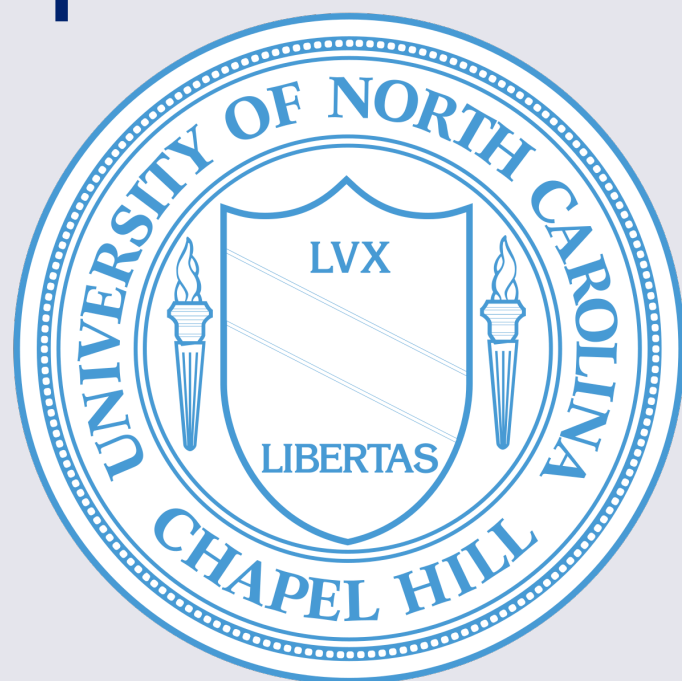


How to Copy Files

Yang Zhan^{1,2}, **Alex Conway**^{3,4}, Ian Groombridge⁵, Yizheng Jiao¹,
Nirjhar Mukherjee¹, Michael A. Bender⁶, Martín Farach-Colton³,
William Jannen⁷, Rob Johnson⁴, Donald E. Porter¹, Jun Yuan⁵

1



2



3



4



5



6



7



Copying is Ubiquitous and Important

`cp -r`

`vmrun start`

container instantiation

backup

Physical Copy

Copying

Physical Copy

High Latency

High Space Use

Existing Logical Copy Implementations

BTRFS

Leverages the underlying copy-on-write B-tree to implement
`cp --reflink`

XFS

Uses an update-in-place B-tree but supports sharing data
blocks with copy-on-write via `cp --reflink`

ZFS

Implements a limited version of copy-on-write copying via
`zfs clone`

Copying

Physical Copy

Copy on Write

High Latency

High Space Use

Copying

Physical Copy

High Latency

High Space Use

Copy on Write

Low Latency

Better Space Use

High Fragmentation

Copying

Physical Copy

High Latency

High Space Use

Copy on Write

Low Latency

Better Space Use

High Fragmentation

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

Init: 64 4MiB files with random data.

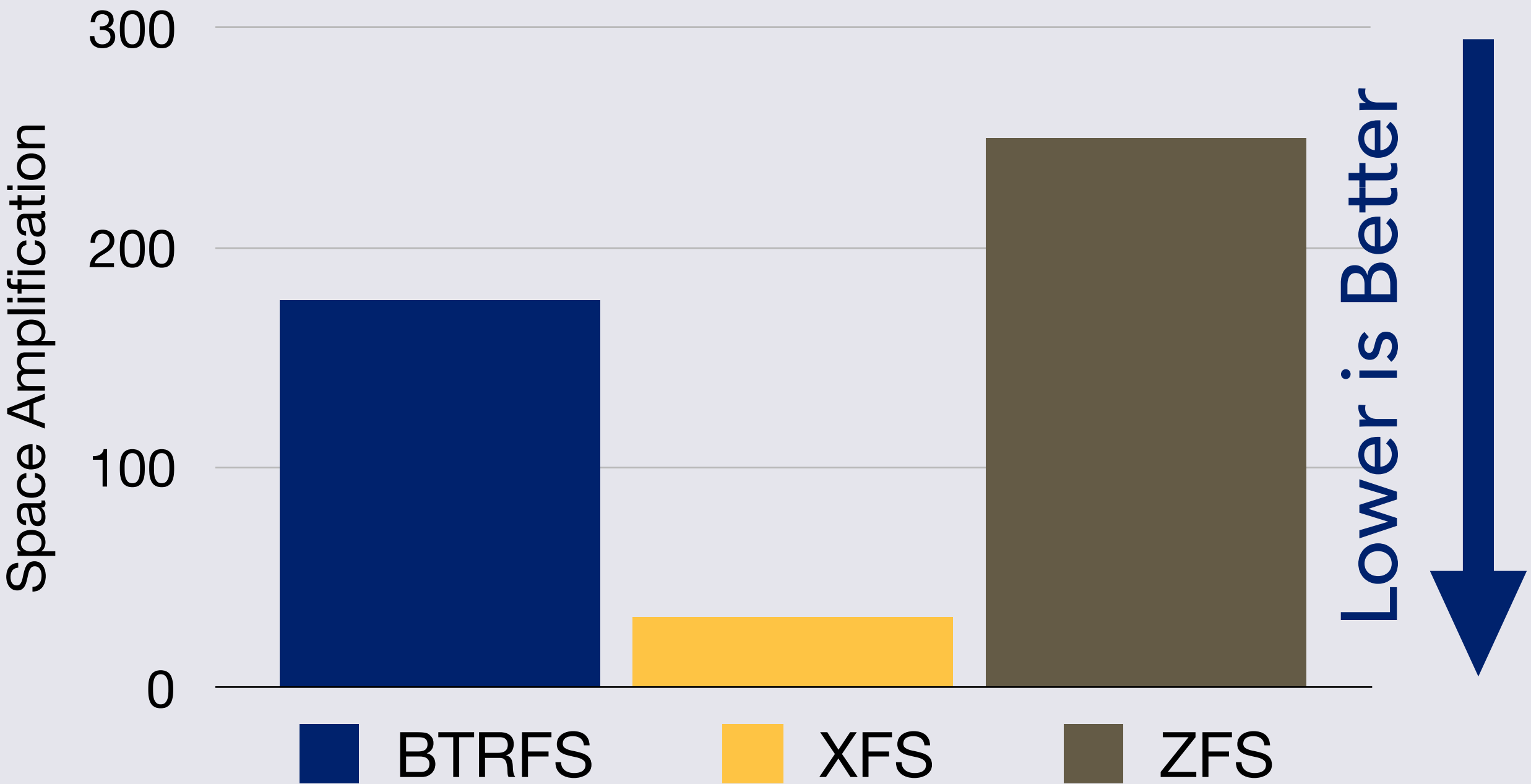
Each round: logically copy all files, then change 16B in each file (1KiB total)

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

Init: 64 4MiB files with random data.

Each round: logically copy all files, then change 16B in each file (1KiB total)



Space Amplification

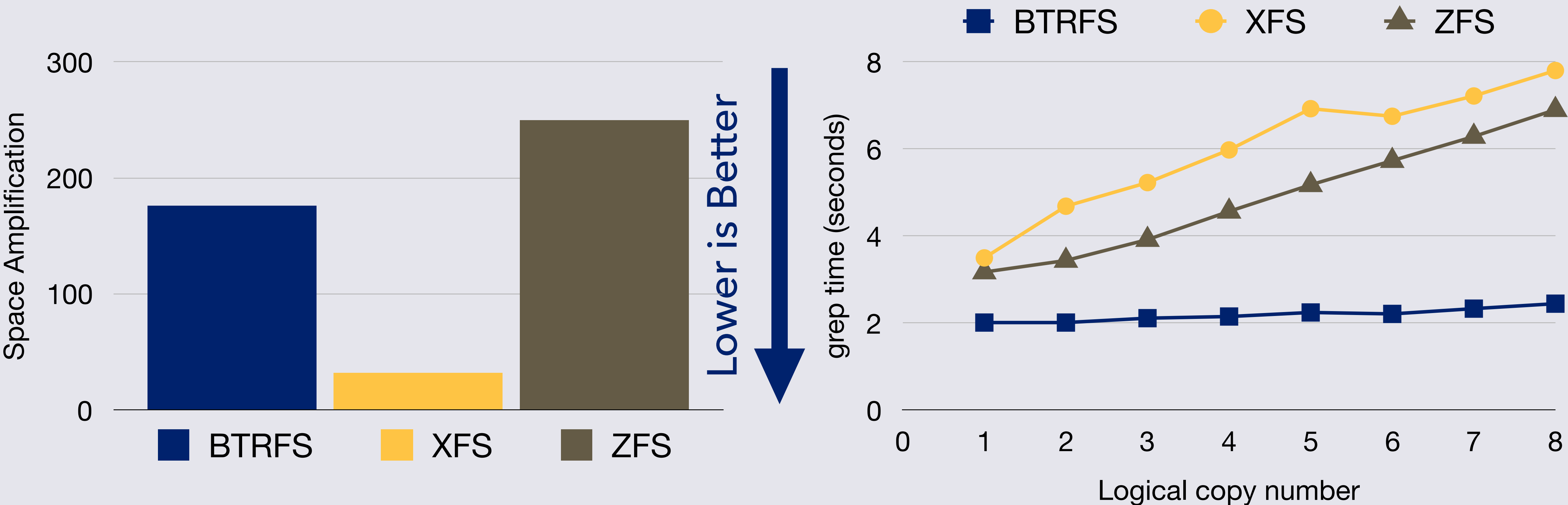
additional file system size / added data (1KiB)

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

Init: 64 4MiB files with random data.

Each round: logically copy all files, then change 16B in each file (1KiB total)



Space Amplification

additional file system size / added data (1KiB)

Fragmentation

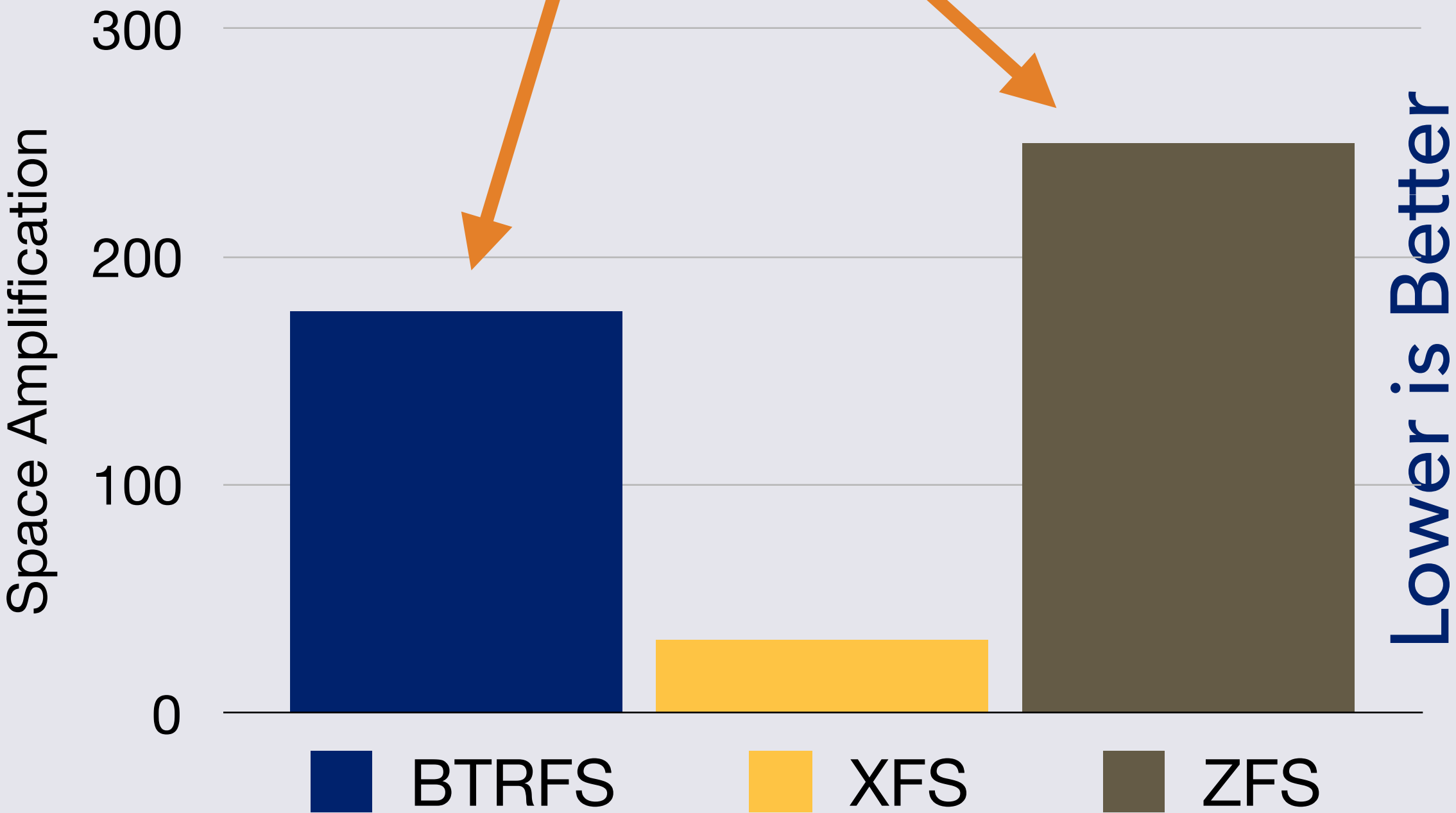
measured by timing a grep over the latest copy

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

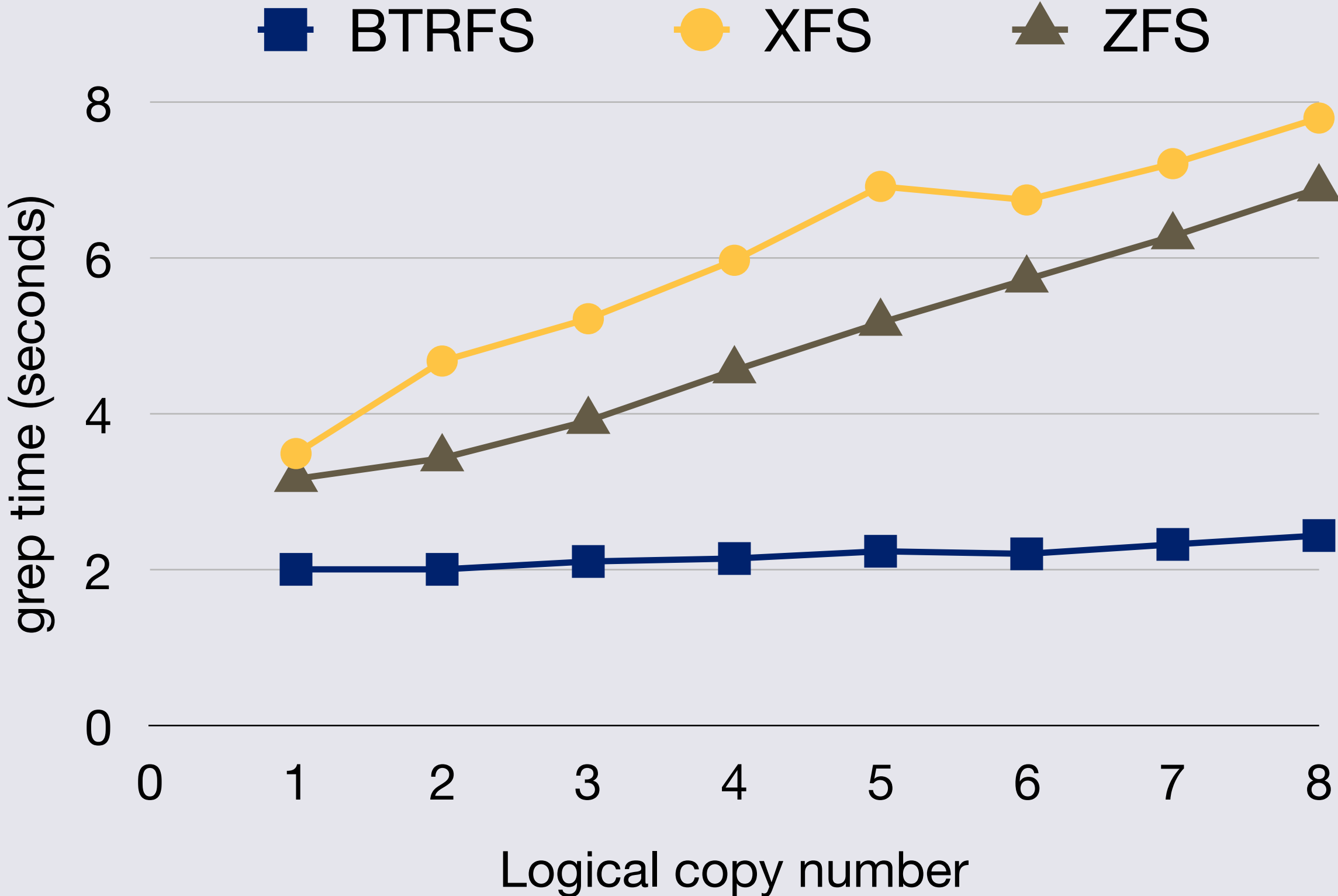
Initial 44 4MiB files with random data.
Each file is 16B, then change 16B in each file (1KiB total)

Have large space amplification



Space Amplification

additional file system size / added data (1KiB)



Fragmentation

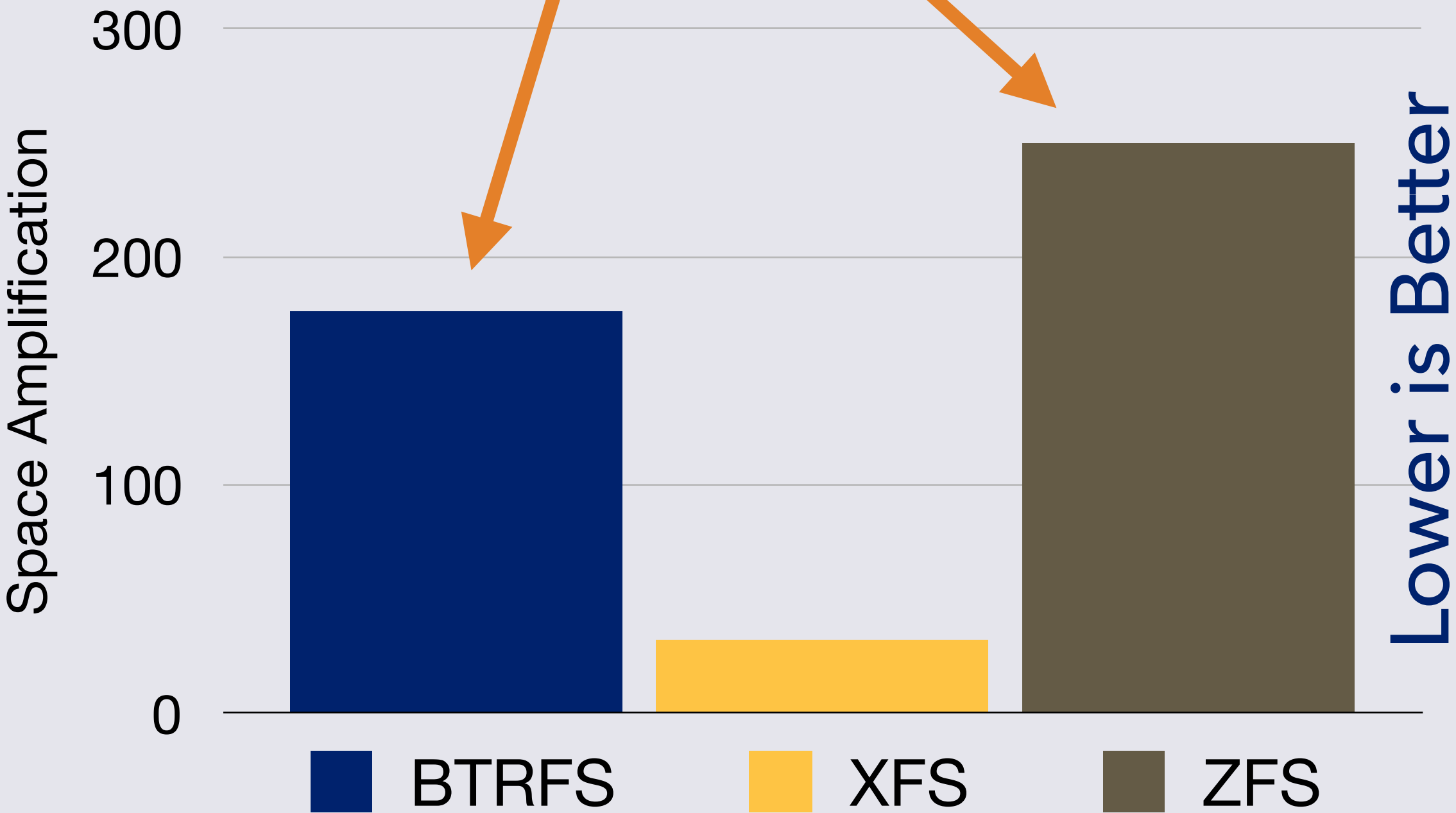
measured by timing a grep over the latest copy

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

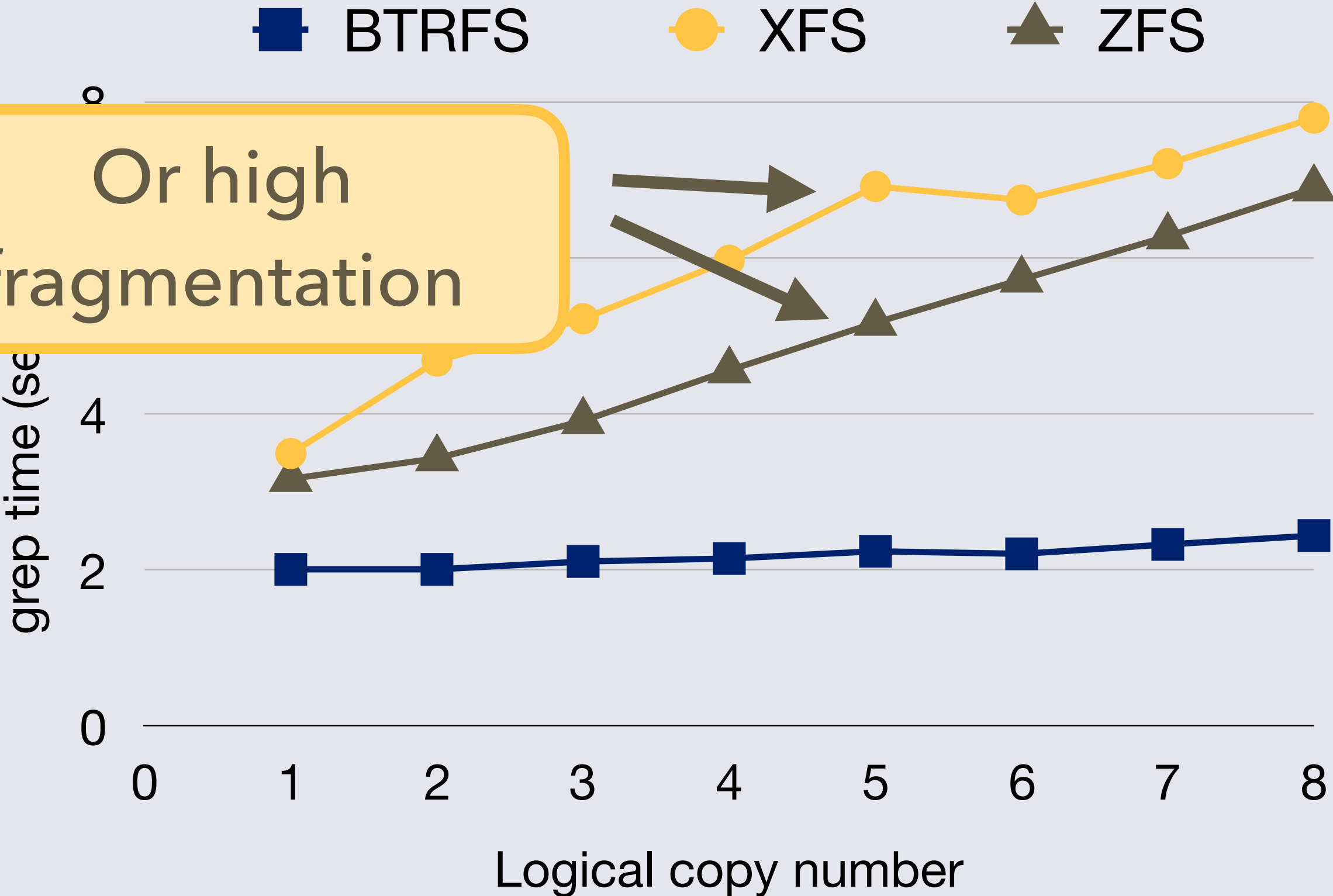
Initial 44 4MiB files with random data.
Each file has 16B of data. Change 16B in each file (1KiB total)

Have large space amplification



Space Amplification

additional file system size / added data (1KiB)



Fragmentation

measured by timing a grep over the latest copy

Copy Performance Goals

Copy Performance Goals

Space efficient

Low latency

Specific to Copying

Copy Performance Goals

Space efficient

Low latency

Specific to Copying

Fast writes

Fast reads

General file system

Copy Performance Goals

Space efficient

Low latency

Specific to Copying

Locality

Fast writes

Fast reads

General file system

Contributions of this Paper

A high performance logical copy implementation...

