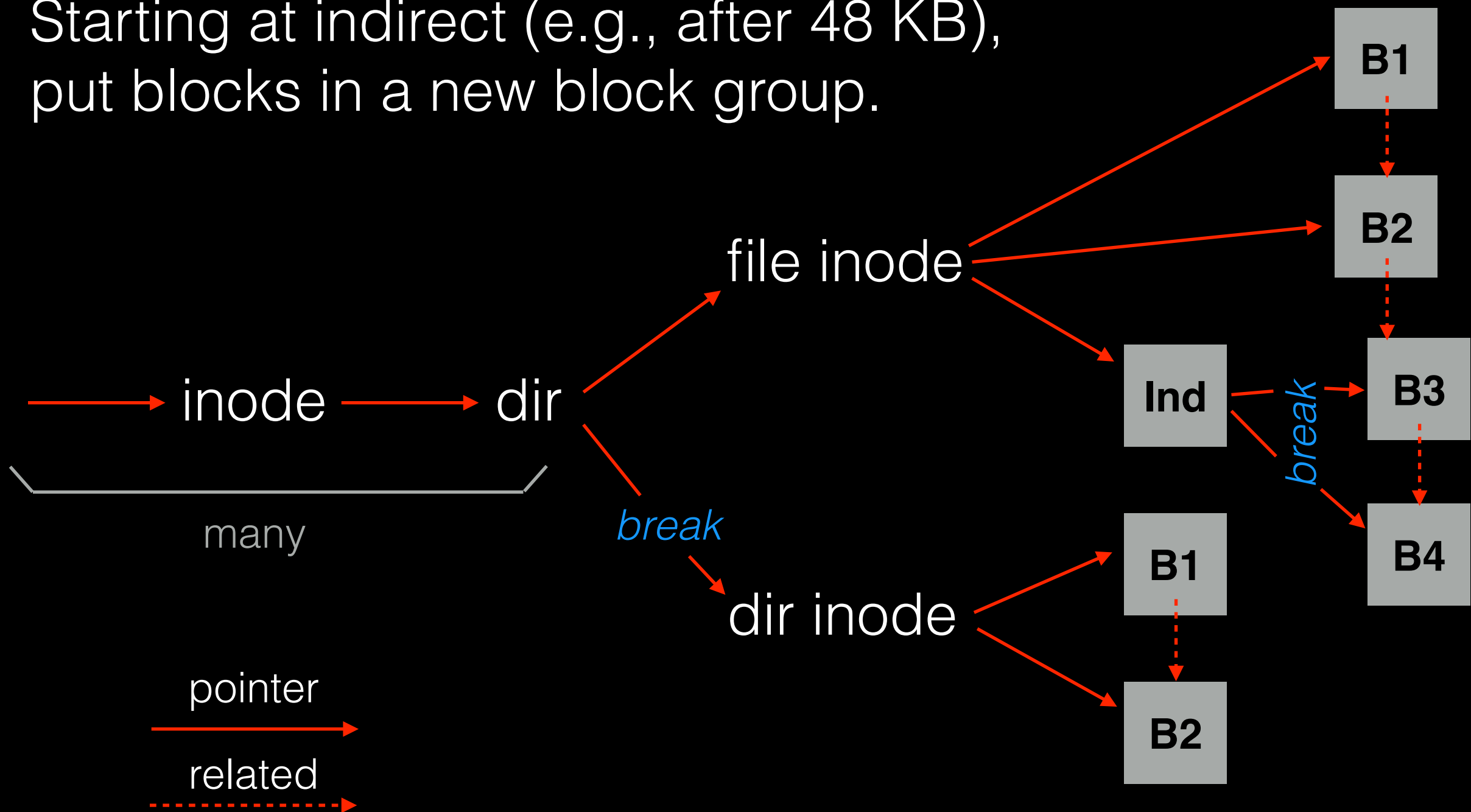It's better to do one seek for the large file than one seek for each of many small files.

Define "large" as requiring an indirect.

inode → dir

many

file inode

dir inode

B1

B2

Ind

B3

B4

B1

B2

break

pointer

related

Starting at indirect (e.g., after 48 KB),
put blocks in a new block group.



inode → dir

many

*break*

file inode

dir inode

B1

B2

B3

B4

Ind

*break*

B1

B2

pointer

related

# Preferences

**File inodes**: allocate in <u>same</u> group with dir

**Dir inodes**: allocate in <u>new</u> group with fewer inodes than the average group

**First data block**: allocate near inode

**Other data blocks**: allocate near previous block

**Large file data blocks**: after 48KB, go to <u>new</u> group.  Move to another group (w/ fewer than avg blocks) every subsequent 1MB.

# Preferences

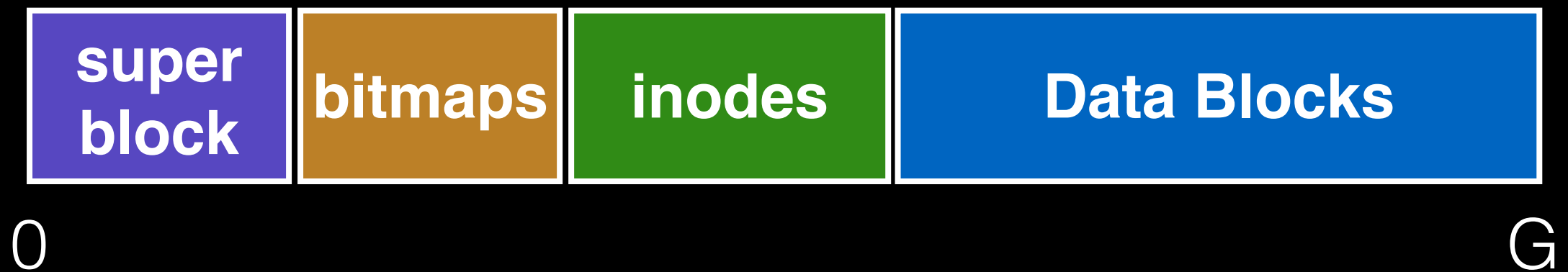**File inodes**: allocate in <u>same</u> group with dir

**Dir inodes**: allocate in <u>new</u> group with <u>fewer inodes than the average group</u>

**First data block**: allocate near inode

**Other data blocks**: allocate near previous block

**Large file data blocks**: after 48KB, go to <u>new</u> group.  Move to another group (w/ <u>fewer than avg blocks</u>) every subsequent 1MB.

# Group Descriptor (aka Summary Block)

# Group Descriptor (aka Summary Block)



how many free inodes, data blocks?

# Techniques

Bitmaps
Locality groups
Rotated super
Large blocks
Fragments
Smart allocation

# Conclusion

First disk-aware file system.

FFS inspired modern files systems, including ext2 and ext3.

FFS also introduced several new features:
 - long file names
 - atomic rename
 - symbolic links

# Advice

All hardware is unique.

Treat disk like disk!

Treat flash like flash!

Treat random-access memory like random-access memory!

# Advice

All hardware is unique.

Treat disk like disk!

Treat flash like flash!

Treat random-access memory like random-access memory!
(actually don't -- the name is a lie)

# Announcements

**Exam** this Friday
 - 7-9pm, CHEM 1351 (same as last time)

**Review** this Wednesday
 - 7-9pm, room TBD.  Bring questions.

**Office hours** Monday (today)
 - after class, in lab