Contributions of this Paper

A high performance logical copy implementation...

In BεtrFS, which leverages the properties of Copy-on-Write B^ε-trees

Contributions of this Paper

A high performance logical copy implementation...

In BetrFS, which leverages the properties of Copy-on-Write B^ε-trees

Space efficient

Low latency

Copy-specific

Fast writes

Fast reads

General file system

Contributions of this Paper

A high performance logical copy implementation...

In BetrFS, which leverages the properties of Copy-on-Write B^ε-trees

Space efficient



Fast writes



Low latency



Fast reads



Copy-specific

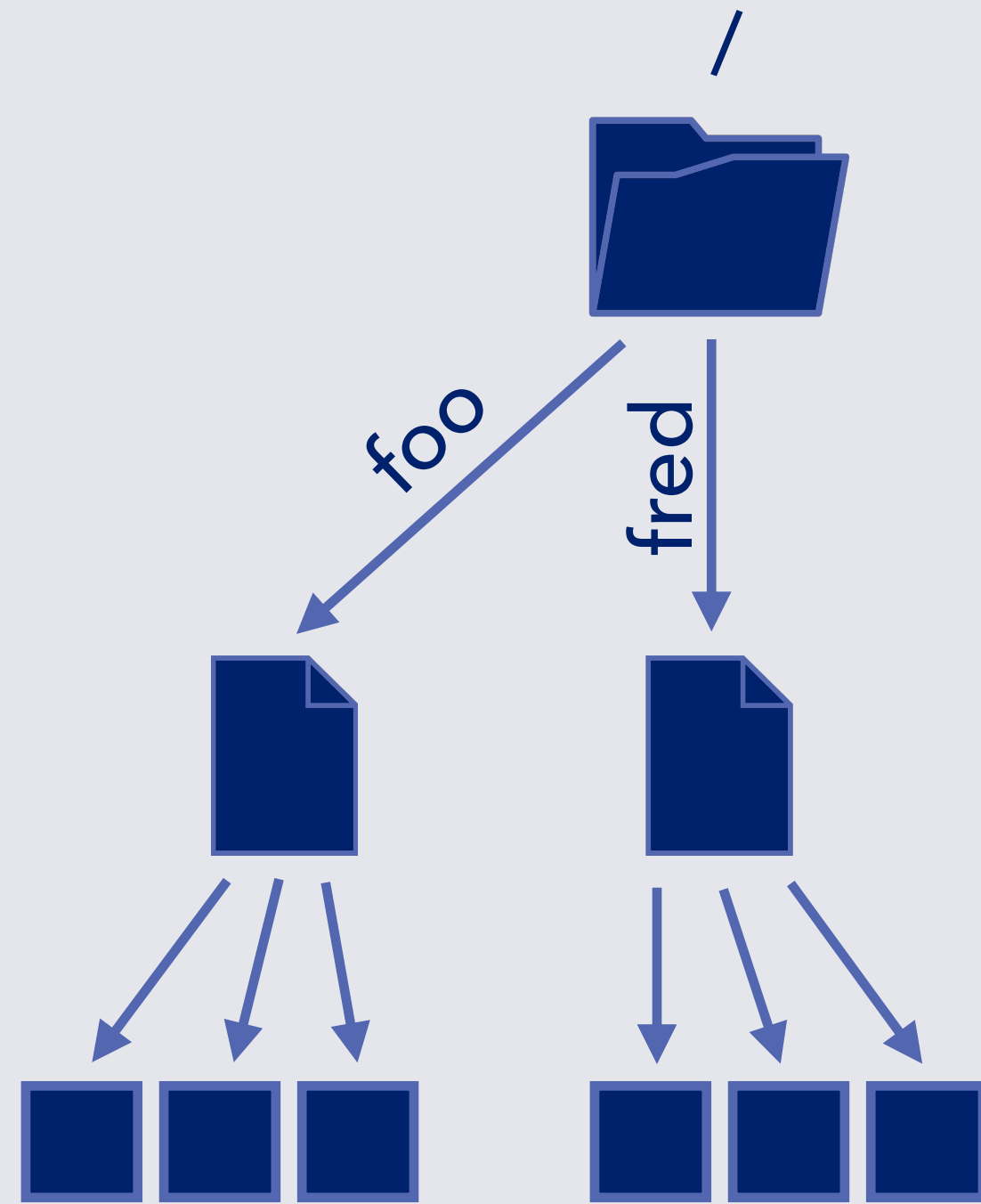
General file system

What is the Challenge of
Logical Copy?

Example: Logical copy in an
inode file system

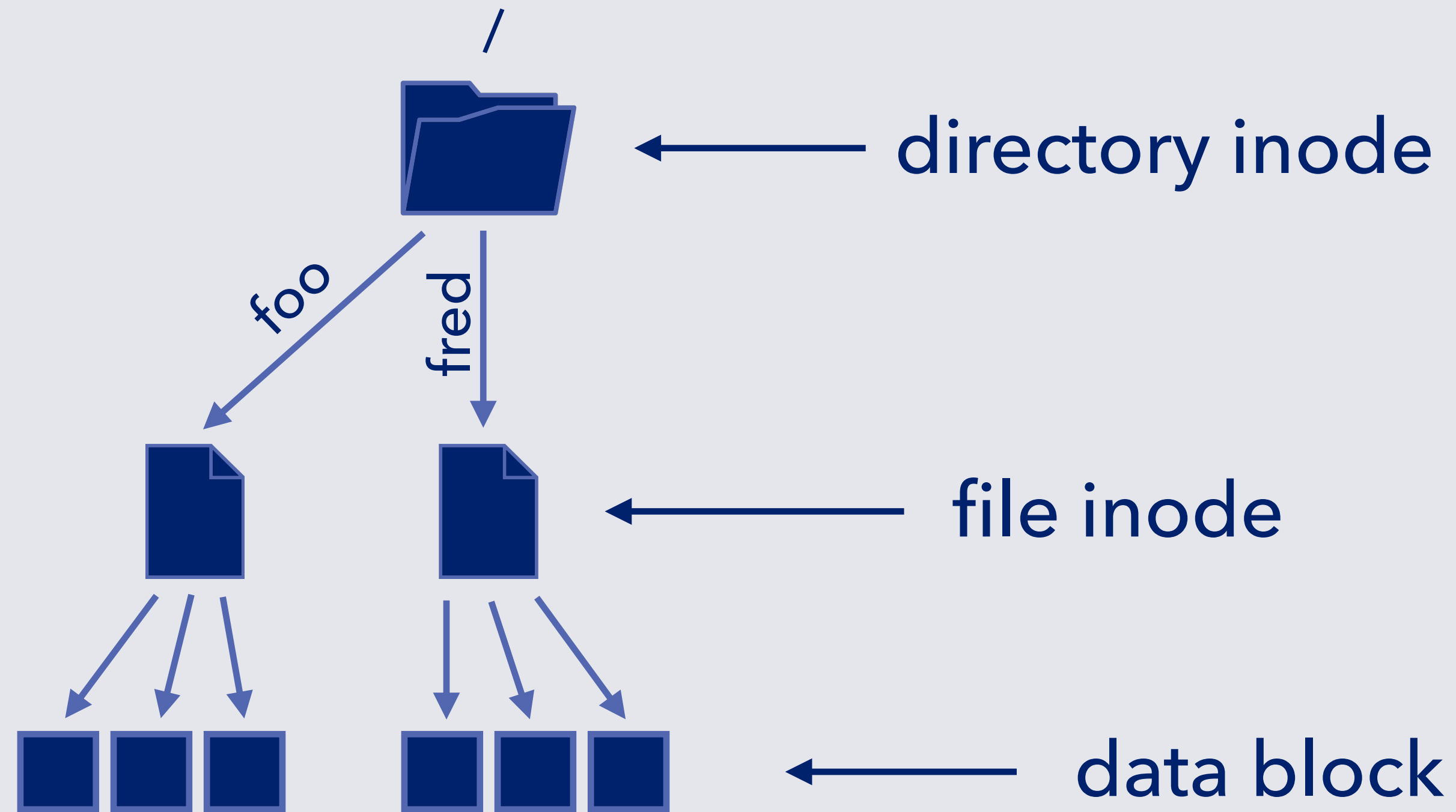
Logical Copy in an Inode File System

Copy /foo to /bar



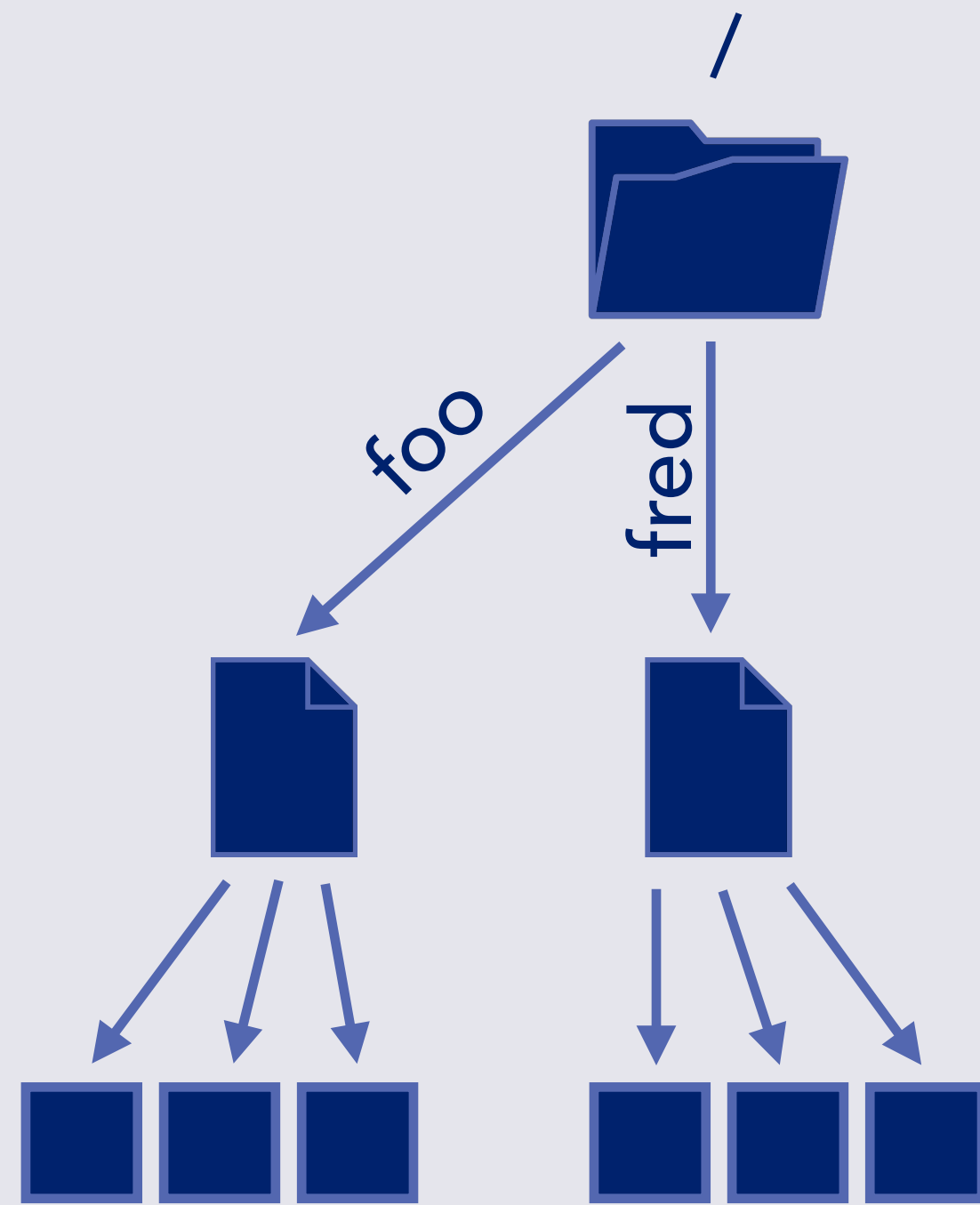
Logical Copy in an Inode File System

Copy /foo to /bar



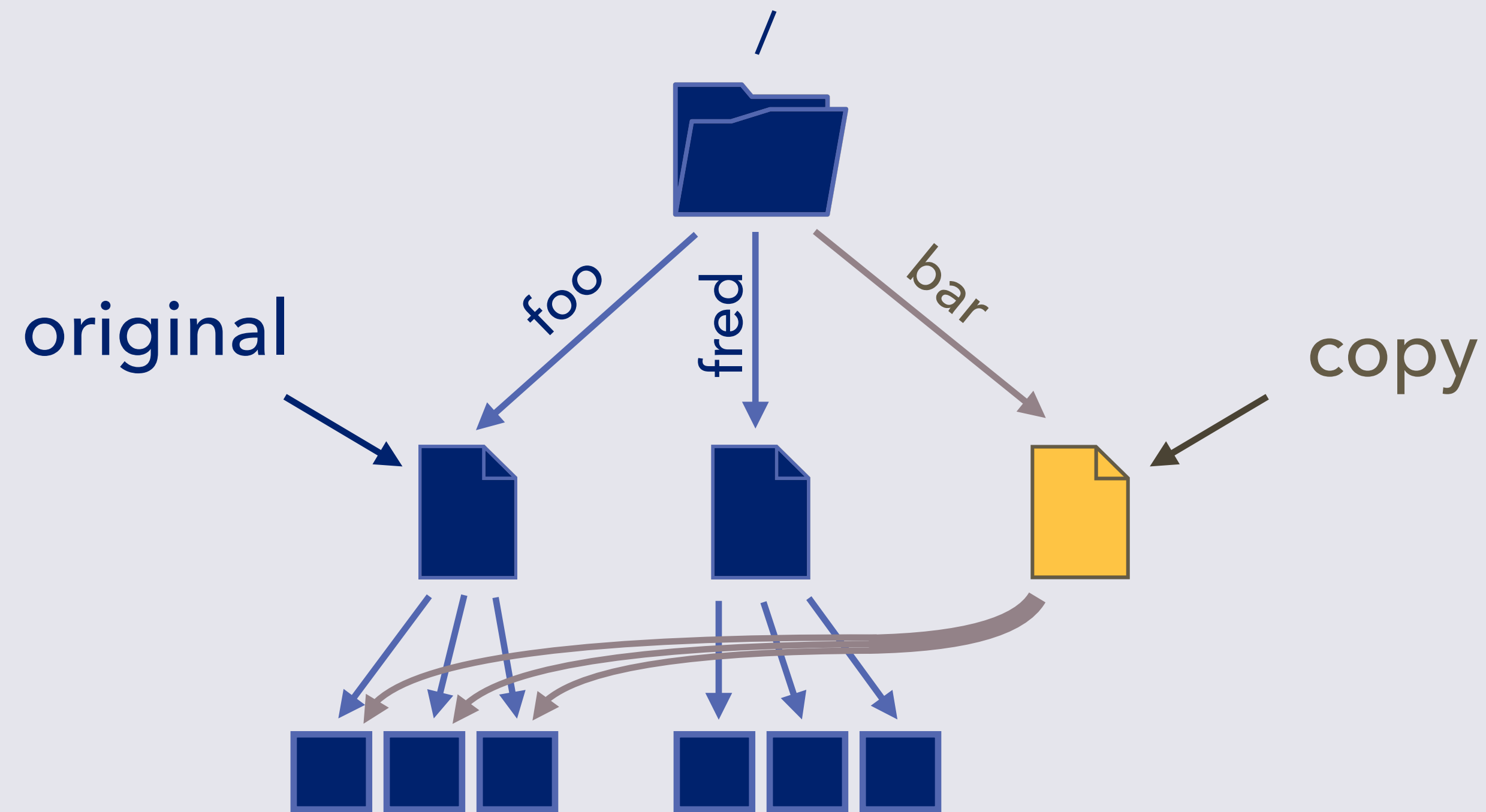
Logical Copy in an Inode File System

Copy /foo to /bar



Logical Copy in an Inode File System

Copy /foo to /bar



Logical Copy in an Inode File System

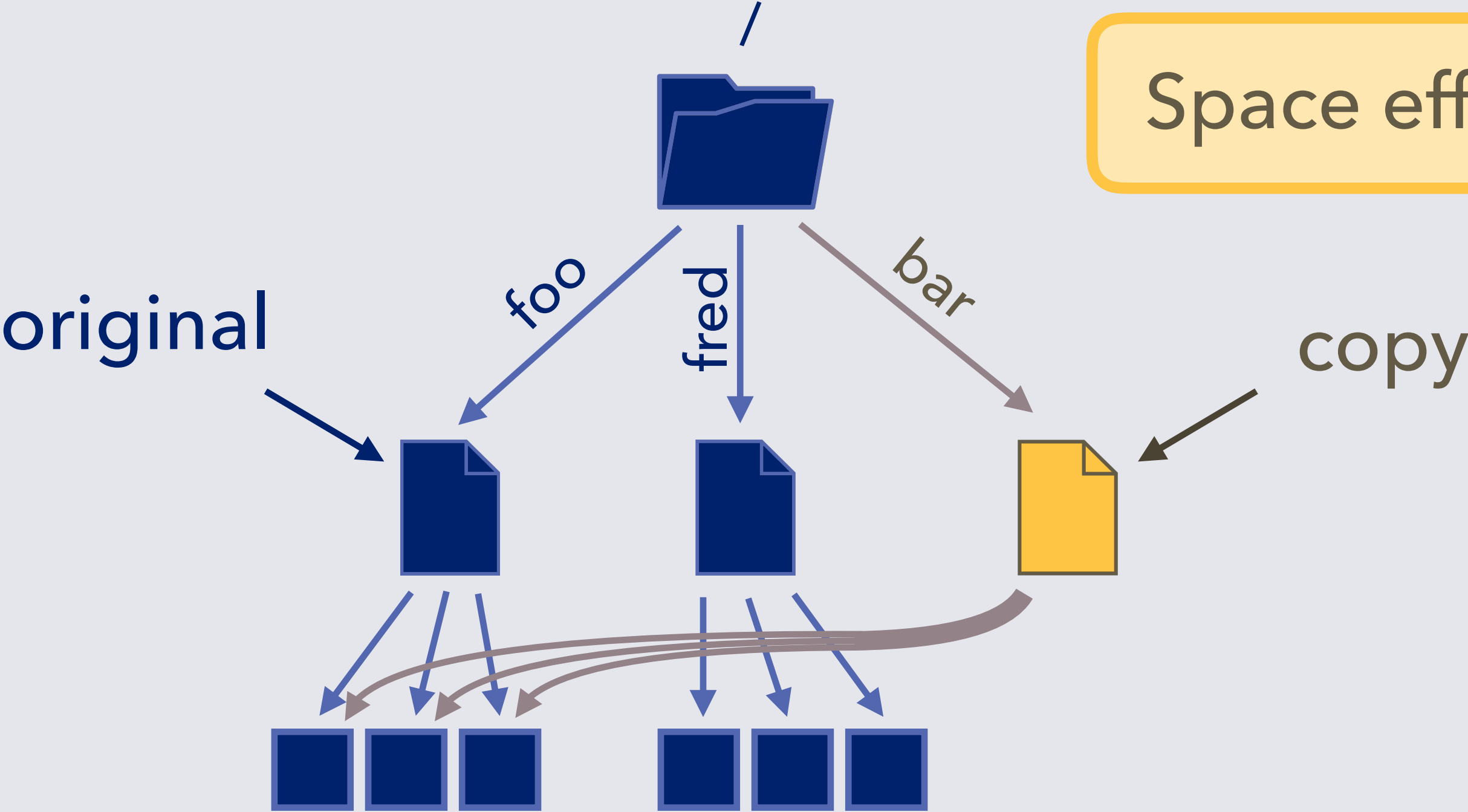
Copy /foo to /bar

Low latency ?

Fast Reads ?

Space efficient ?

Fast Writes ?



Logical Copy in an Inode File System

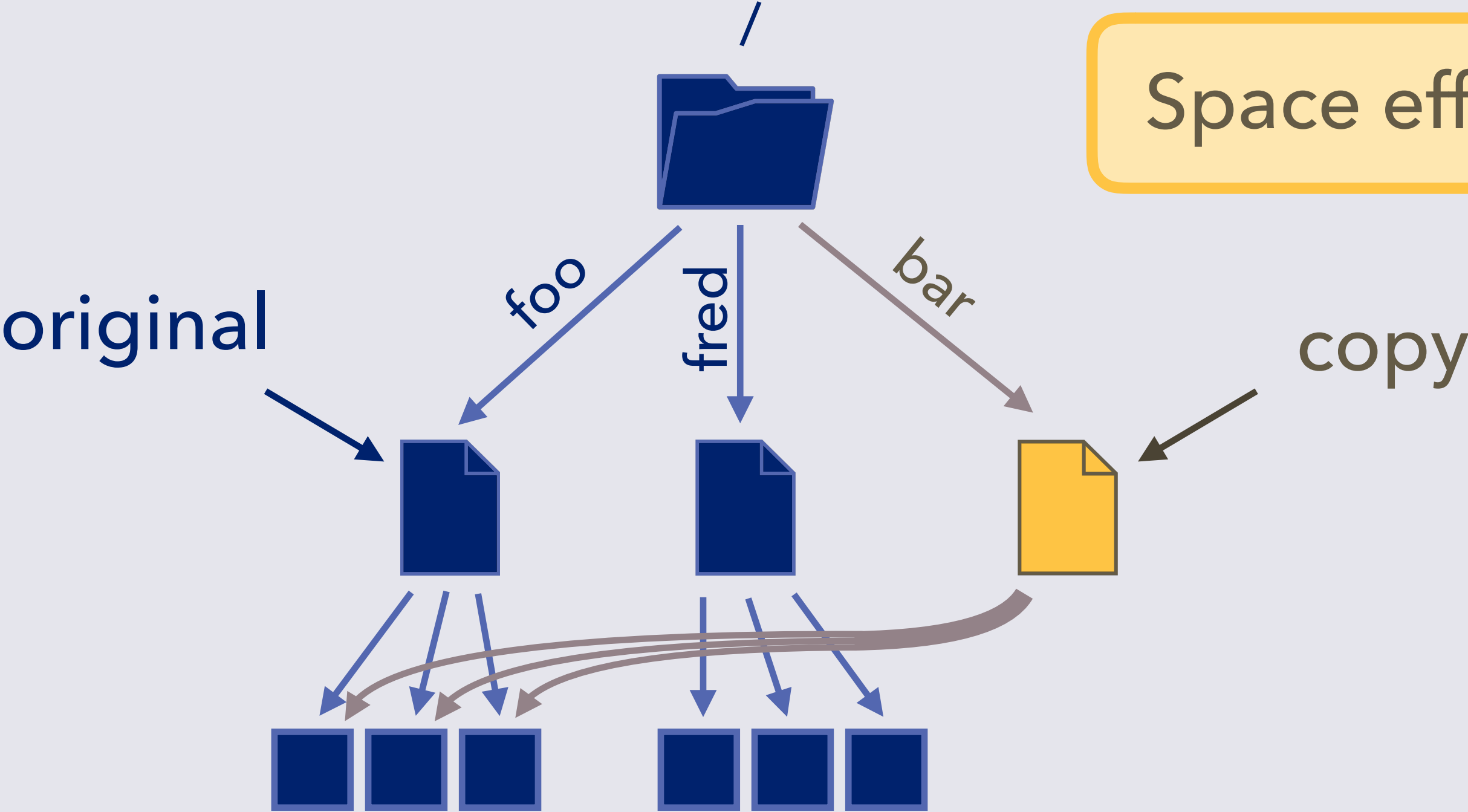
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ?

Fast Writes ?



Logical Copy in an Inode File System

Copy /foo to /bar

Low latency



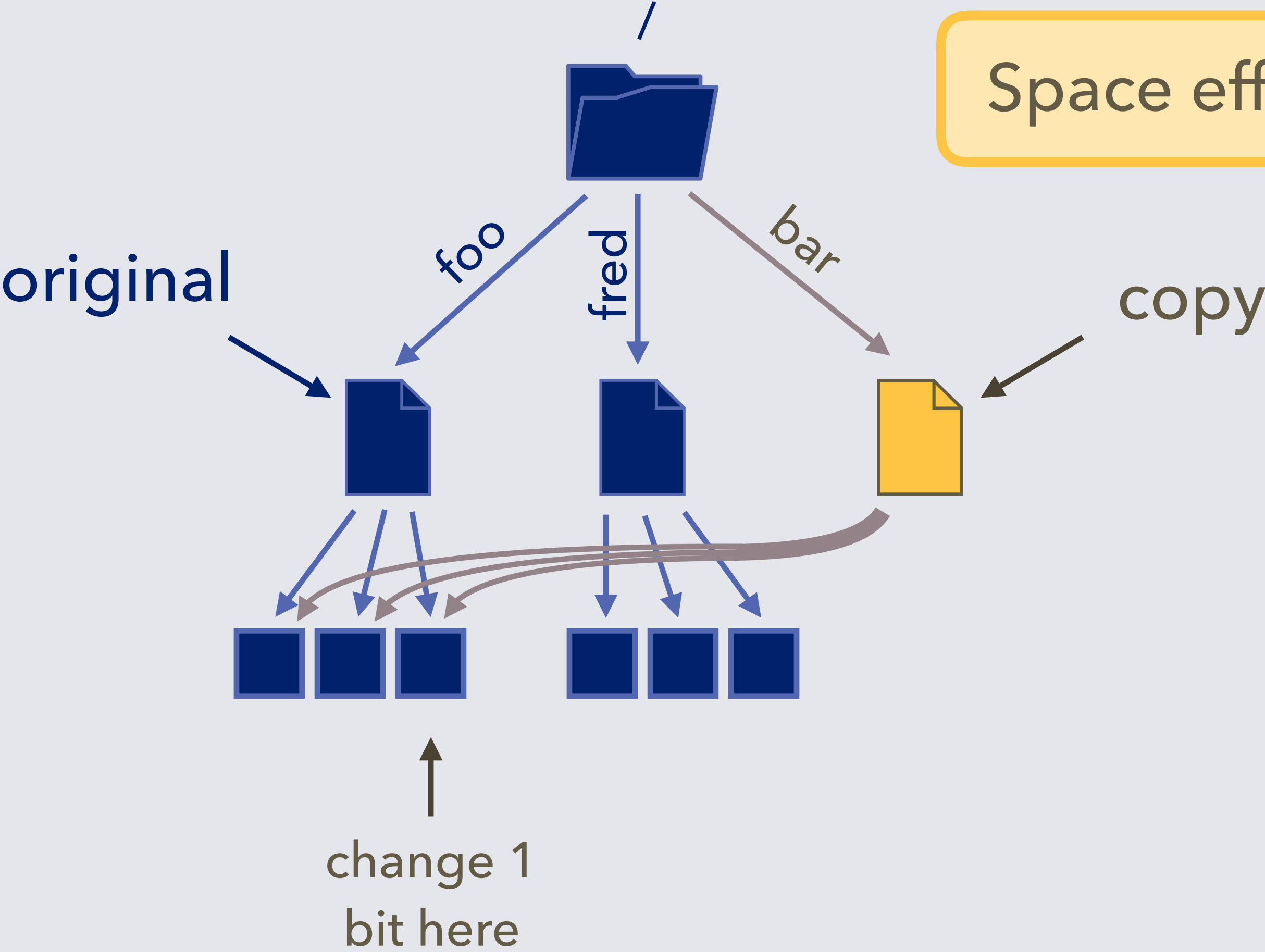
Fast Reads



Space efficient



Fast Writes



Logical Copy in an Inode File System

Copy /foo to /bar

Low latency



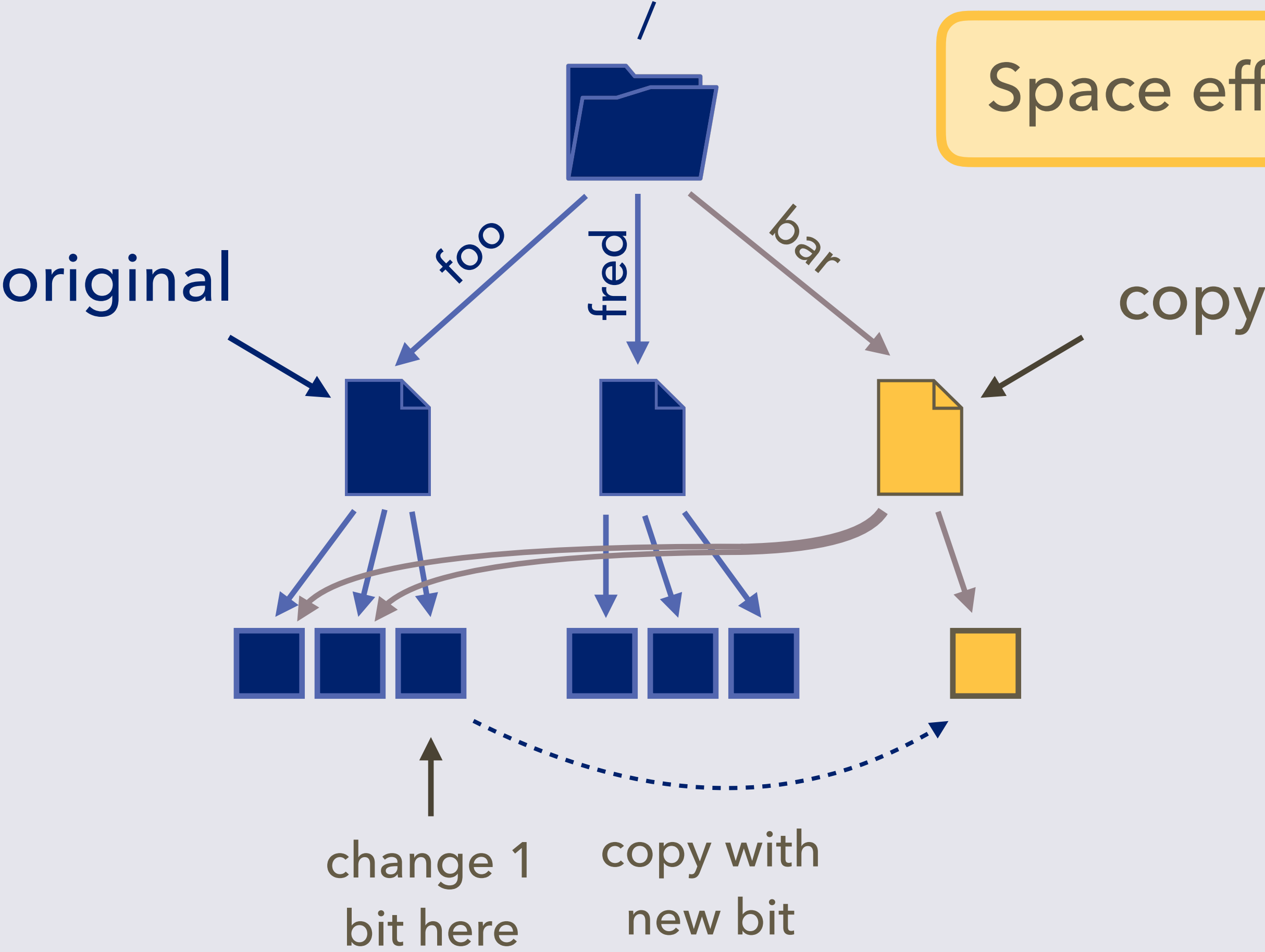
Fast Reads



Space efficient



Fast Writes



Logical Copy in an Inode File System

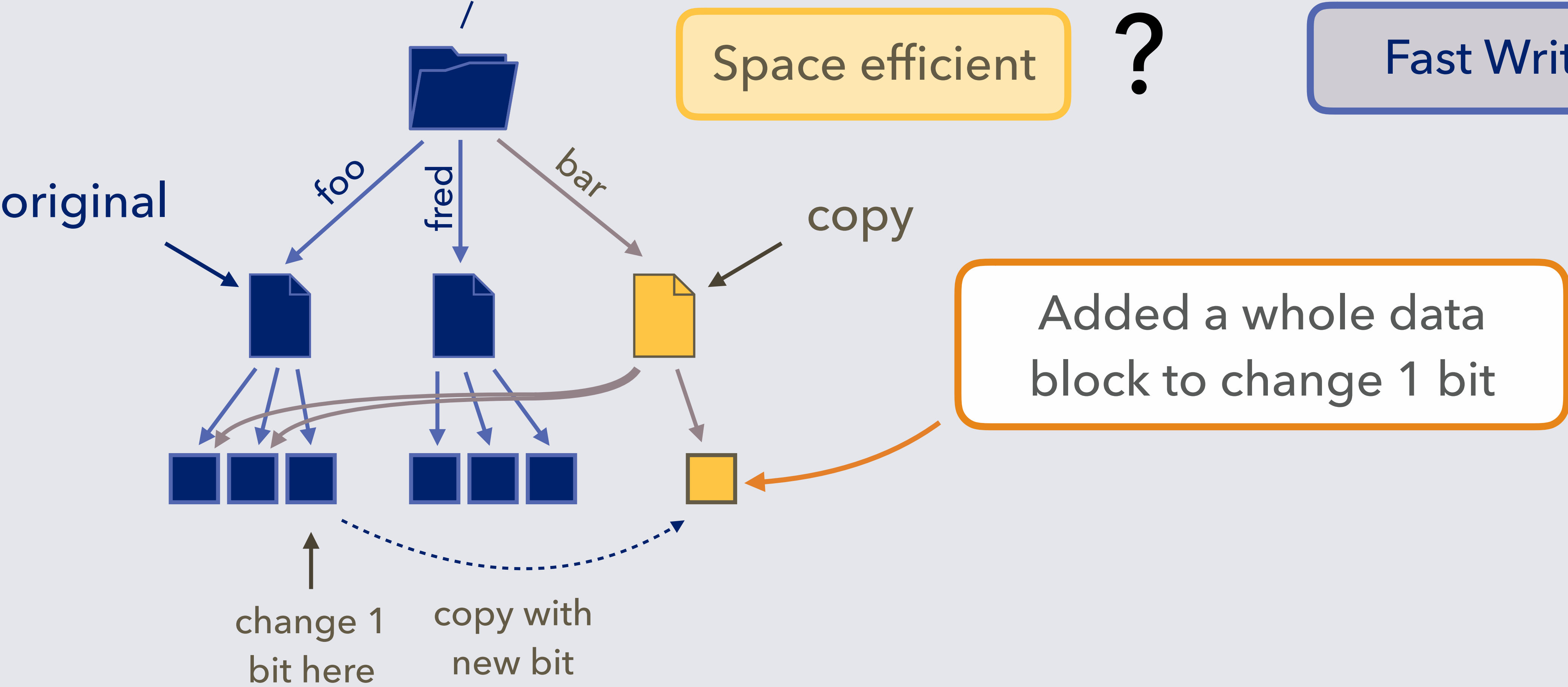
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ?

Fast Writes ?



Logical Copy in an Inode File System

Copy /foo to /bar

Low latency



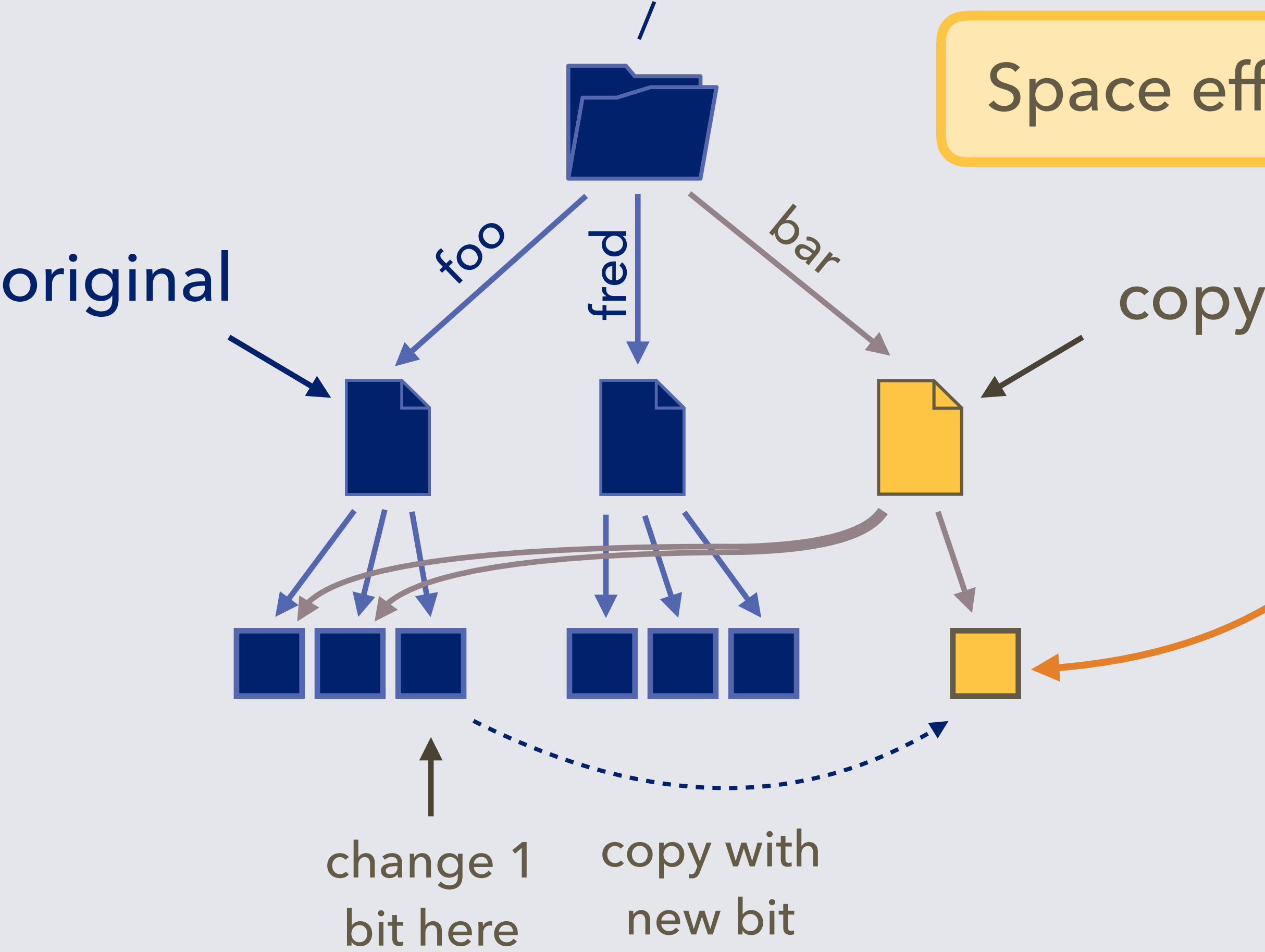
Fast Reads



Space efficient



Fast Writes



Added a whole data block to change 1 bit

This is at least 4KiB and can be more

Logical Copy in an Inode File System

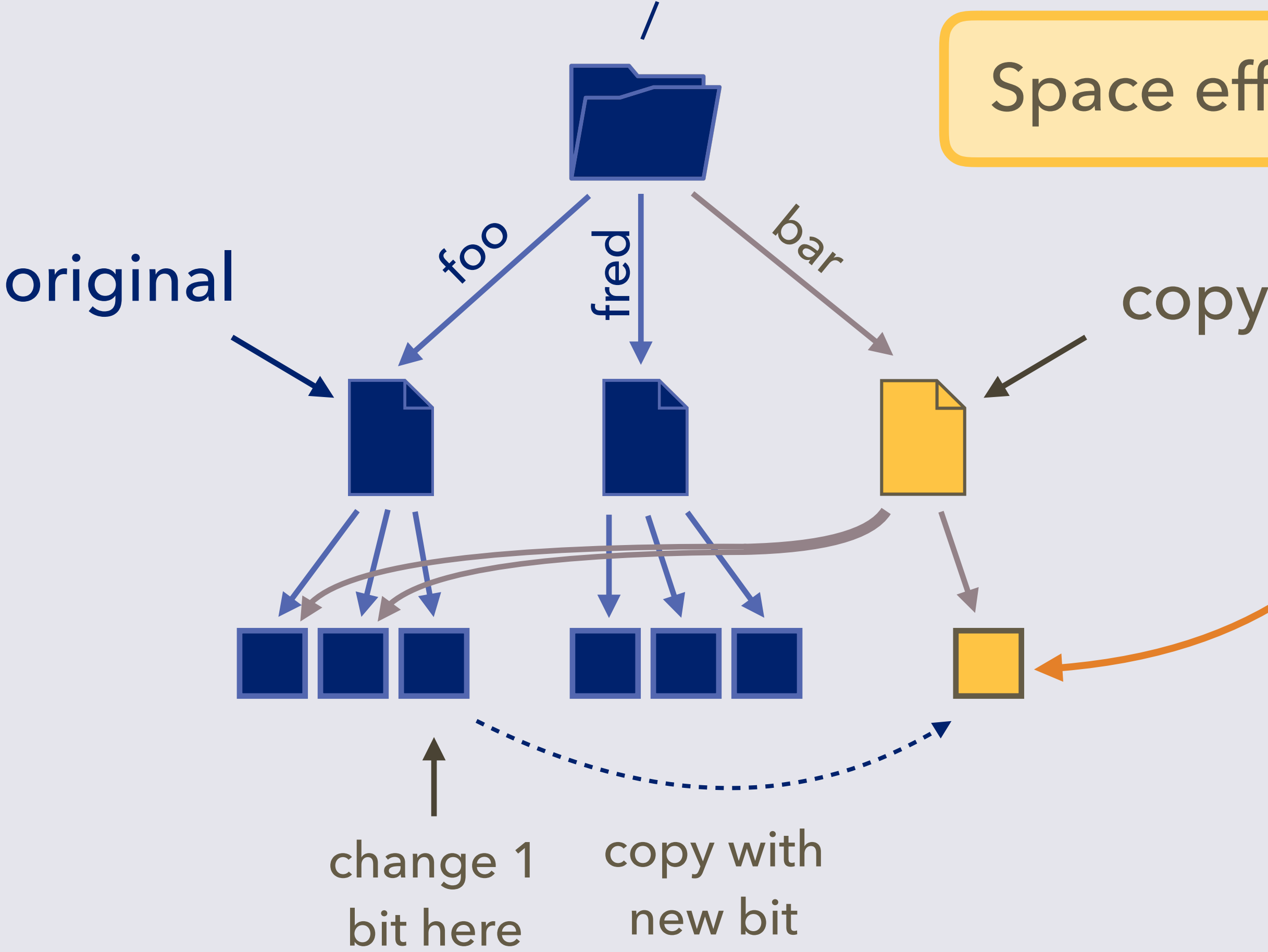
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ✗

Fast Writes ?



Added a whole data block to change 1 bit

This is at least 4KiB and can be more

Logical Copy in an Inode File System

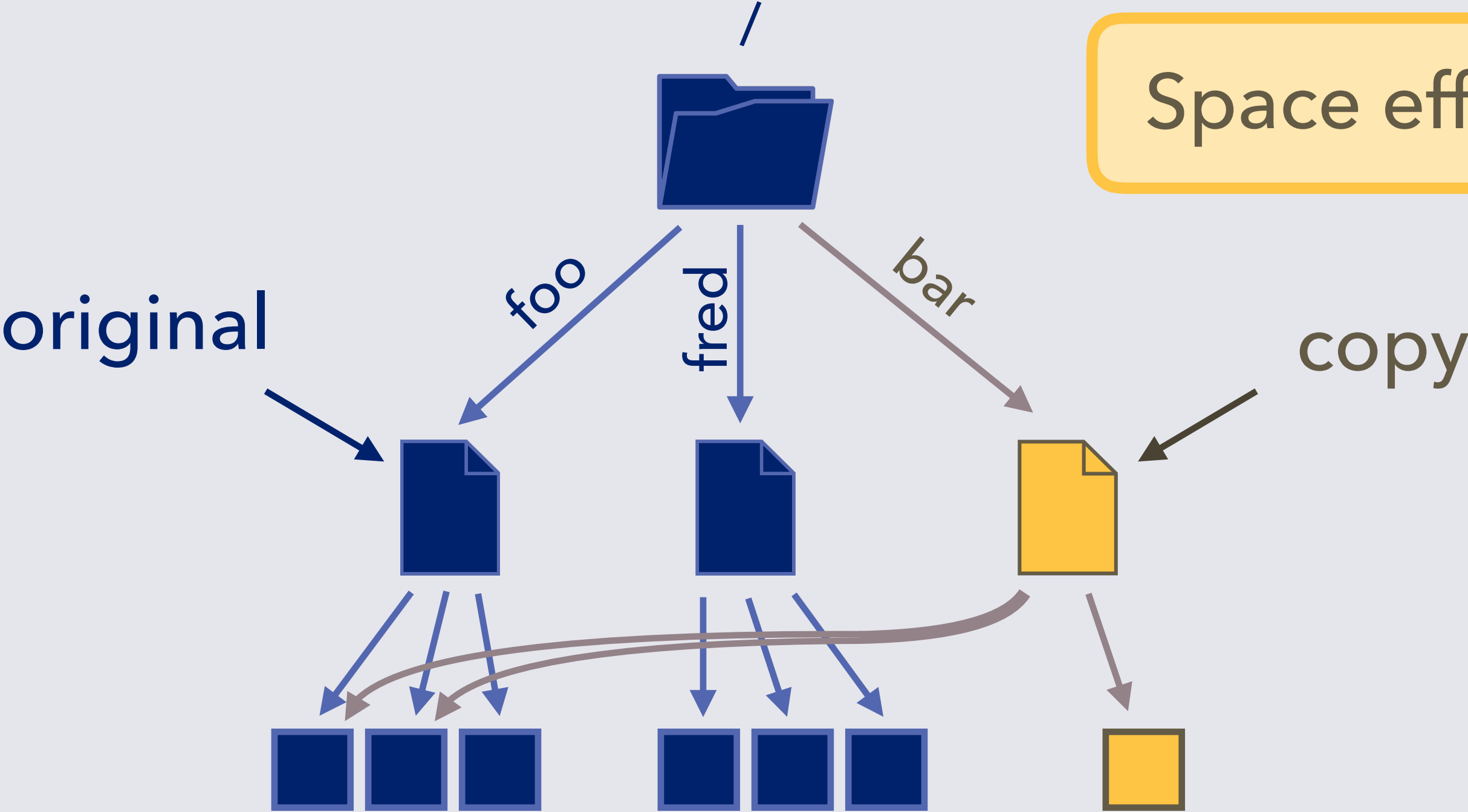
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ✗

Fast Writes ?



Logical Copy in an Inode File System

Copy /foo to /bar

Low latency



Fast Reads

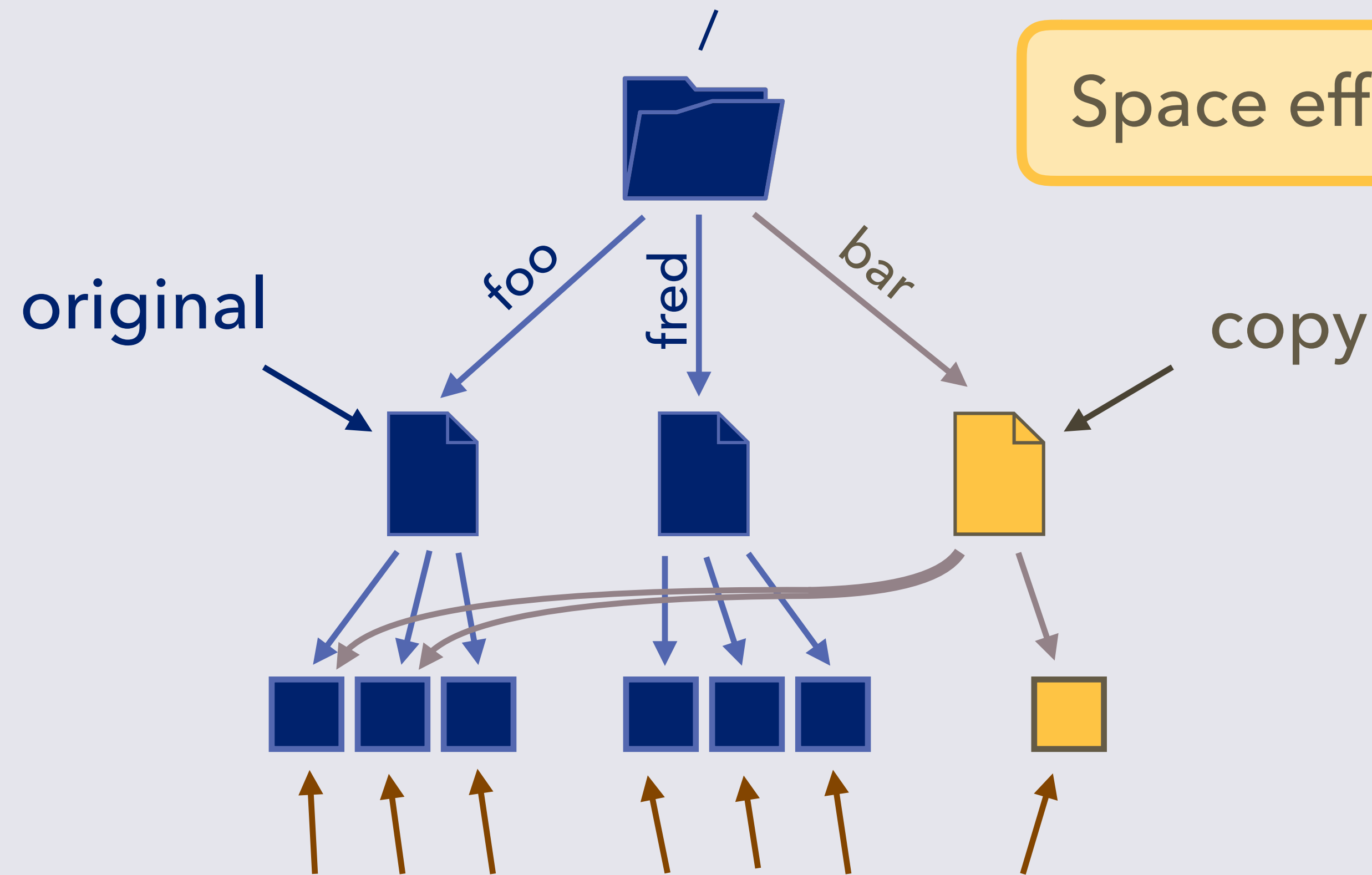
?

Space efficient



Fast Writes

?



no locality guarantees between data blocks

Logical Copy in an Inode File System

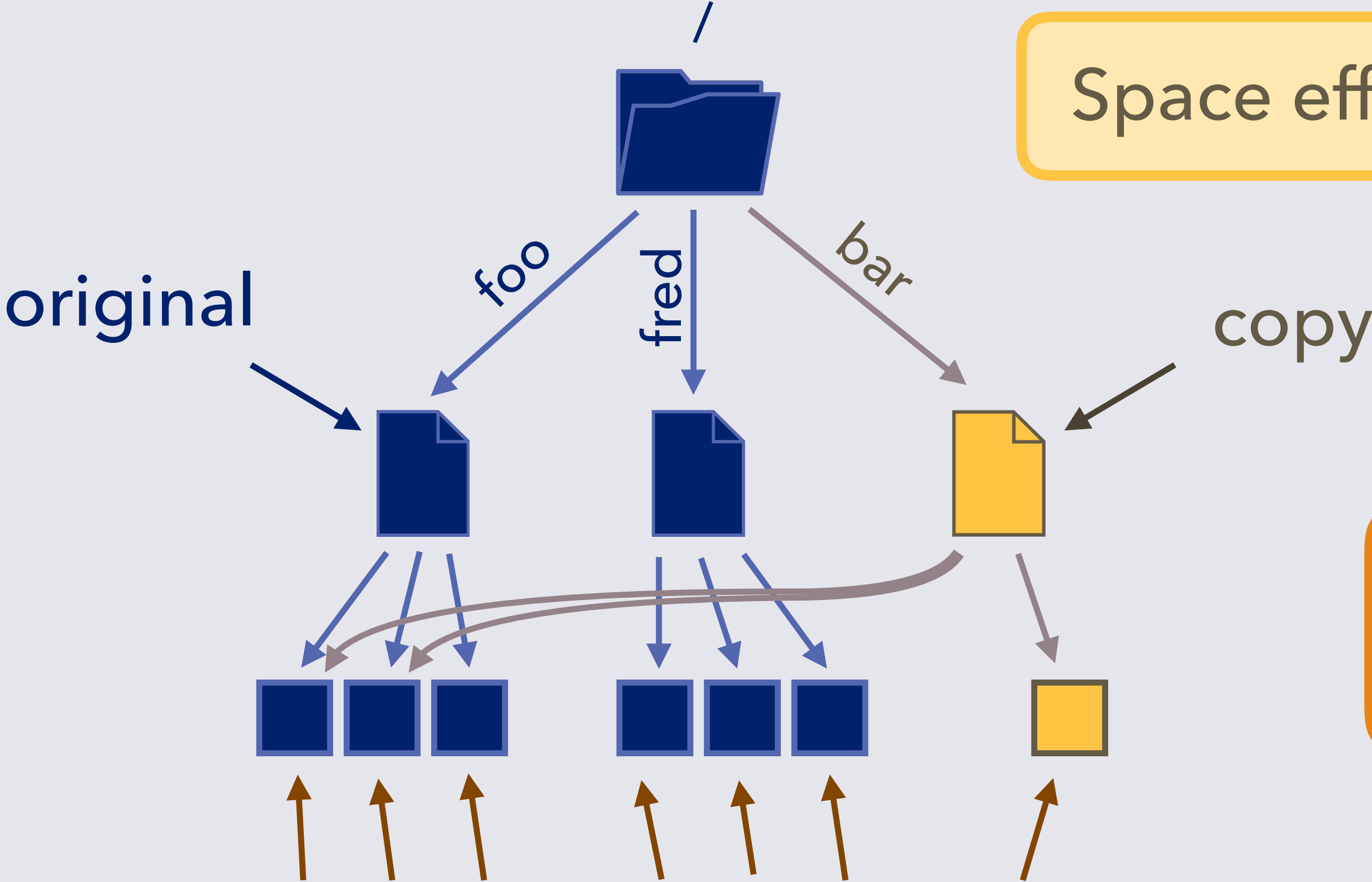
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ✗

Fast Writes ?



only have locality if the blocks are large

no locality guarantees between data blocks

Logical Copy in an Inode File System

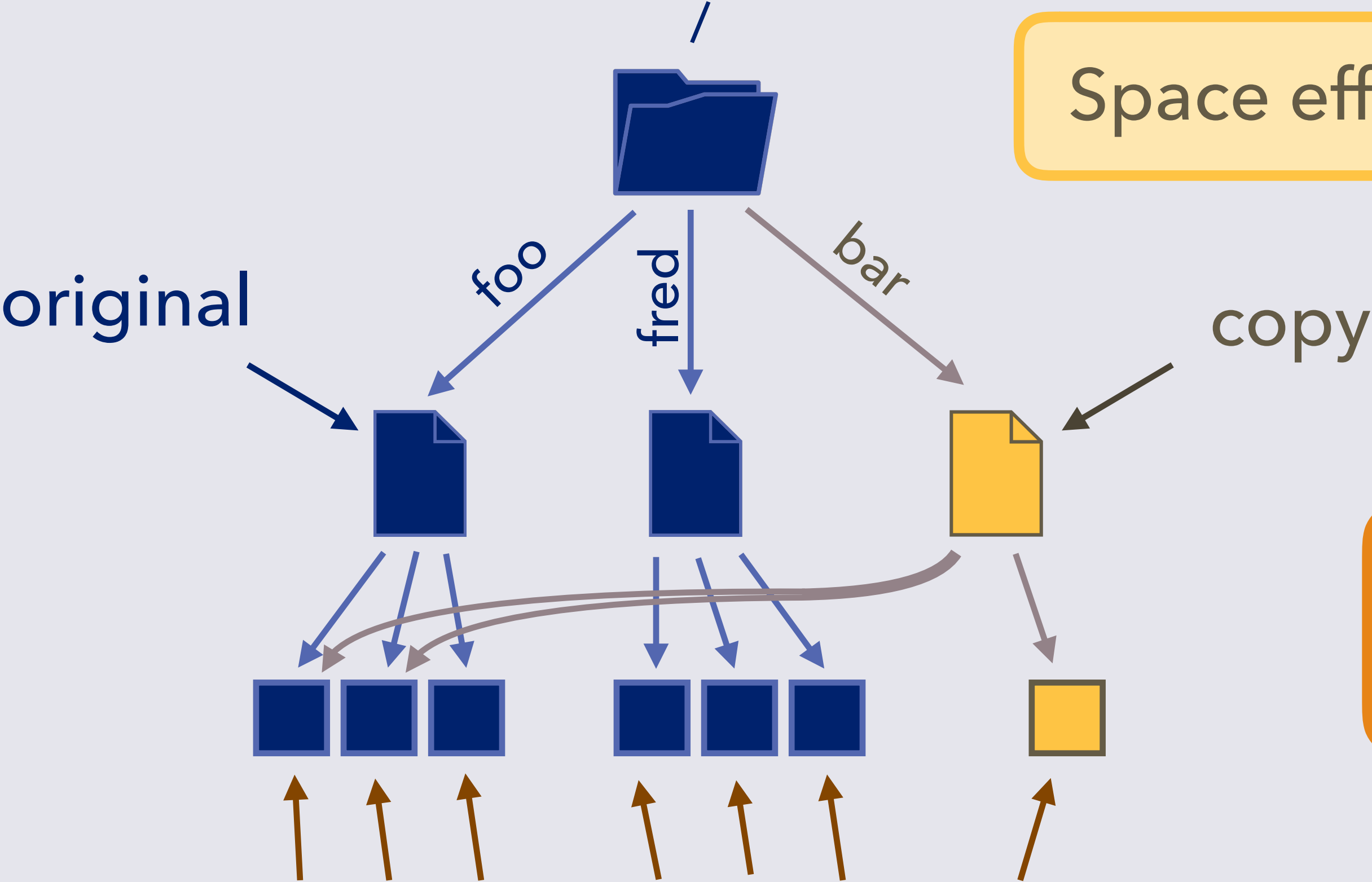
Copy /foo to /bar

Low latency ✓

Fast Reads ?

Space efficient ✗

Fast Writes ?



usually 4KiB \Rightarrow too small for locality

only have locality if the blocks are large

no locality guarantees between data blocks

Logical Copy in an Inode File System

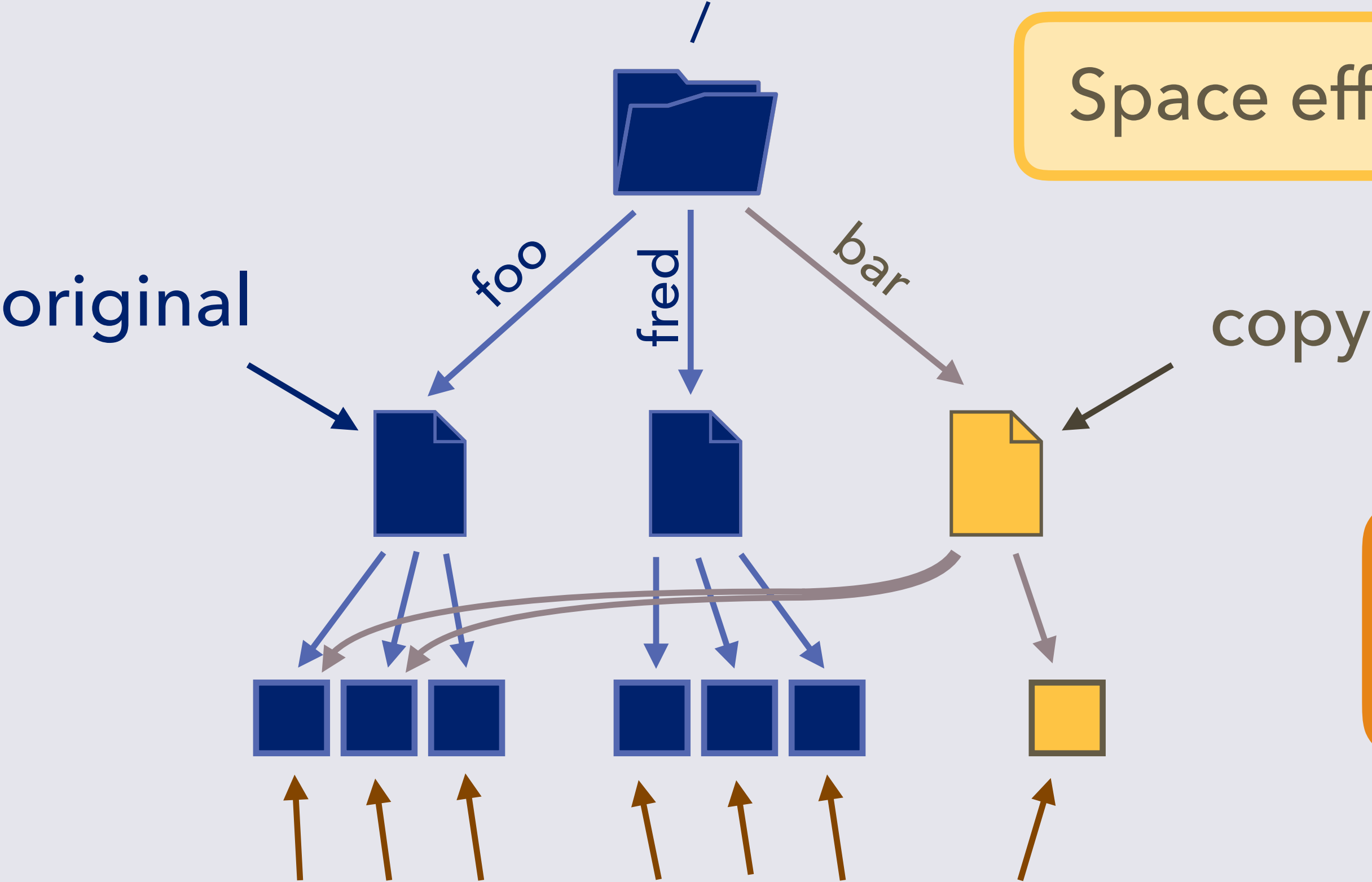
Copy /foo to /bar

Low latency ✓

Fast Reads ✗

Space efficient ✗

Fast Writes ✗



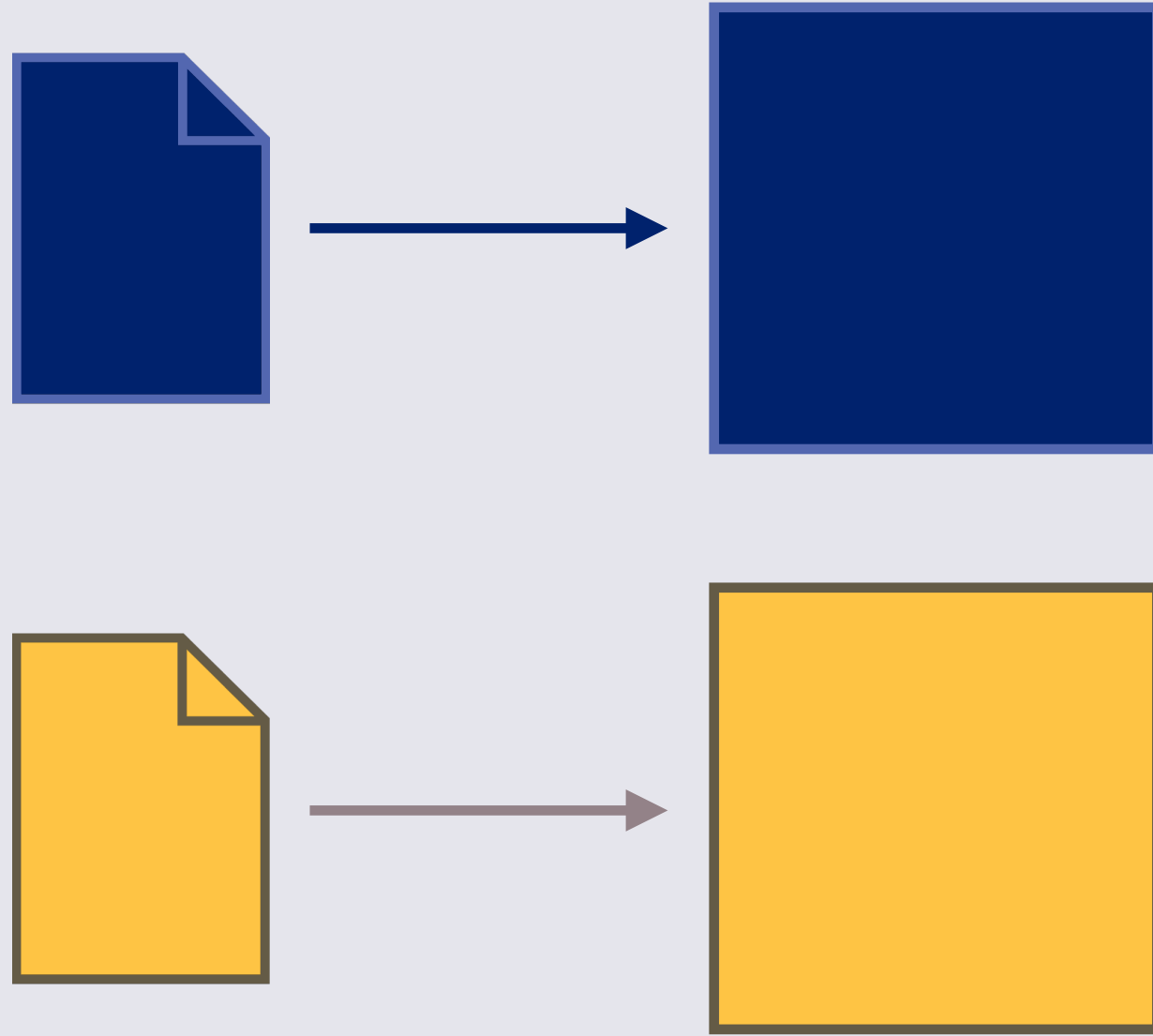
usually 4KiB \Rightarrow too small for locality

only have locality if the blocks are large

no locality guarantees between data blocks

Space-Locality Tradeoff

Larger blocks



Better locality

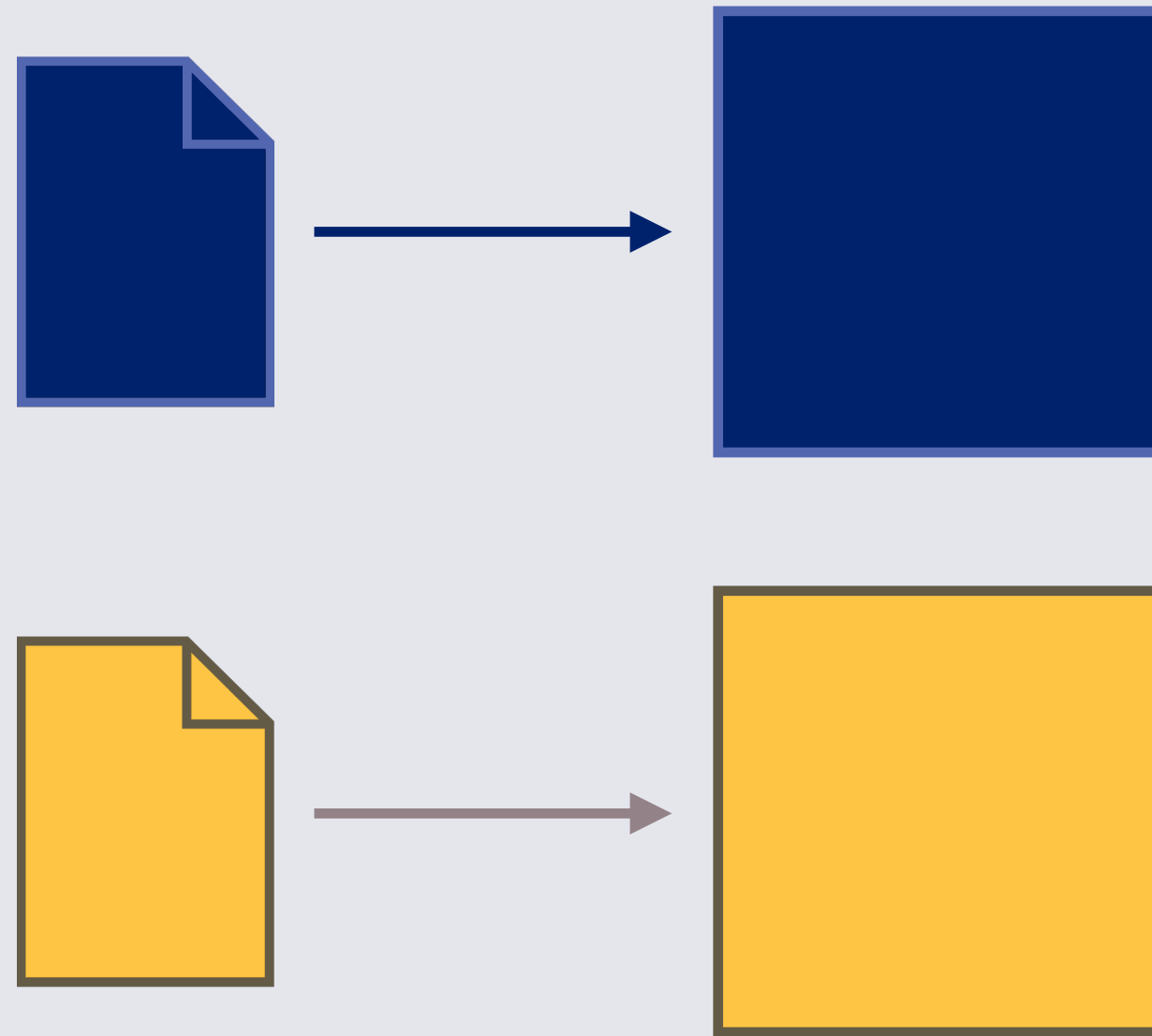


Worse space efficiency



Space-Locality Tradeoff

Larger blocks



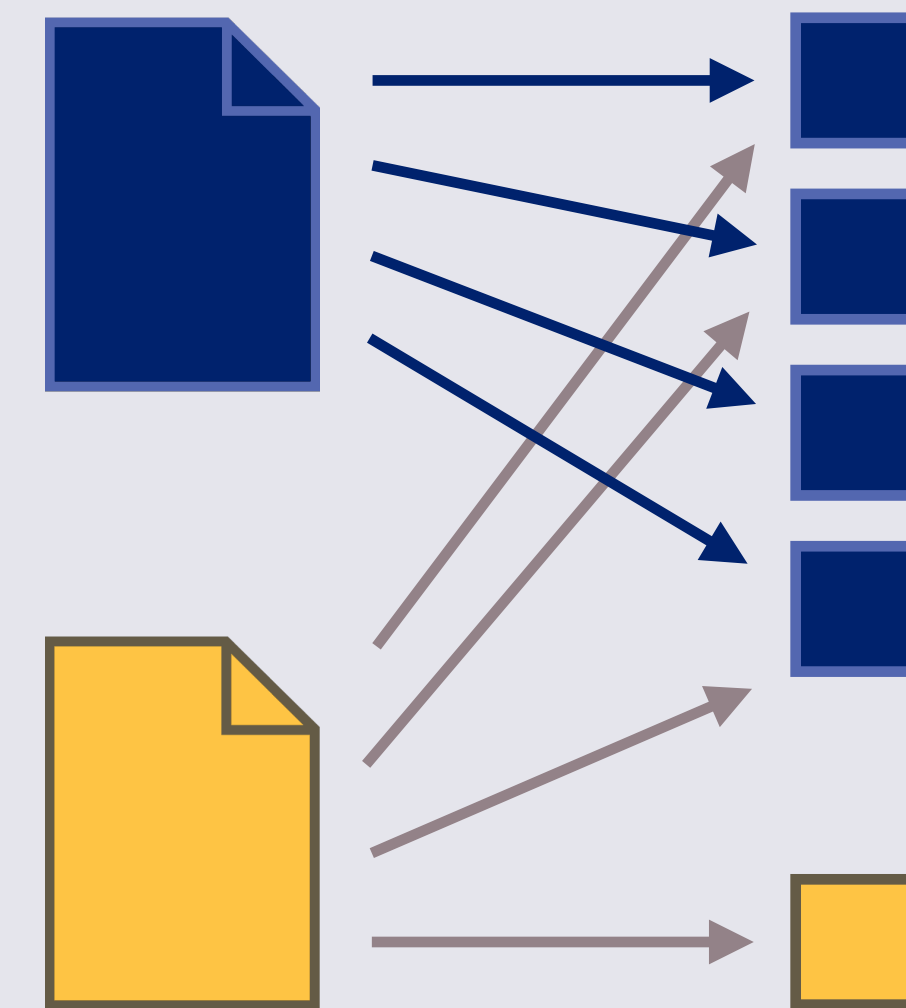
Better locality



Worse space efficiency



Smaller blocks



Worse locality



Better space efficiency

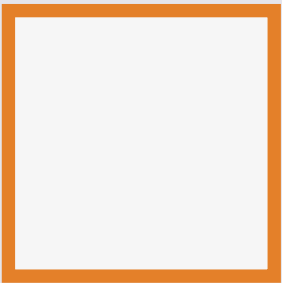
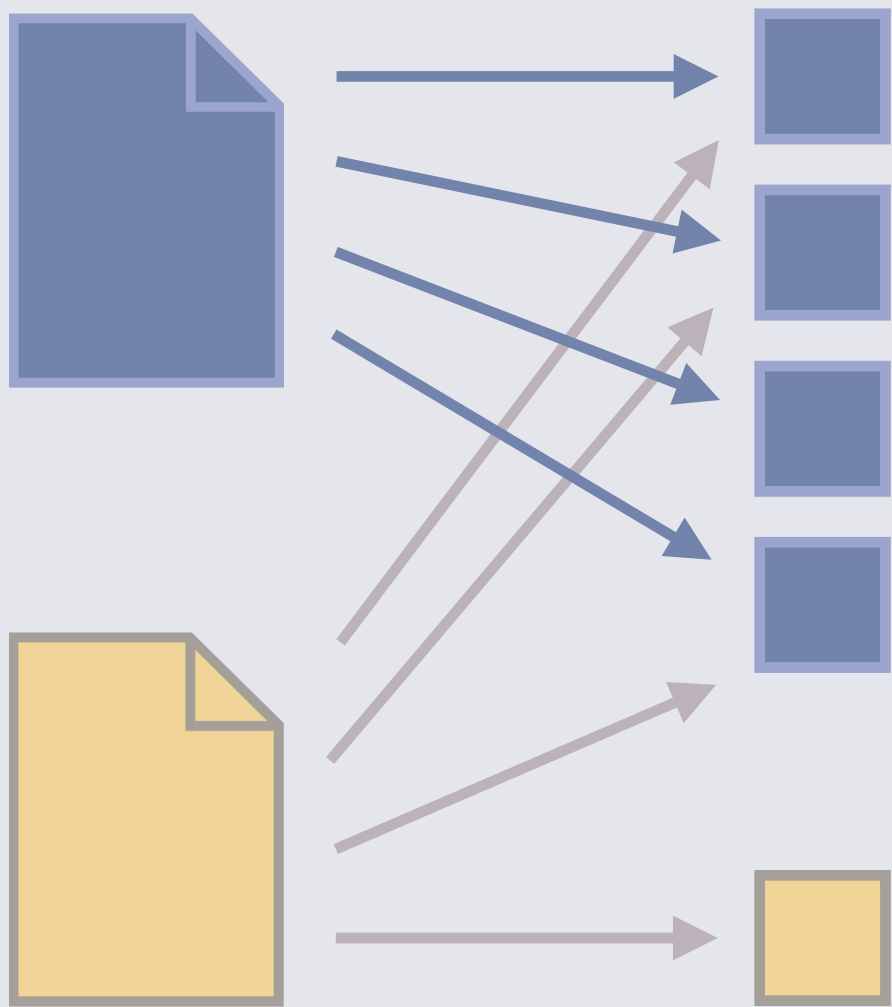
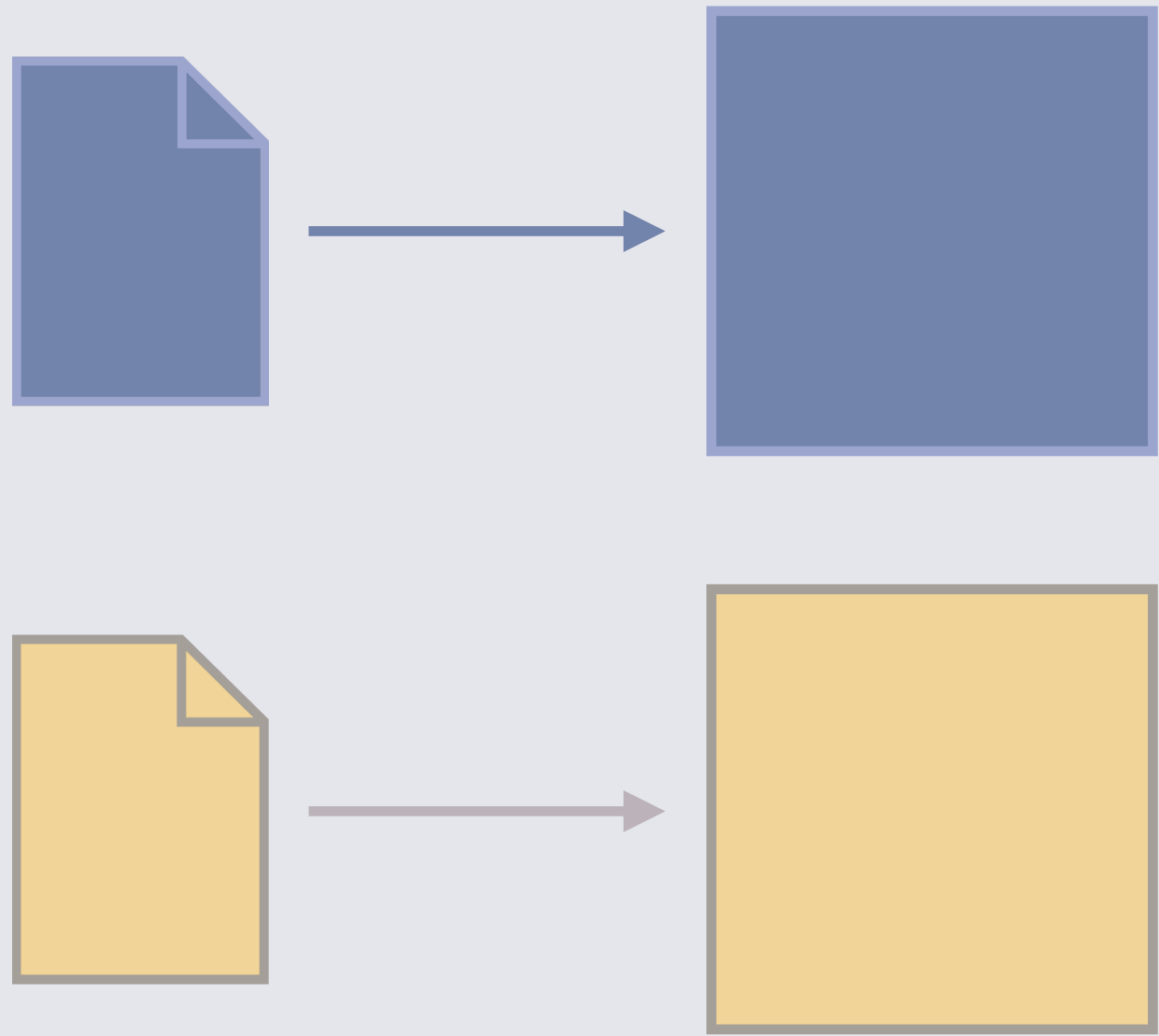


Space-Locality Tradeoff

Larger blocks

Smaller blocks

4KiB blocks



Too large for space Too small for locality



Better locality



Worse locality



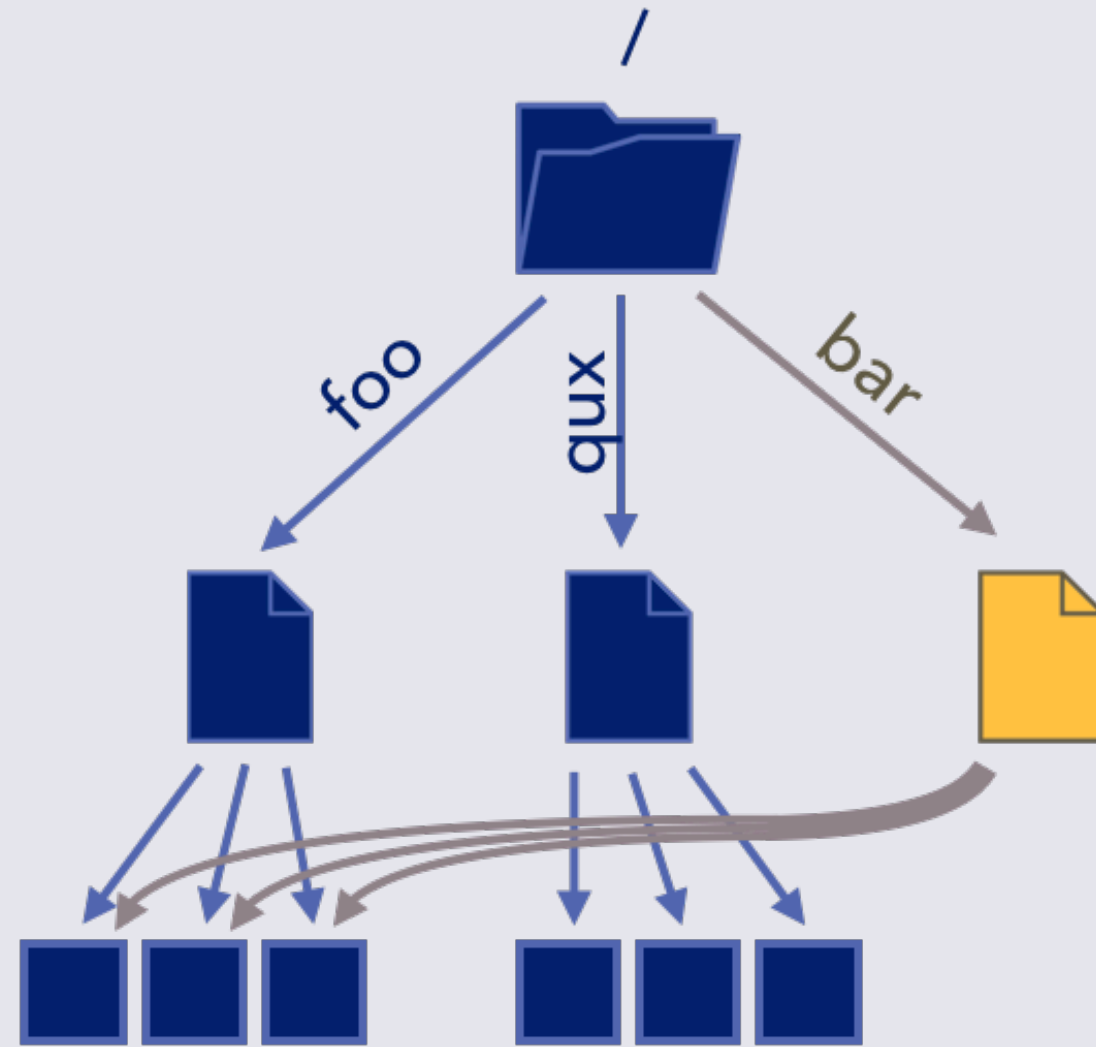
Worse space efficiency



Better space efficiency

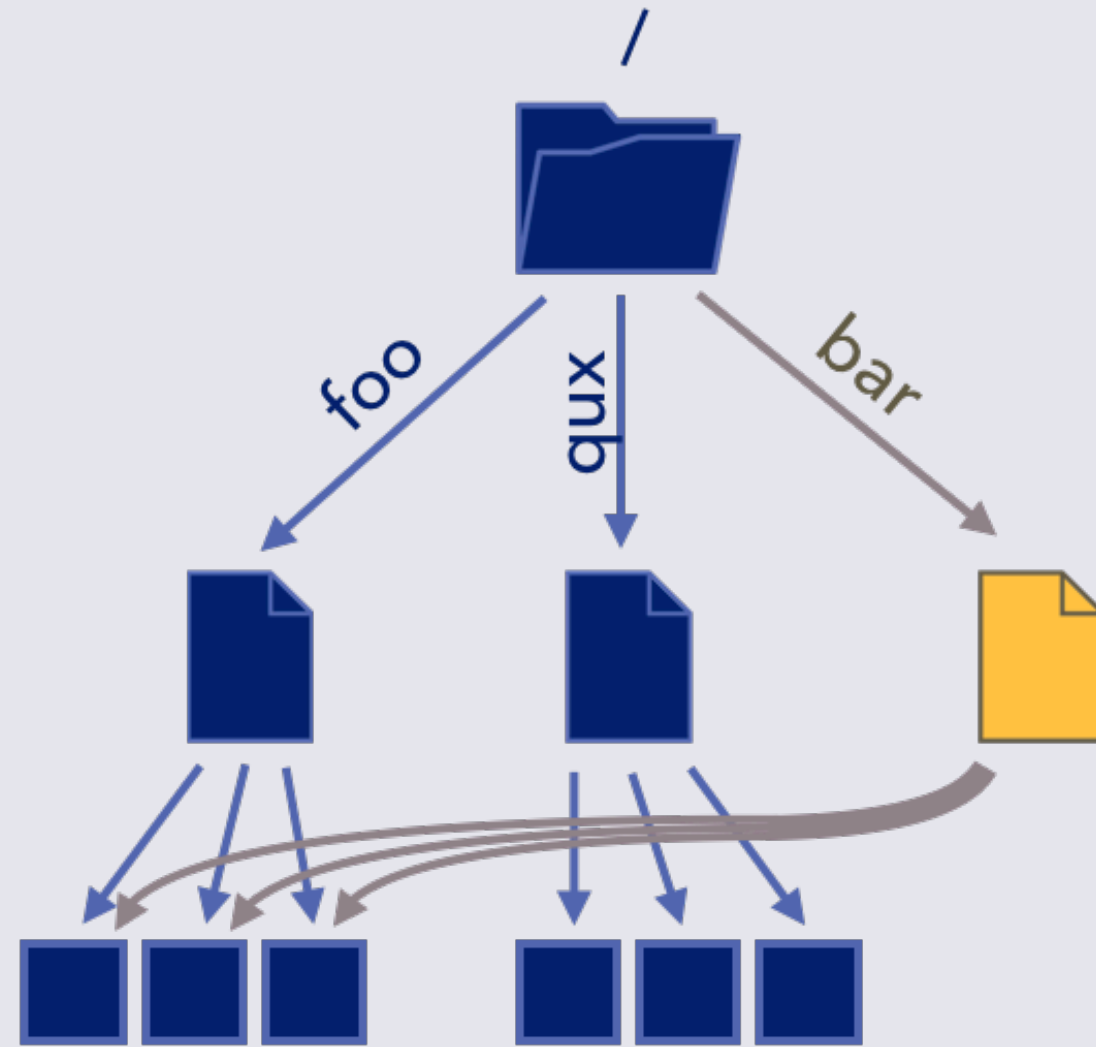


Inode Logical Copy Takeaway

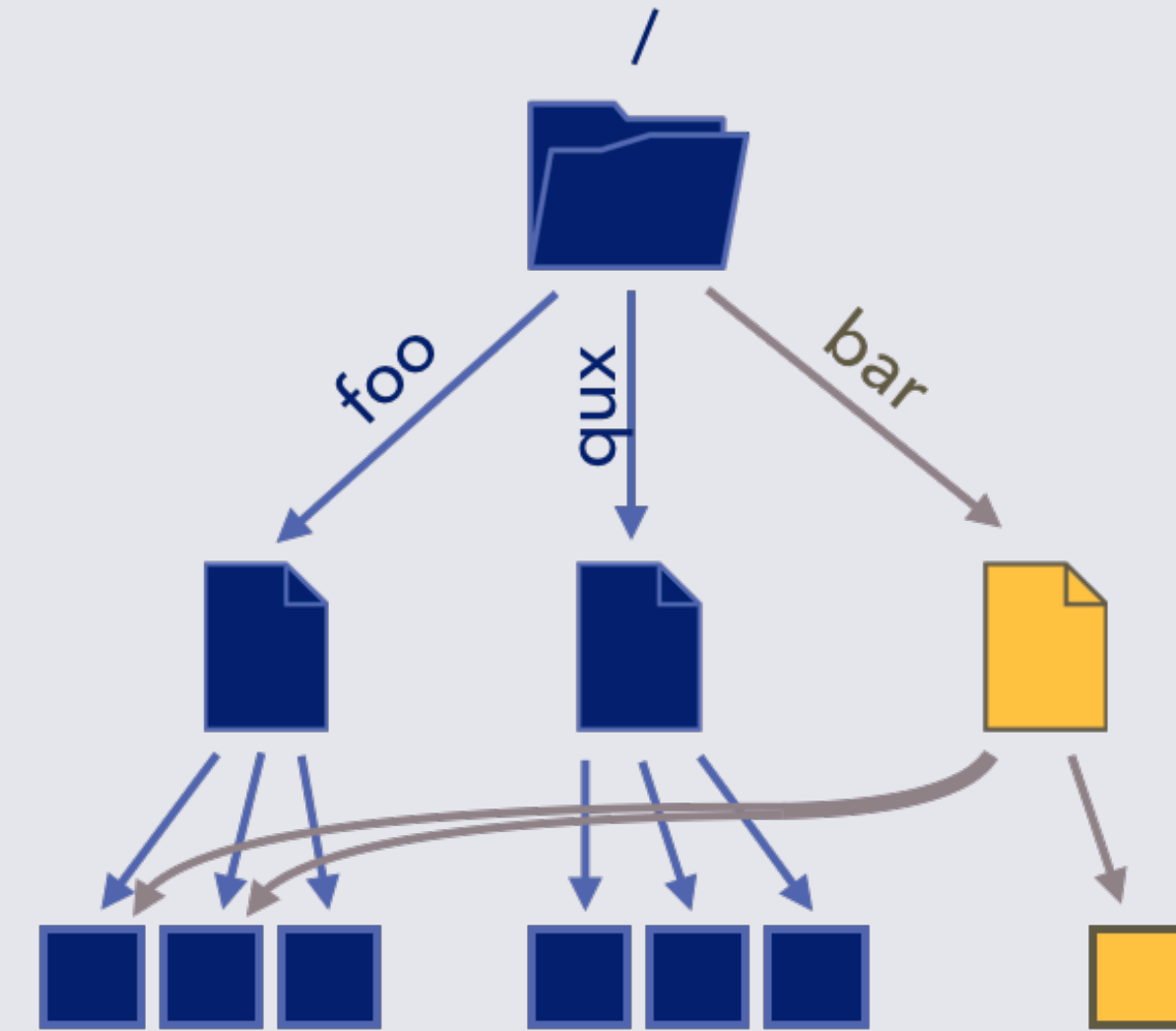


Using a DAG to share data is great for latency

Inode Logical Copy Takeaway



Using a DAG to share data is great for latency

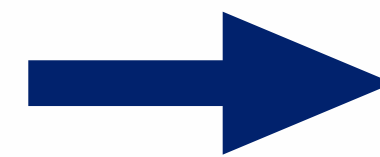


Challenge:
Small writes break sharing

Our Solution: B^ϵ -DAGs

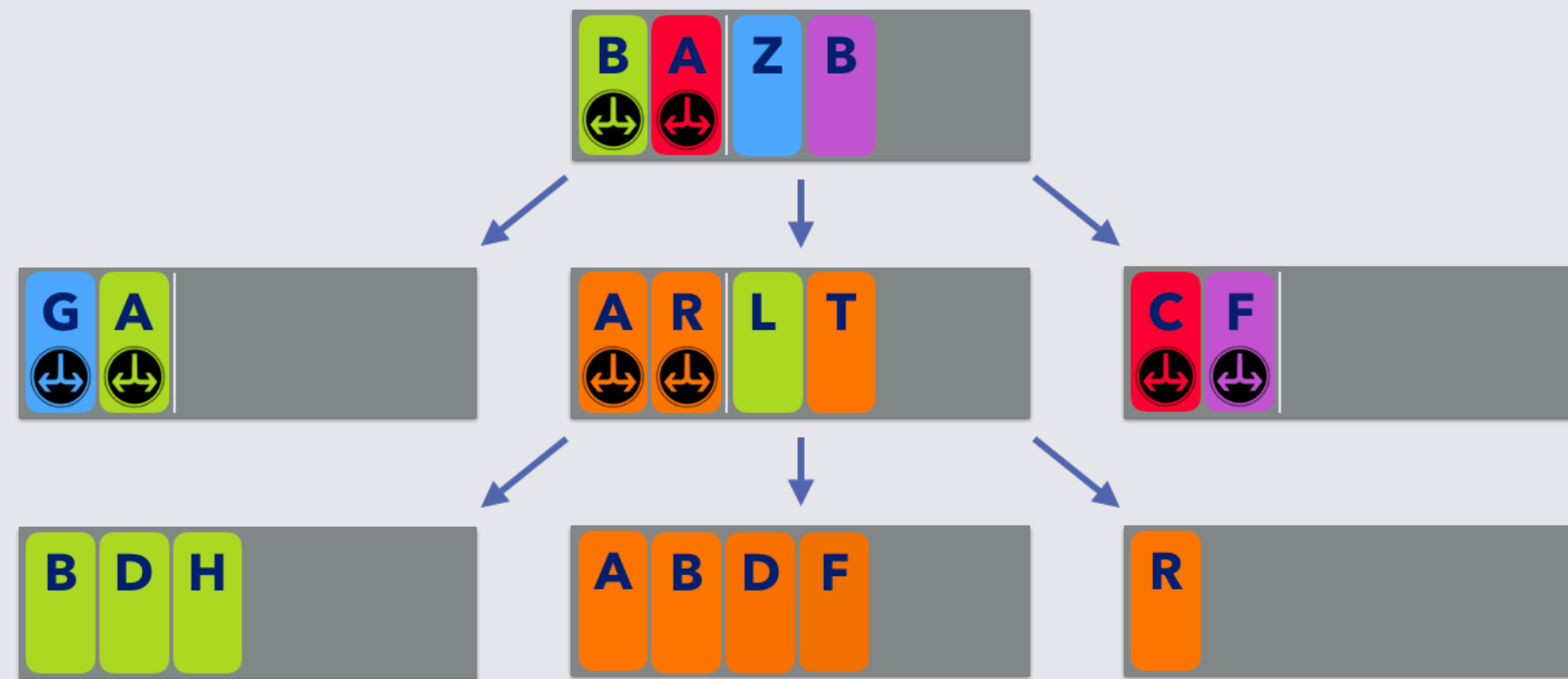
B^ϵ trFS

B^ϵ -Trees



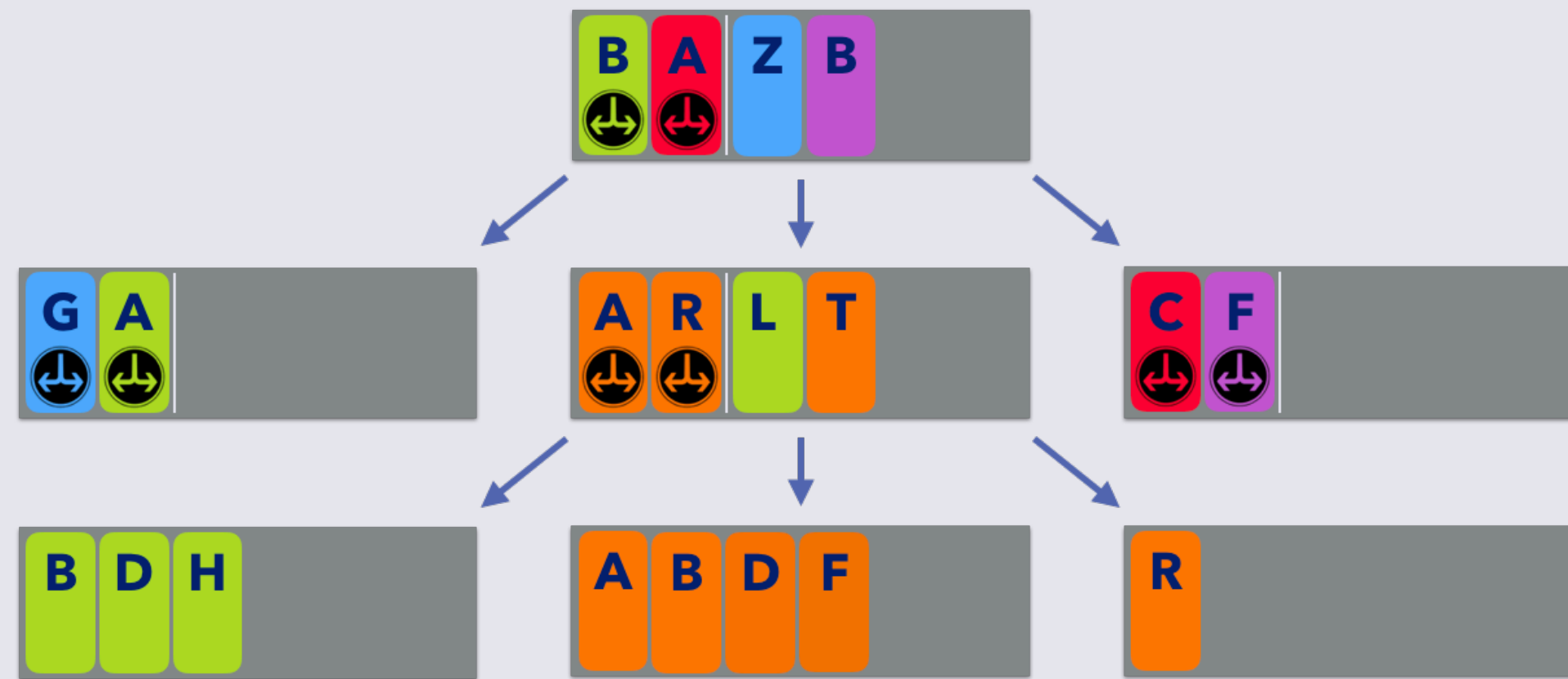
B^ϵ -DAGs

Our Solution: B^ϵ -DAGs

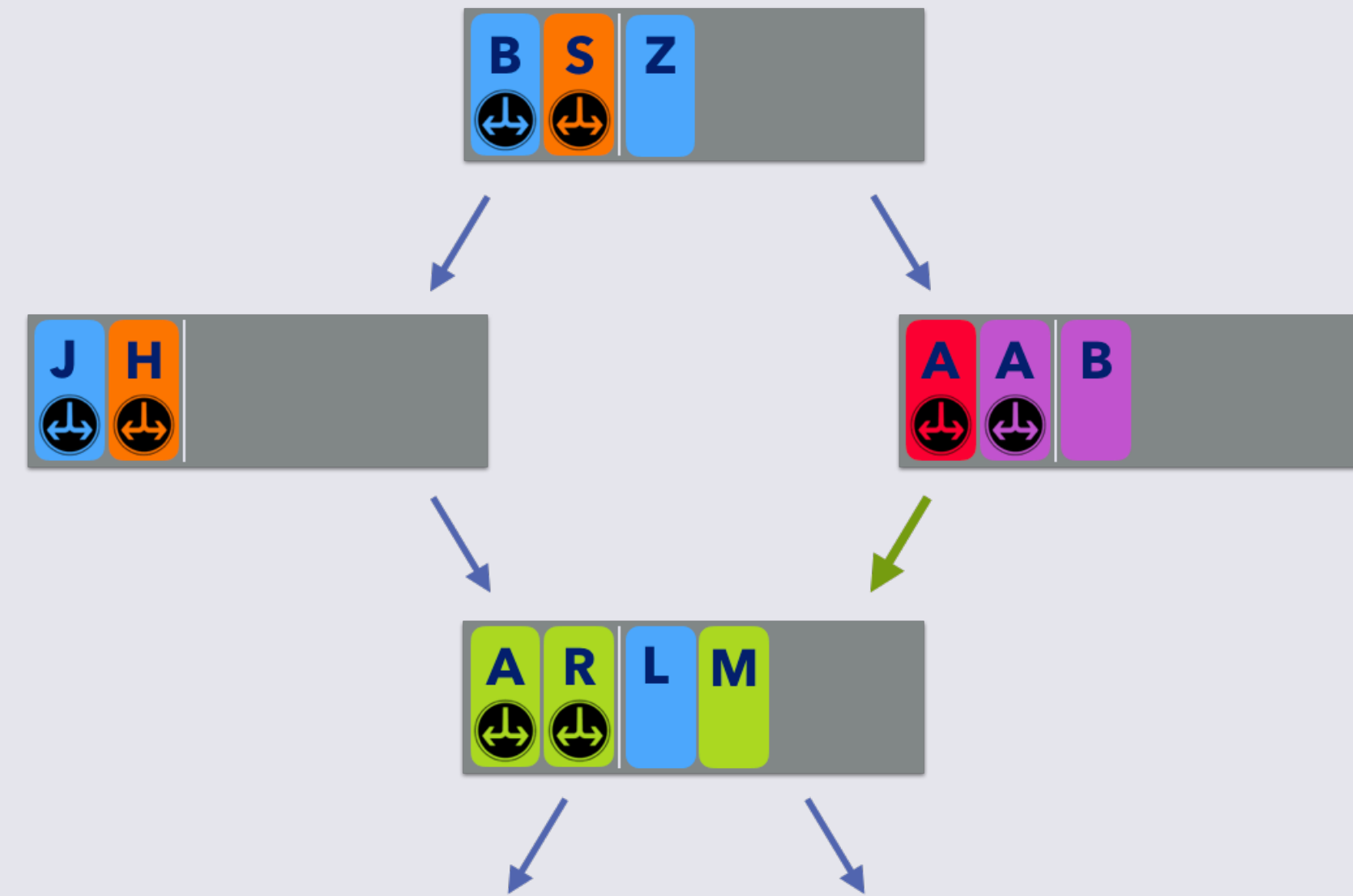


B^ϵ -trees have good locality and
batch together small writes

Our Solution: B^ϵ -DAGs



B^ϵ -trees have good locality and batch together small writes

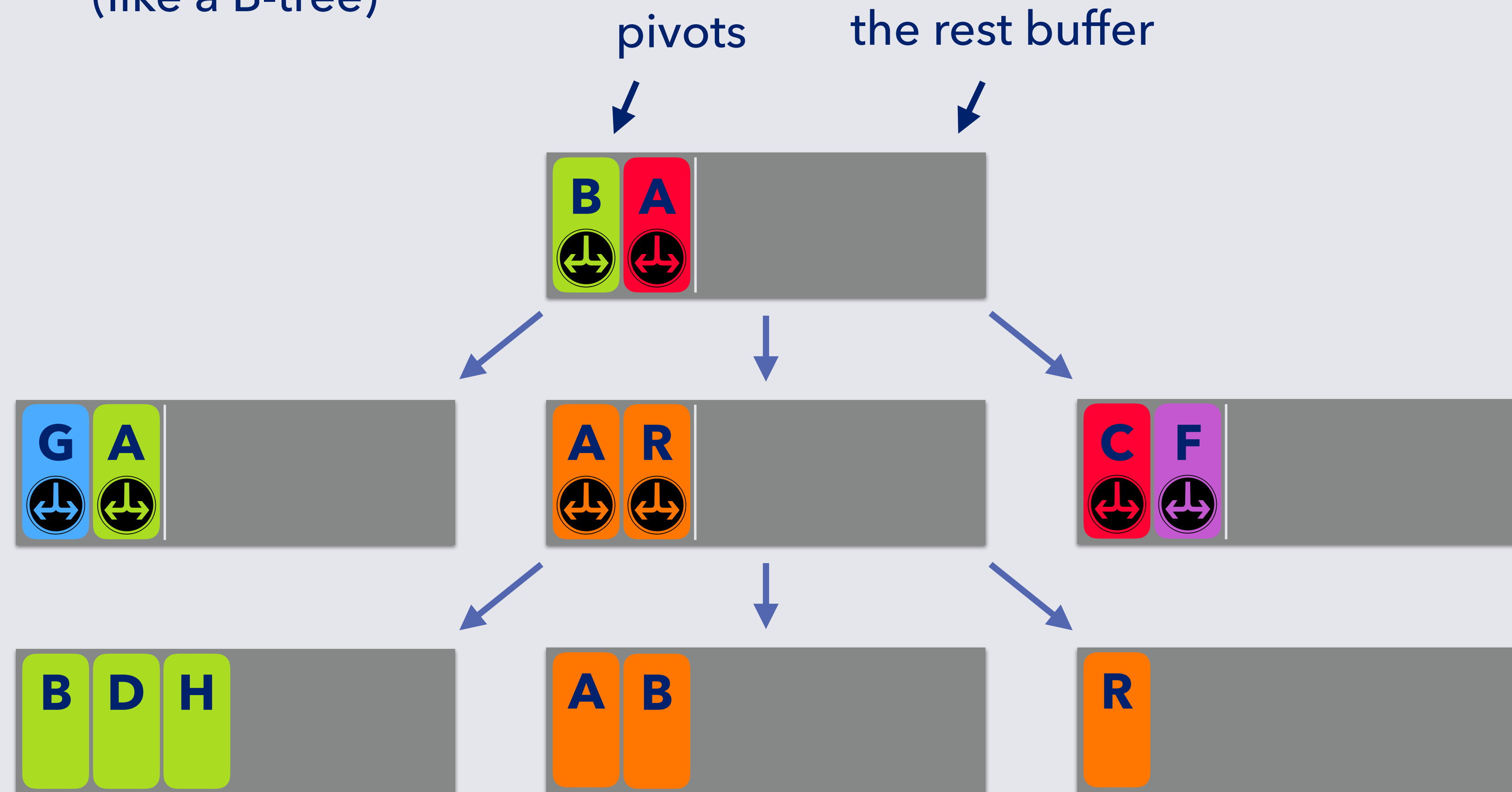


In this paper, we turn B^ϵ -trees into B^ϵ -DAGs to share data between files

B^ϵ -Trees

B ϵ -Trees

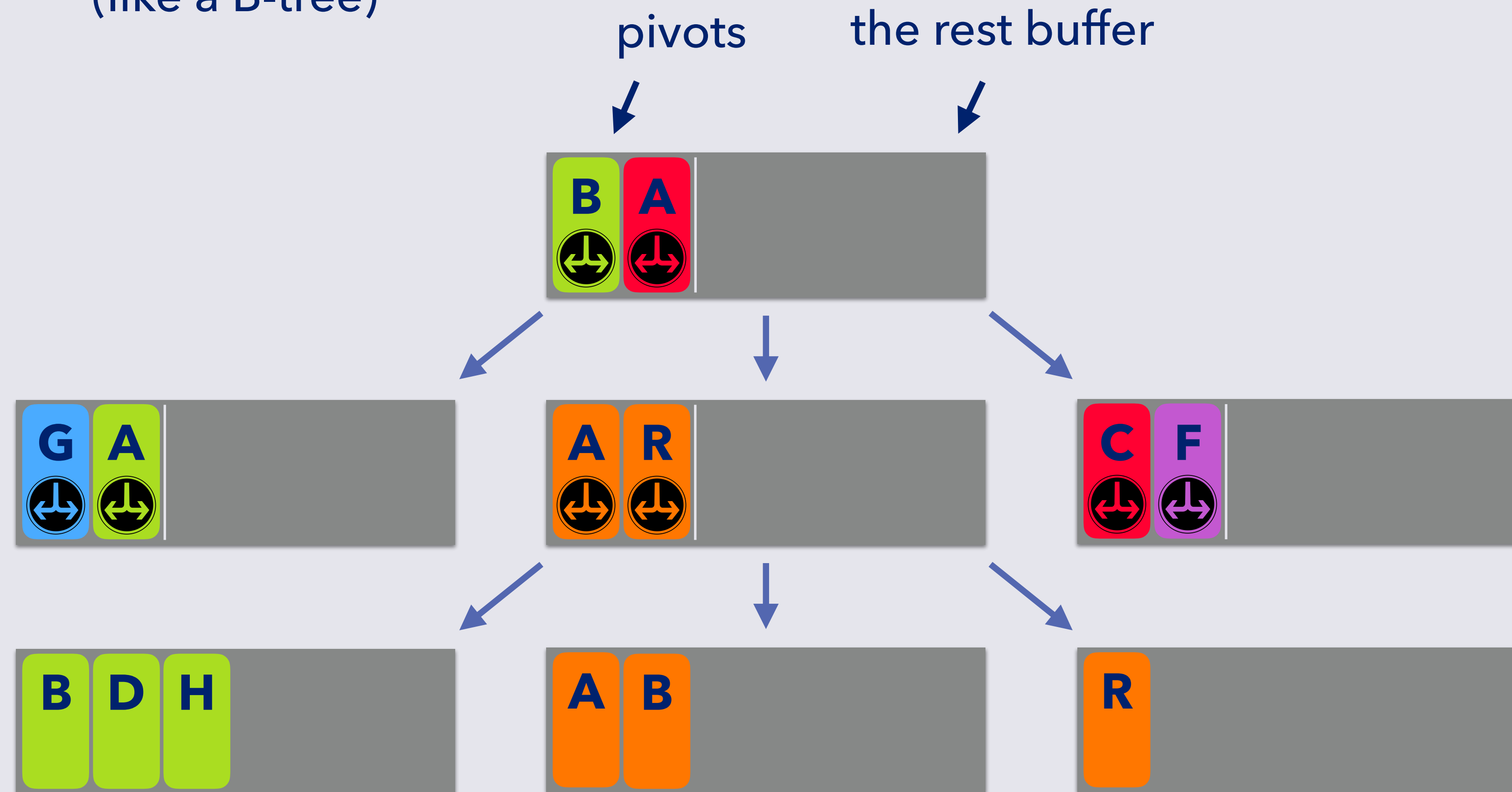
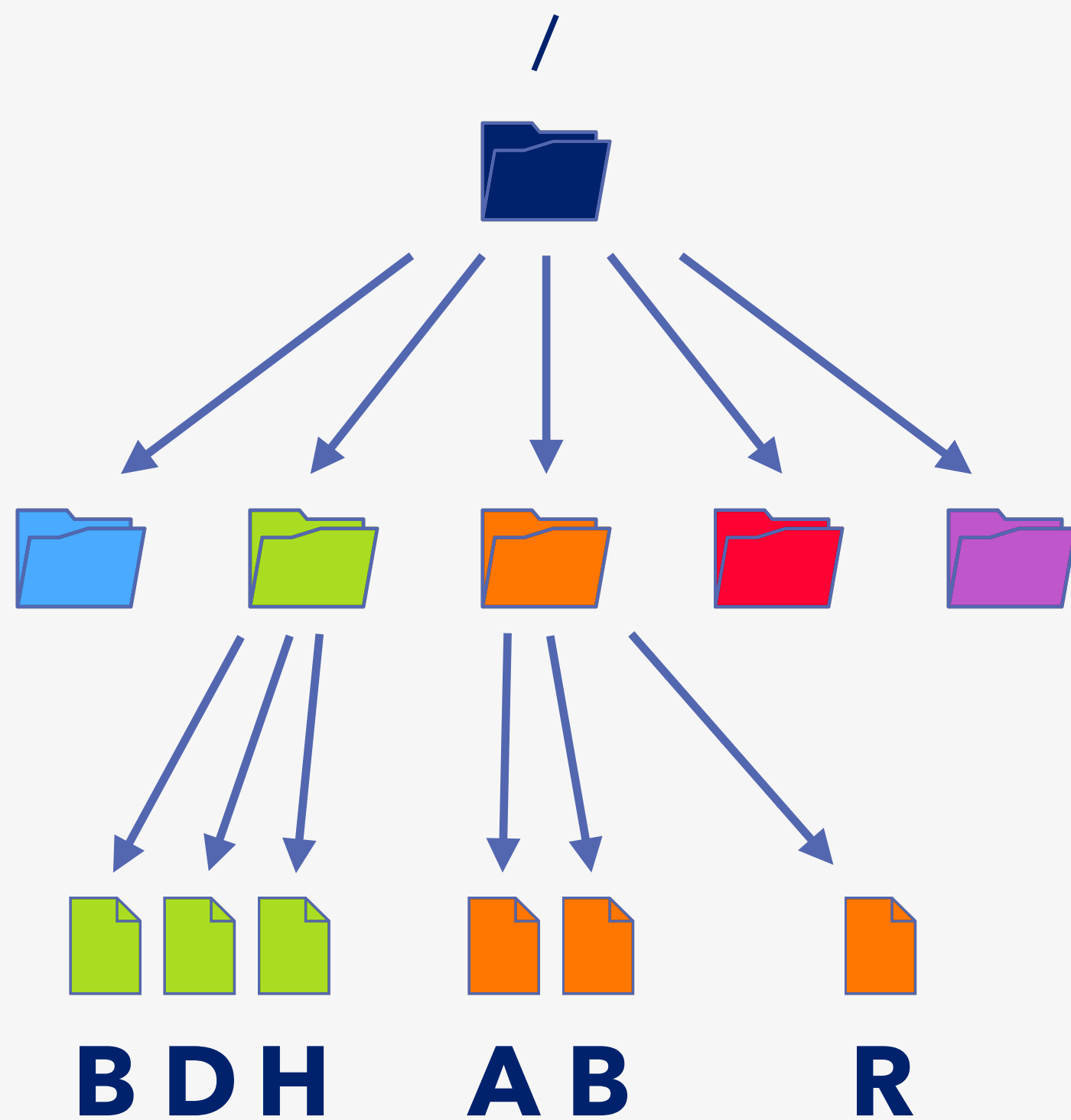
A B ϵ -tree is a search tree
(like a B-tree)



B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

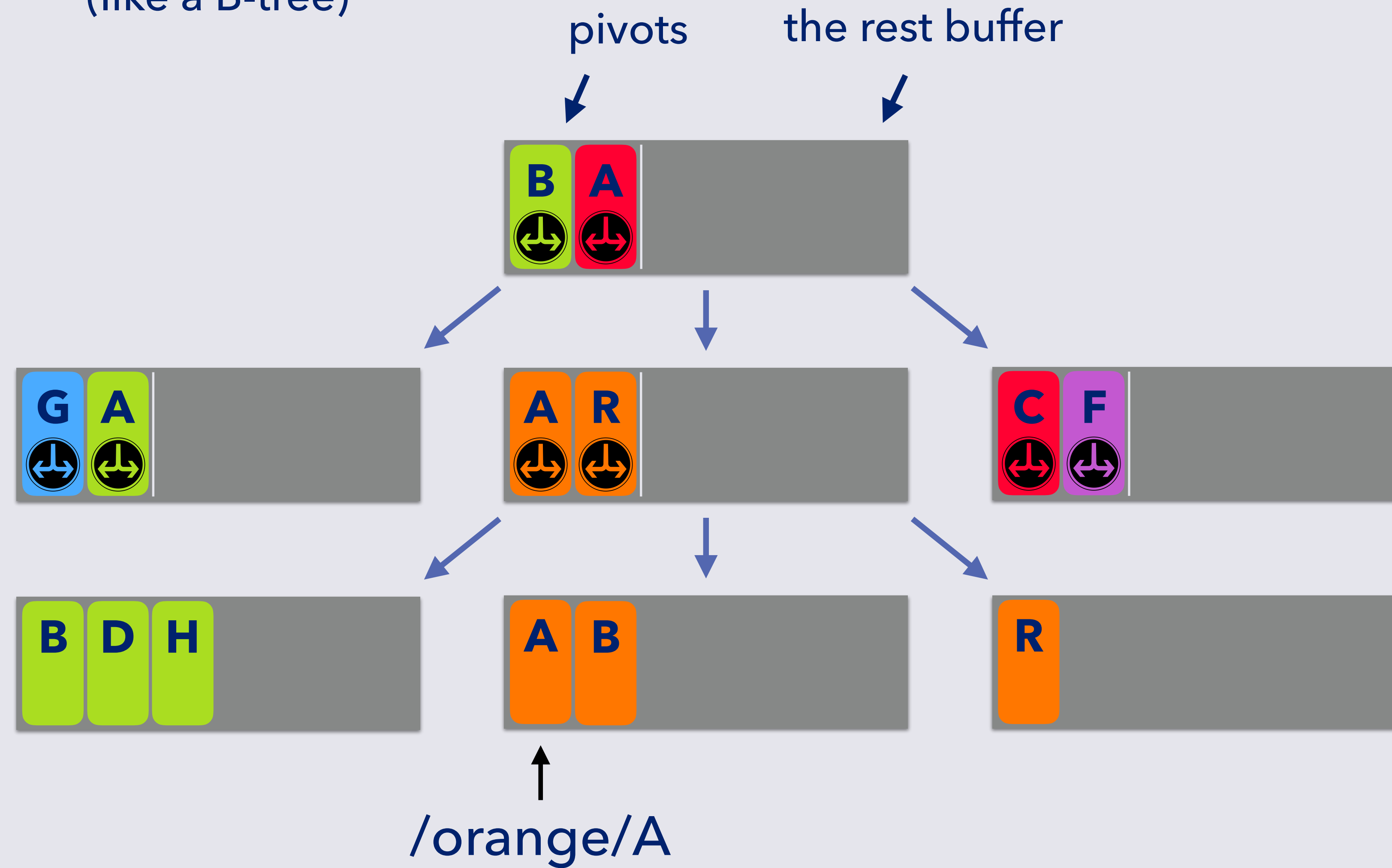
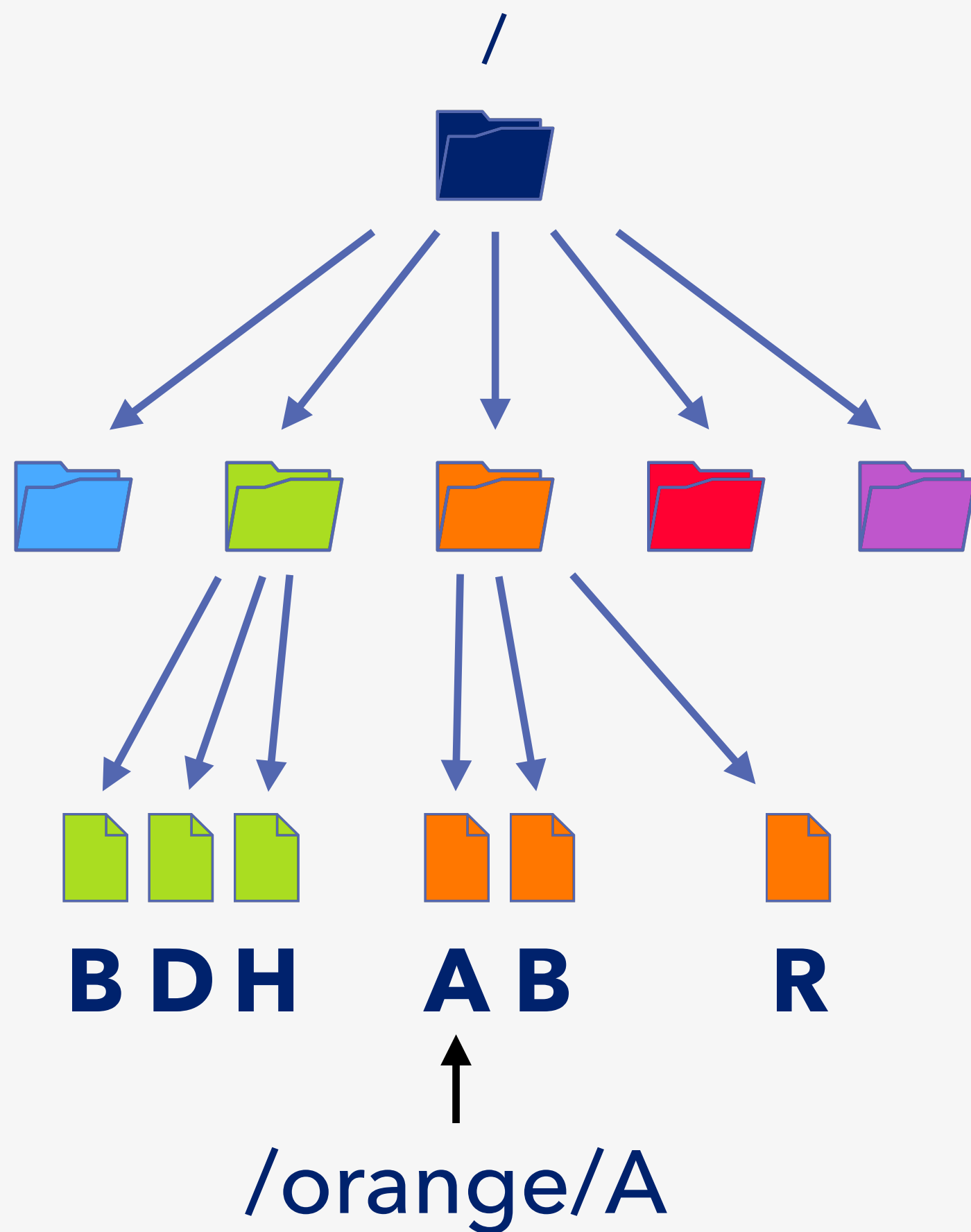
directory tree



B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

directory tree



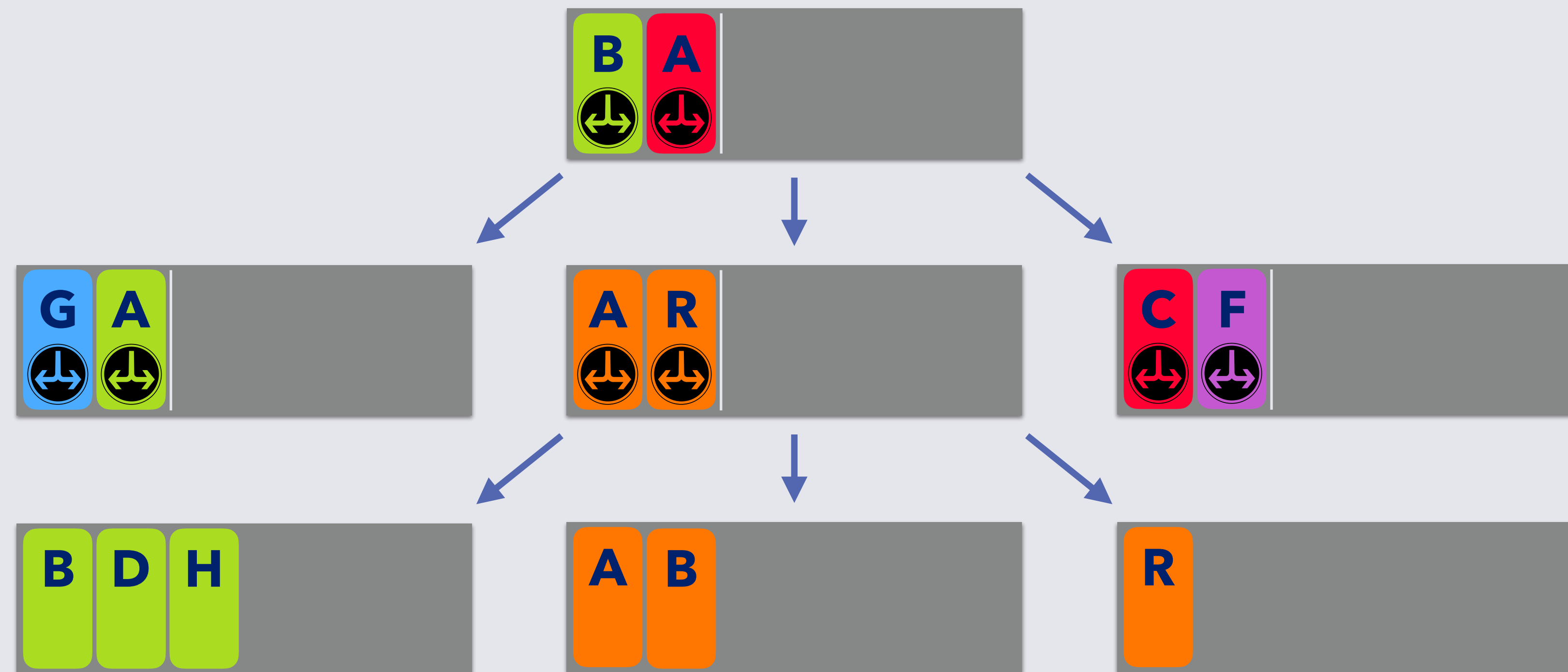
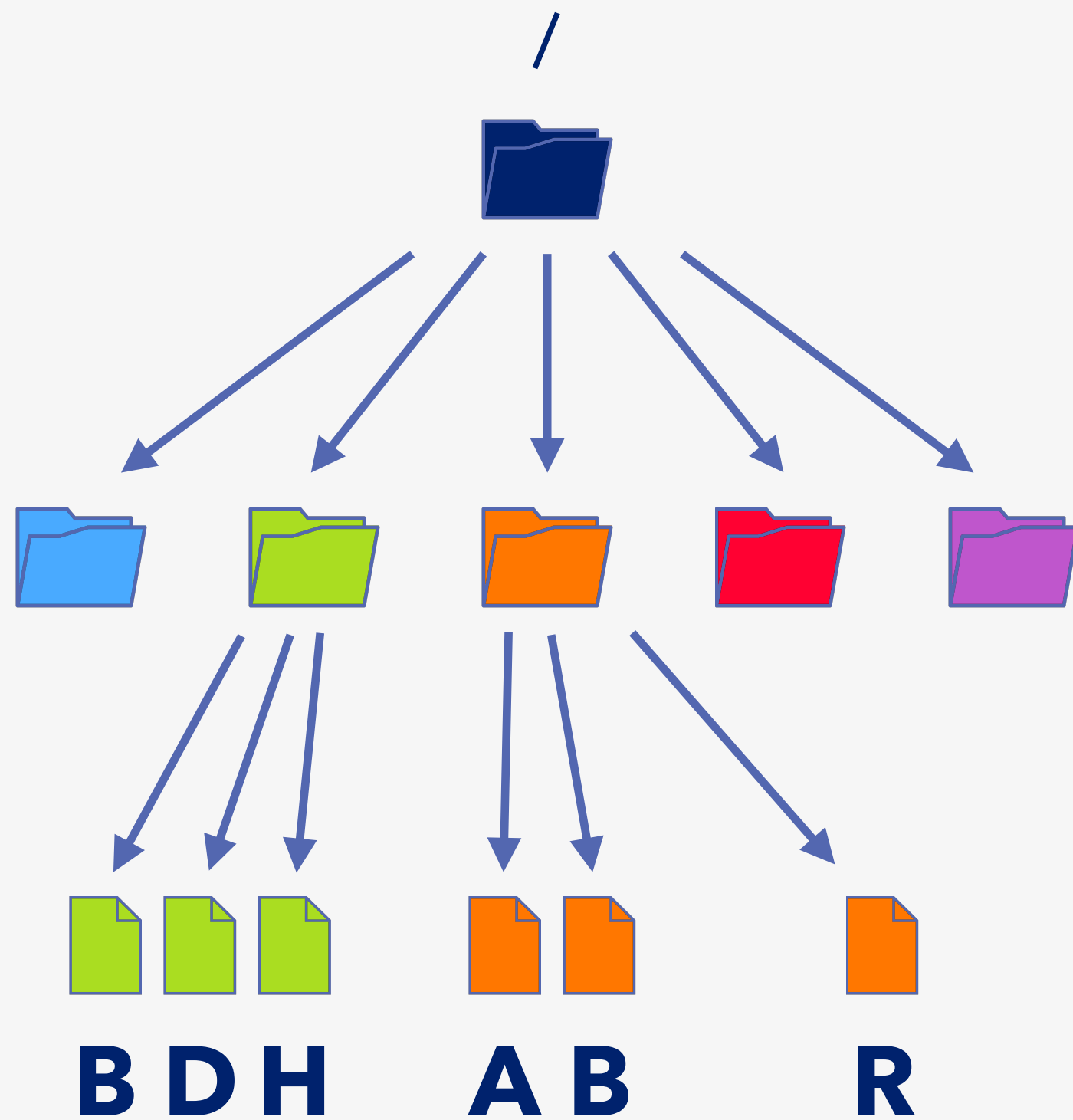
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

D New file:
/orange/D

directory tree

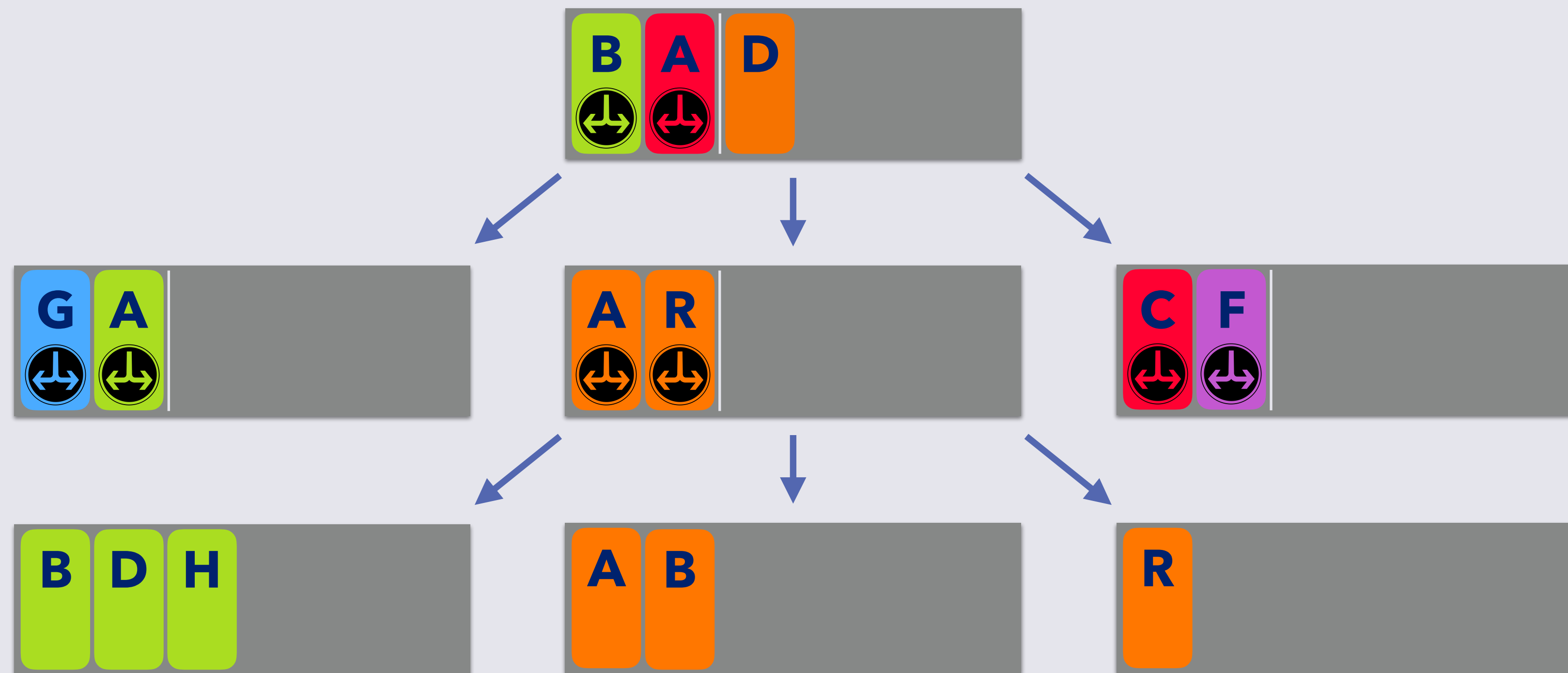
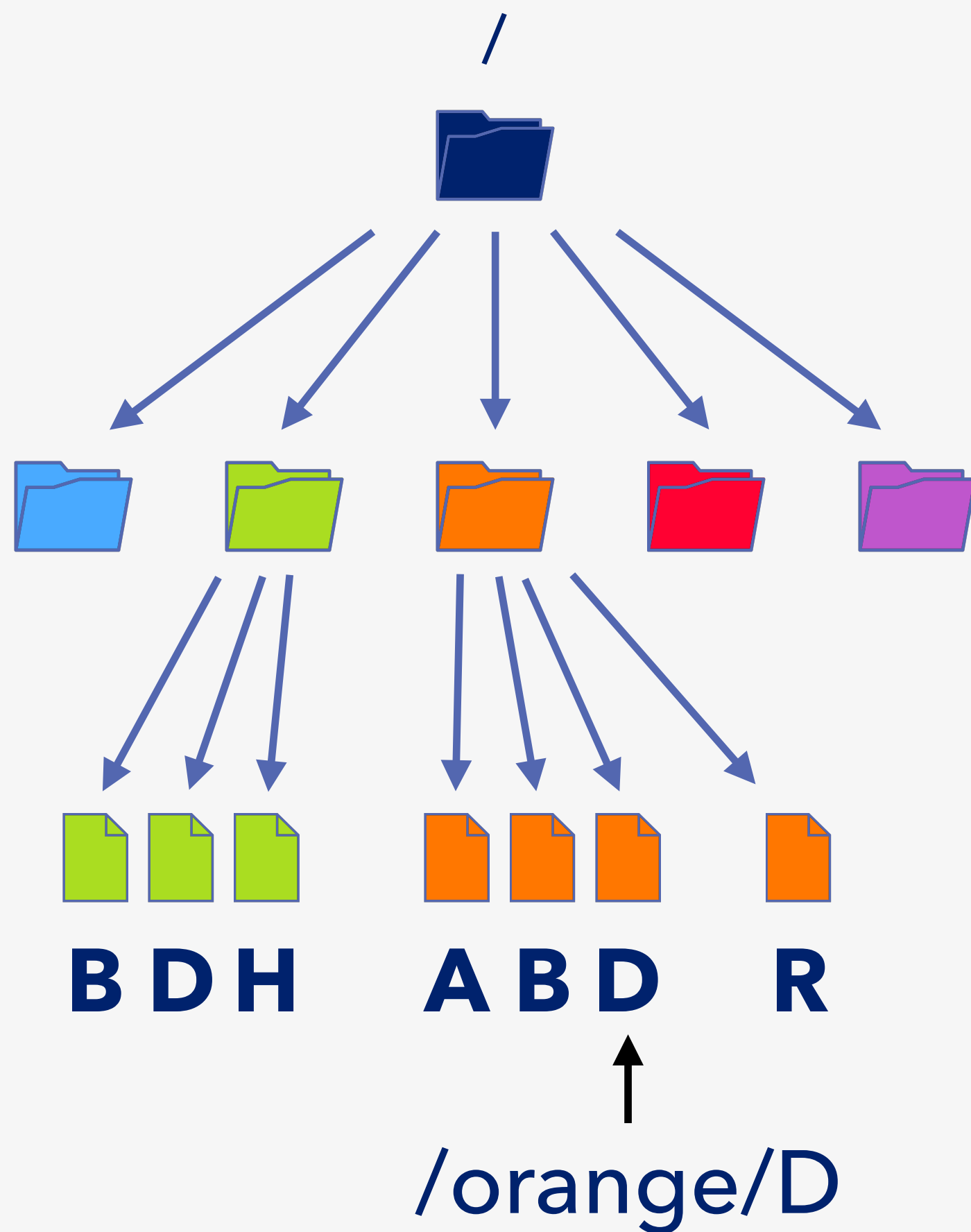


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

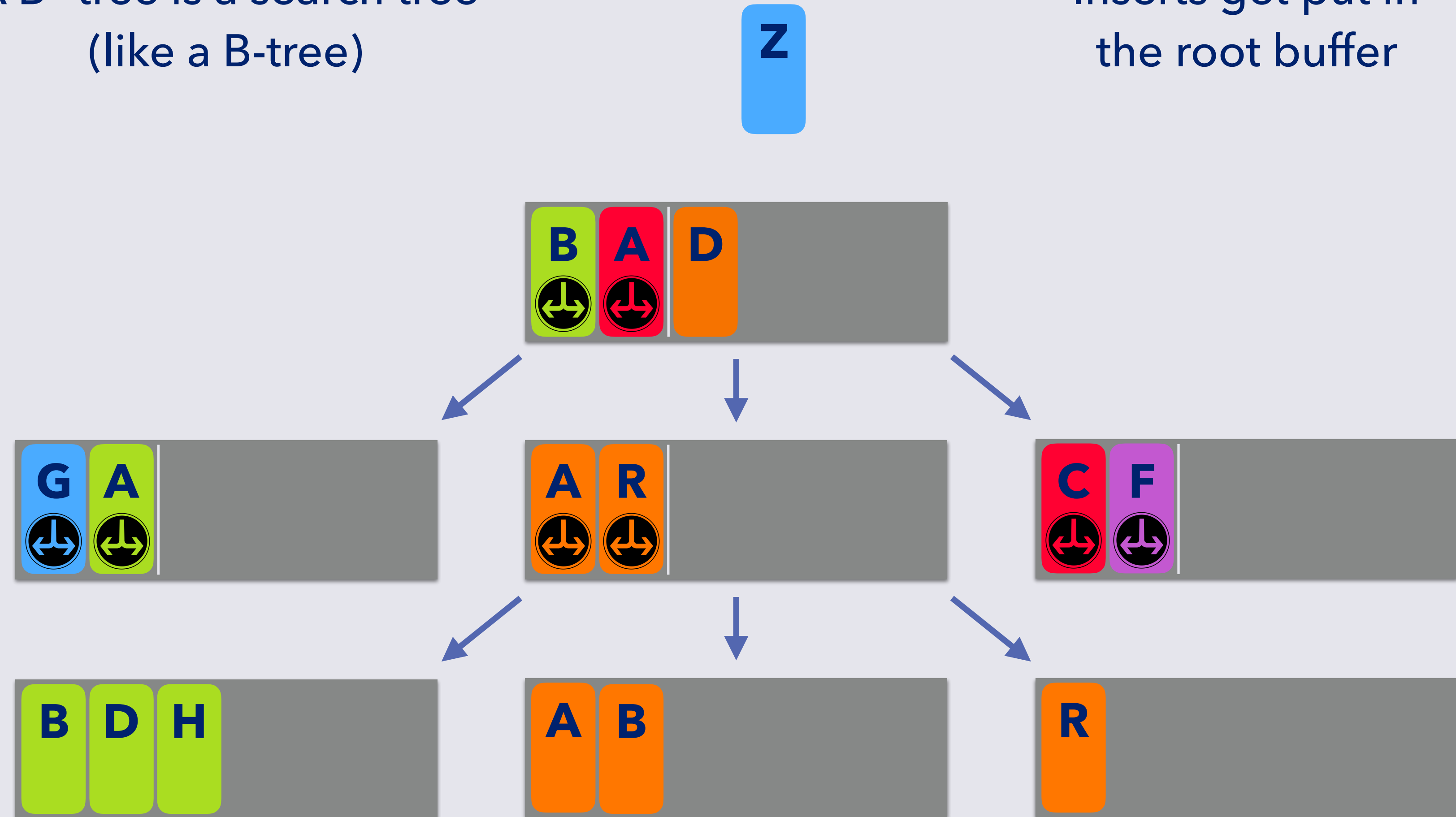
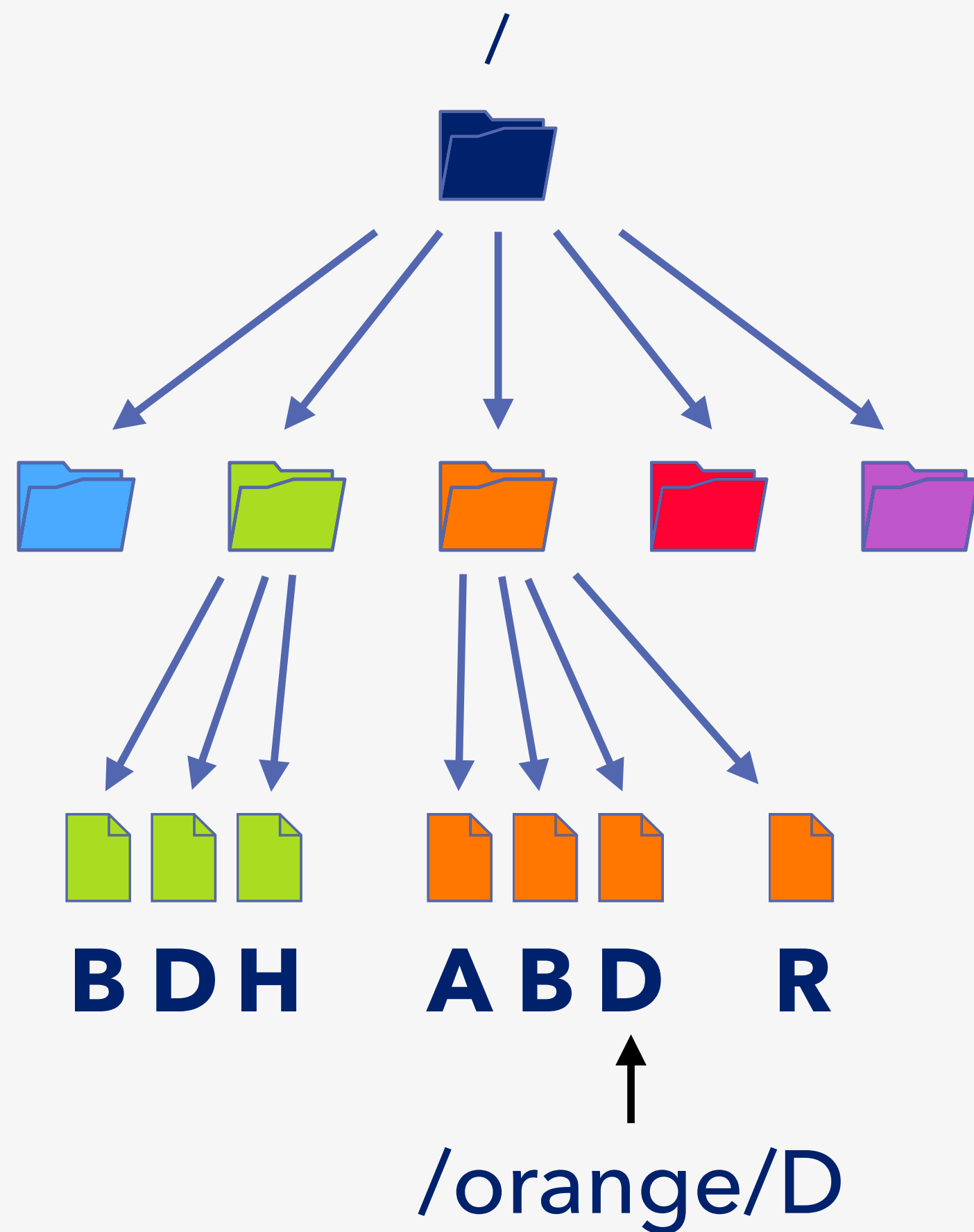


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

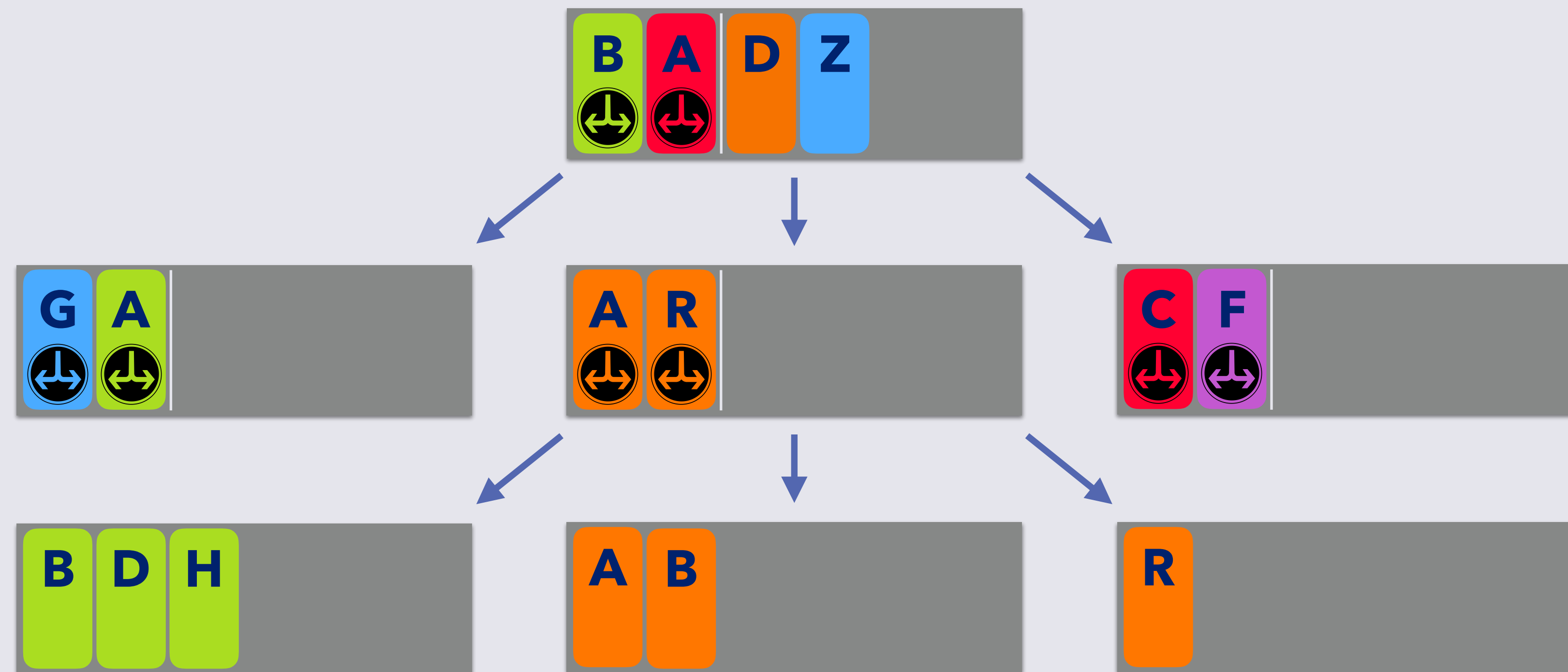
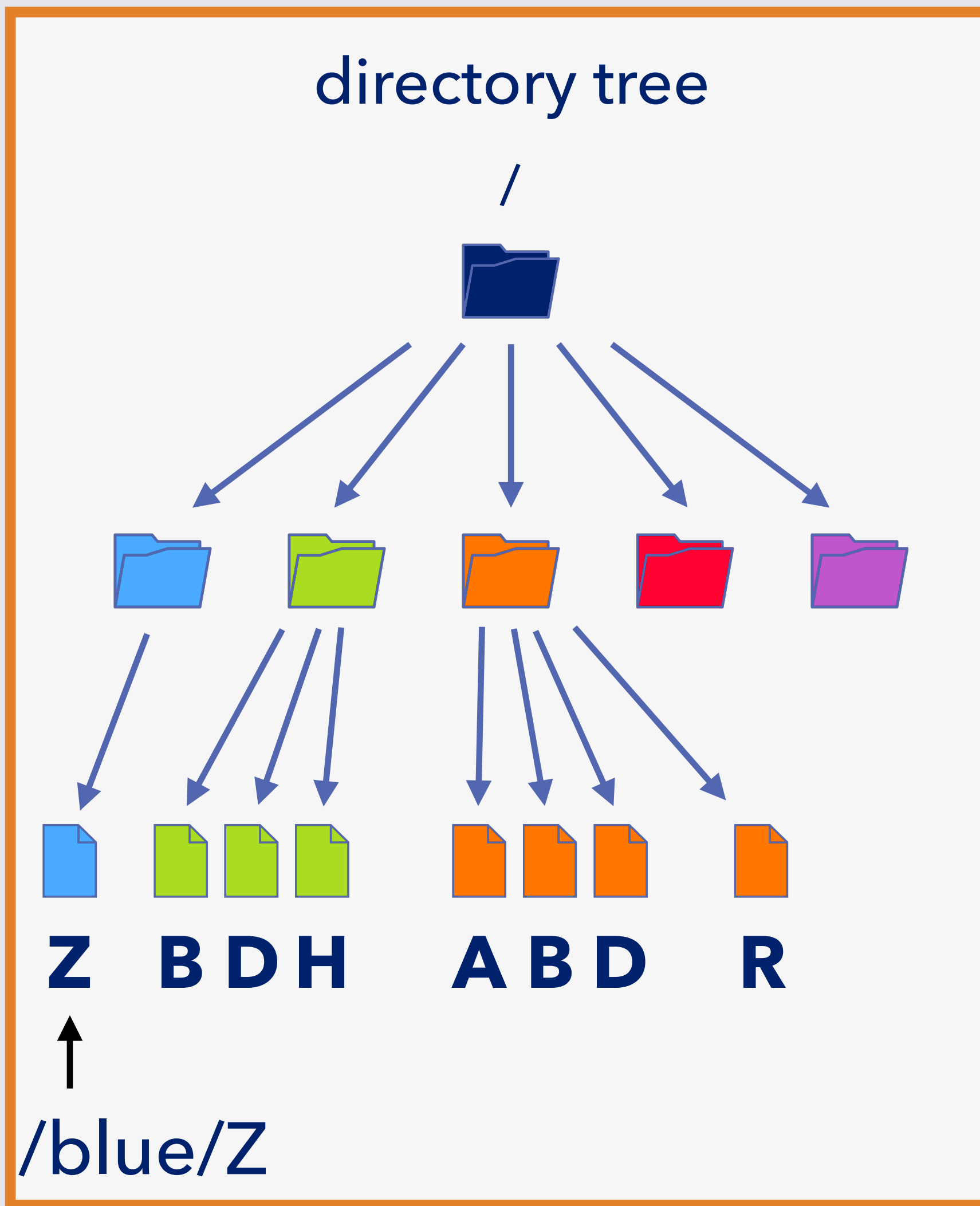


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

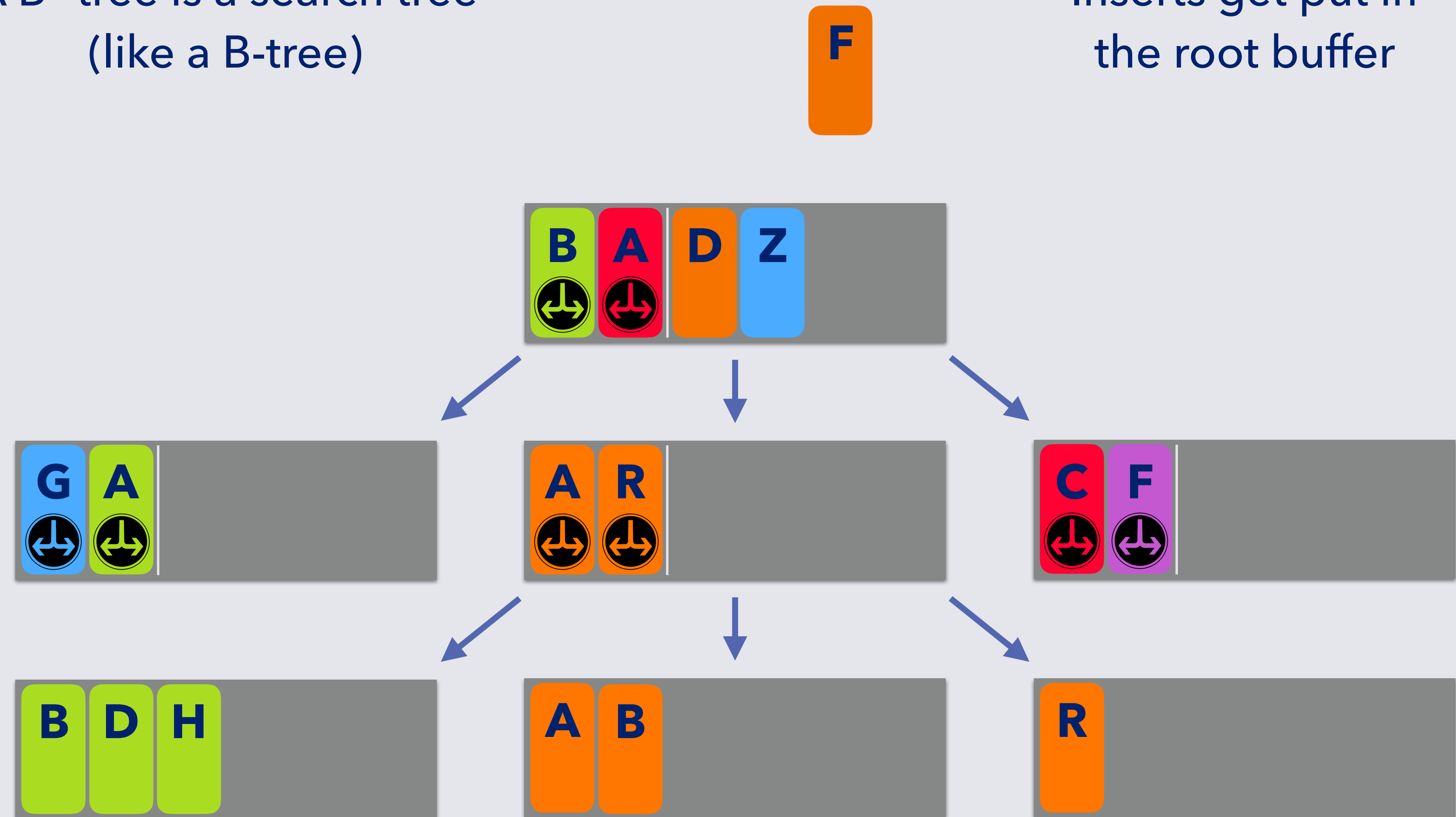
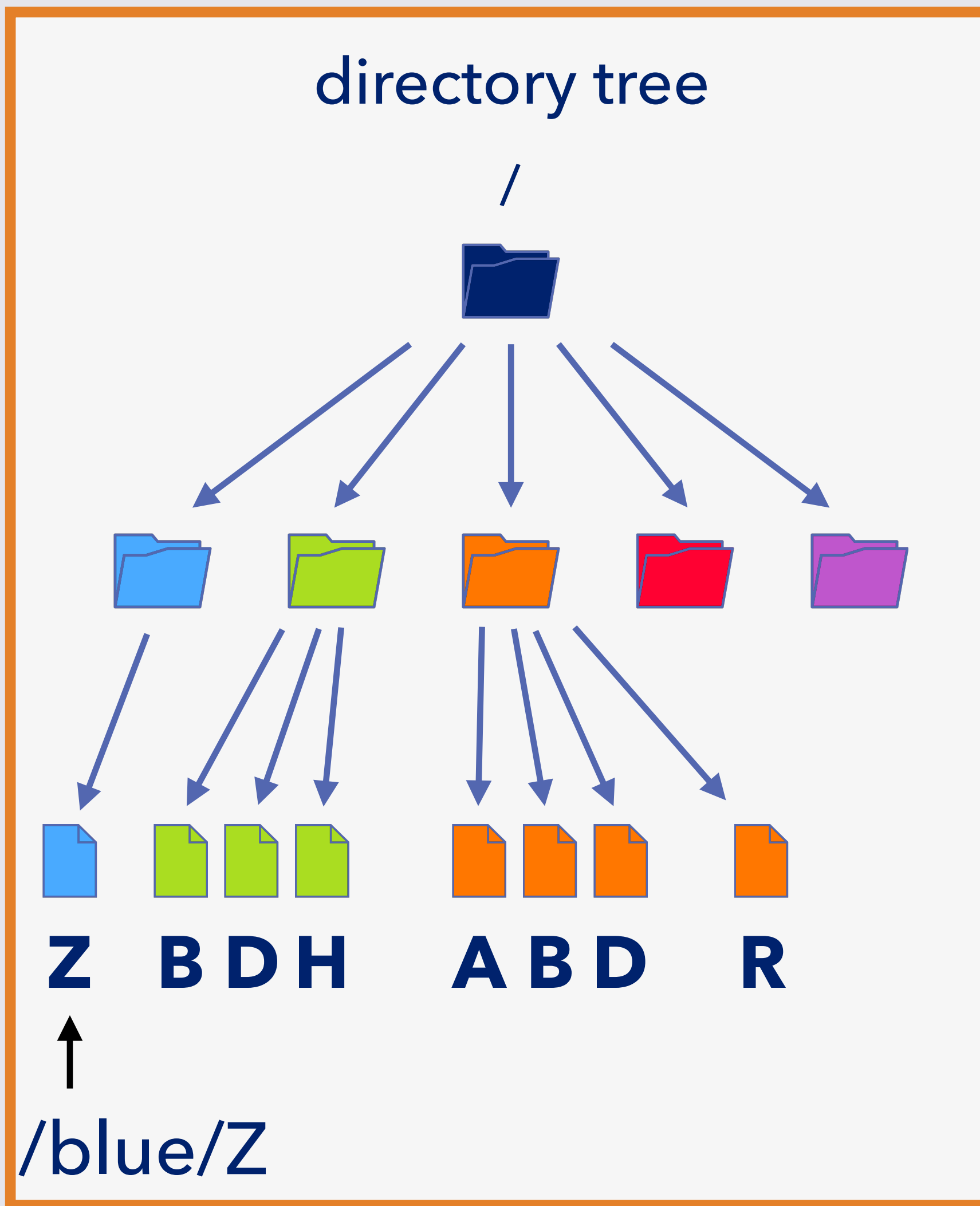


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

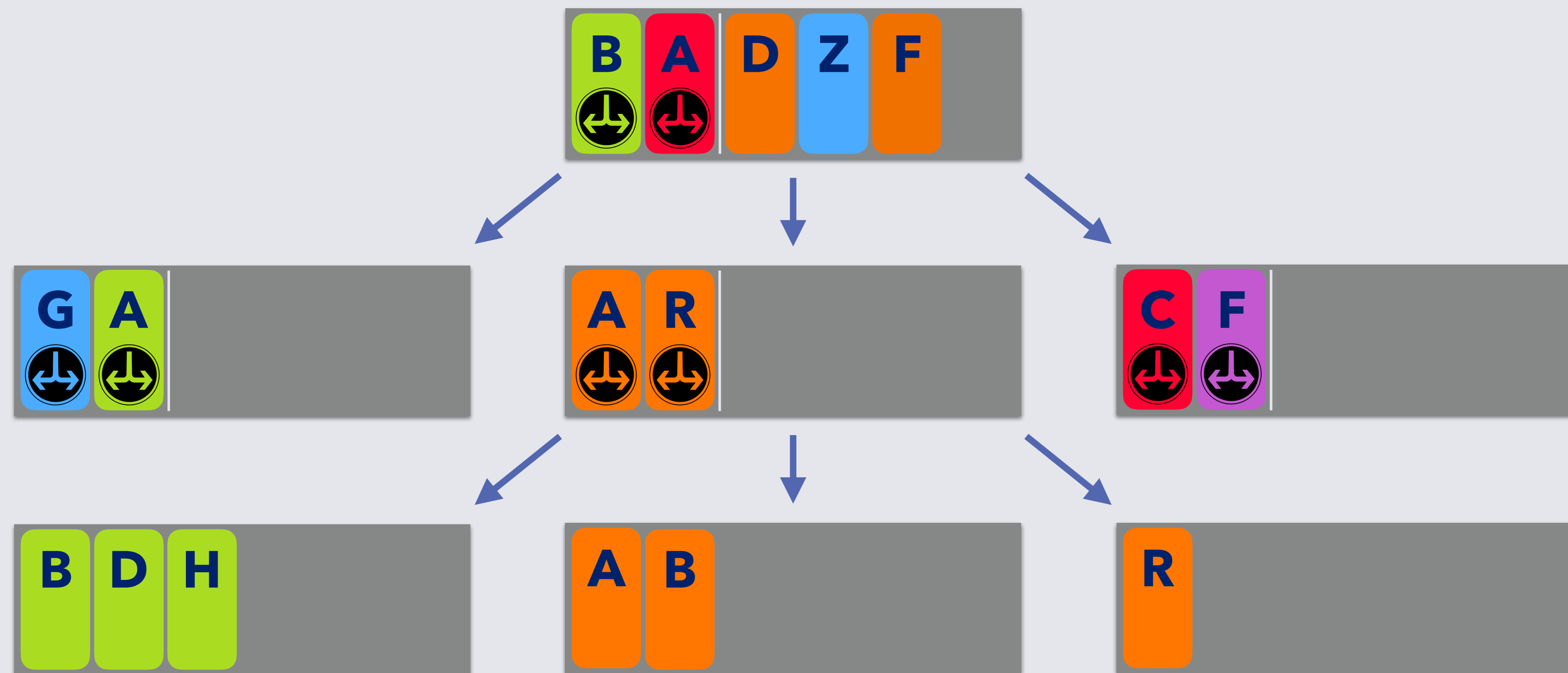
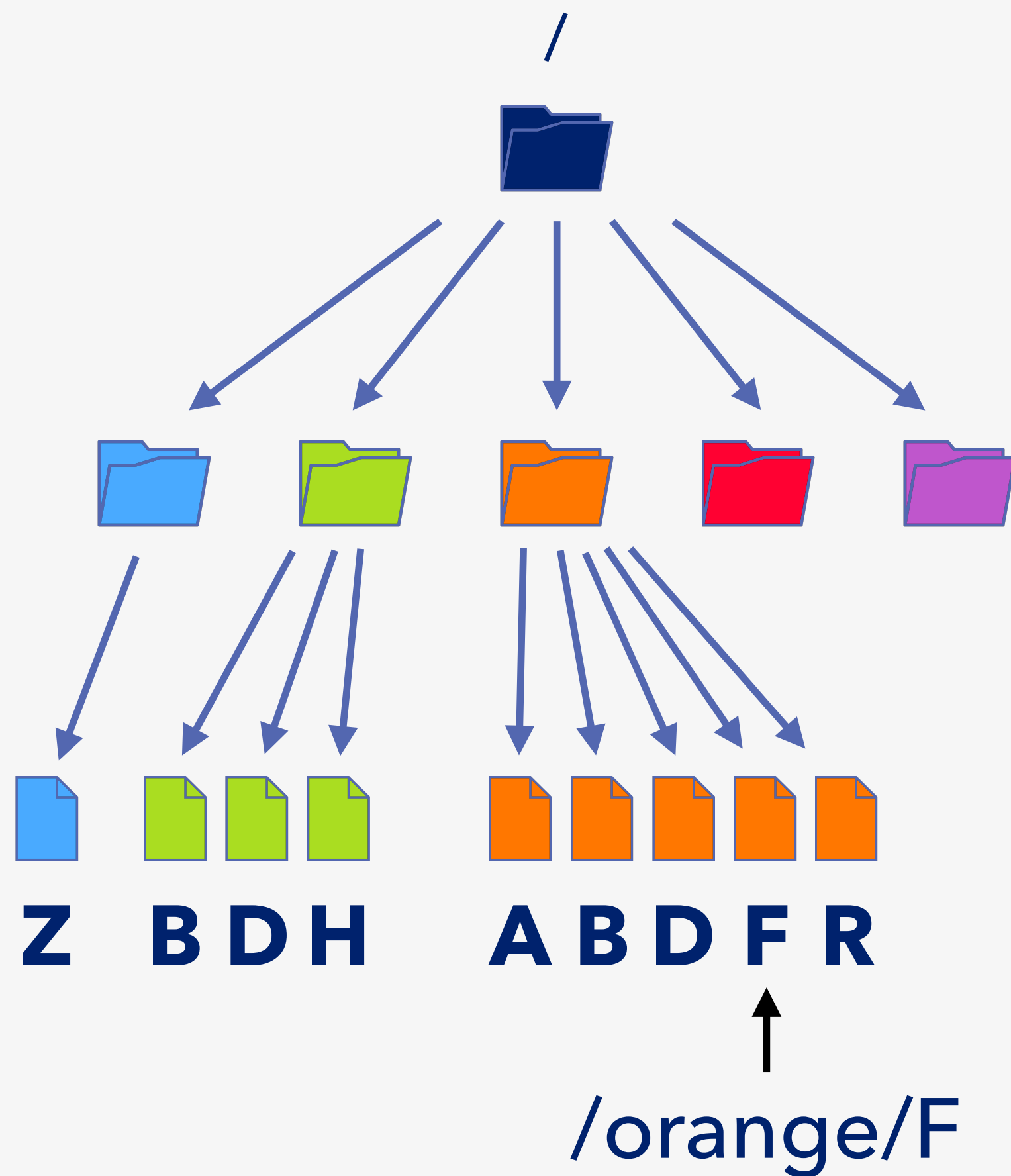


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

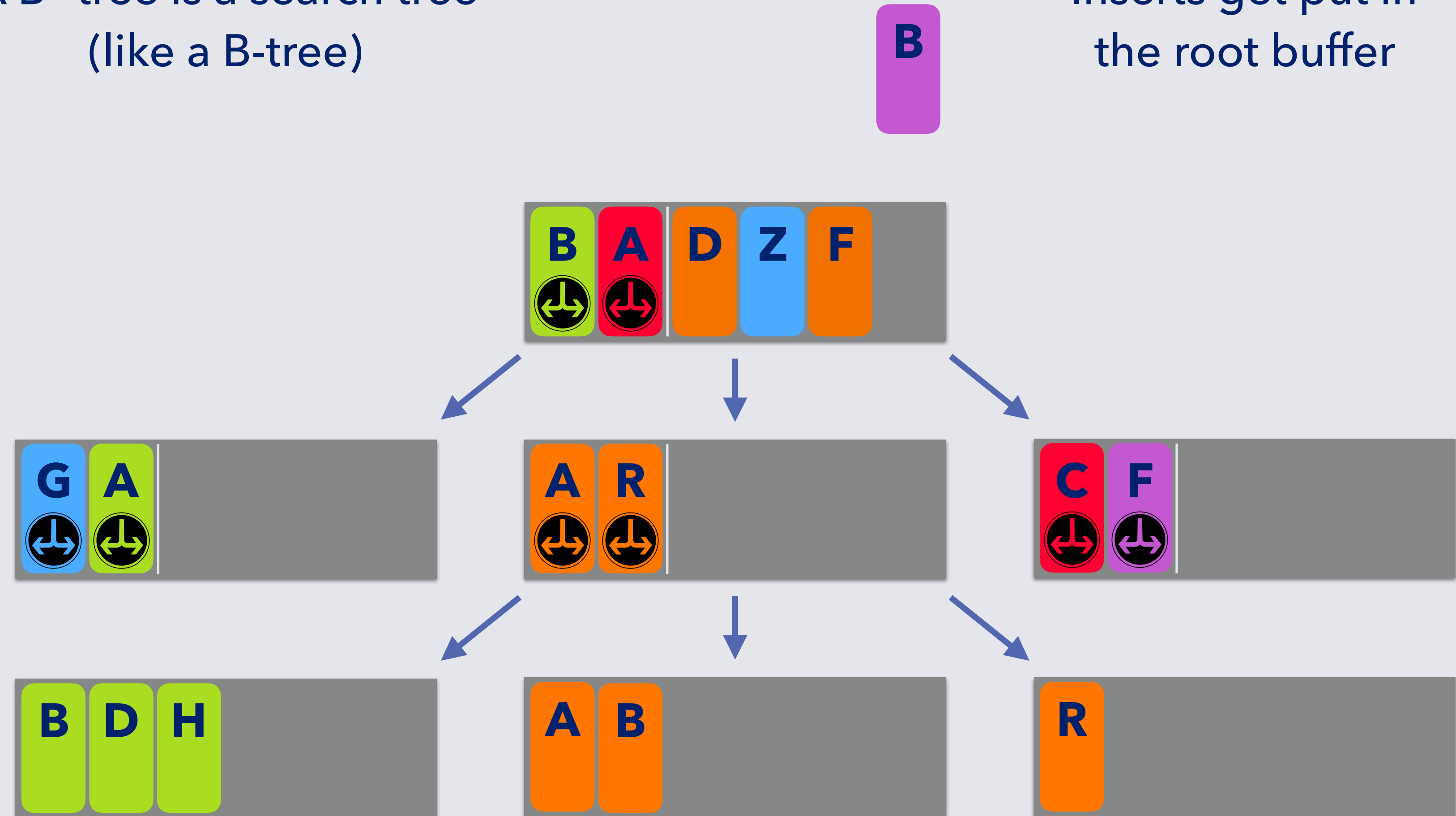
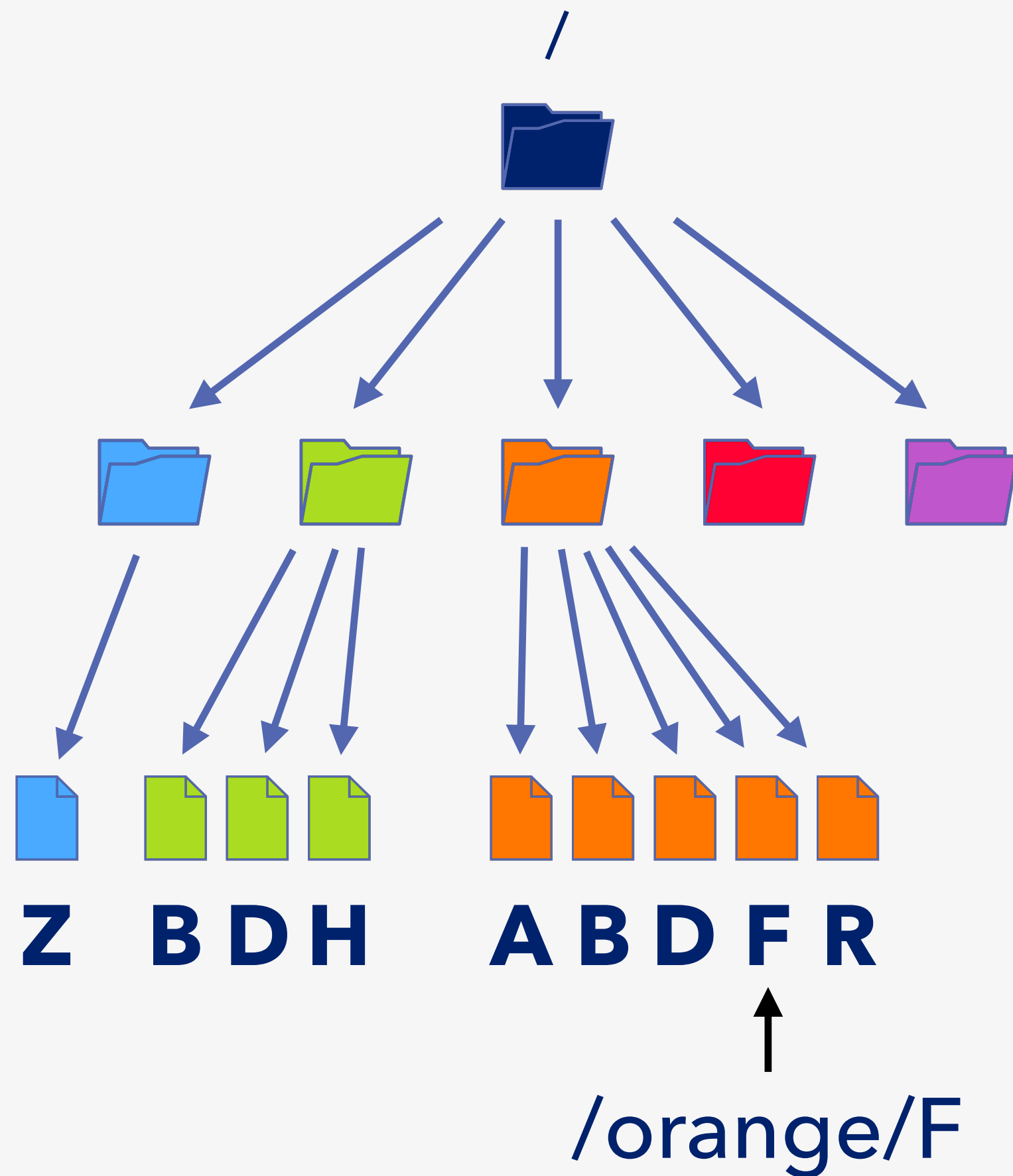


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree



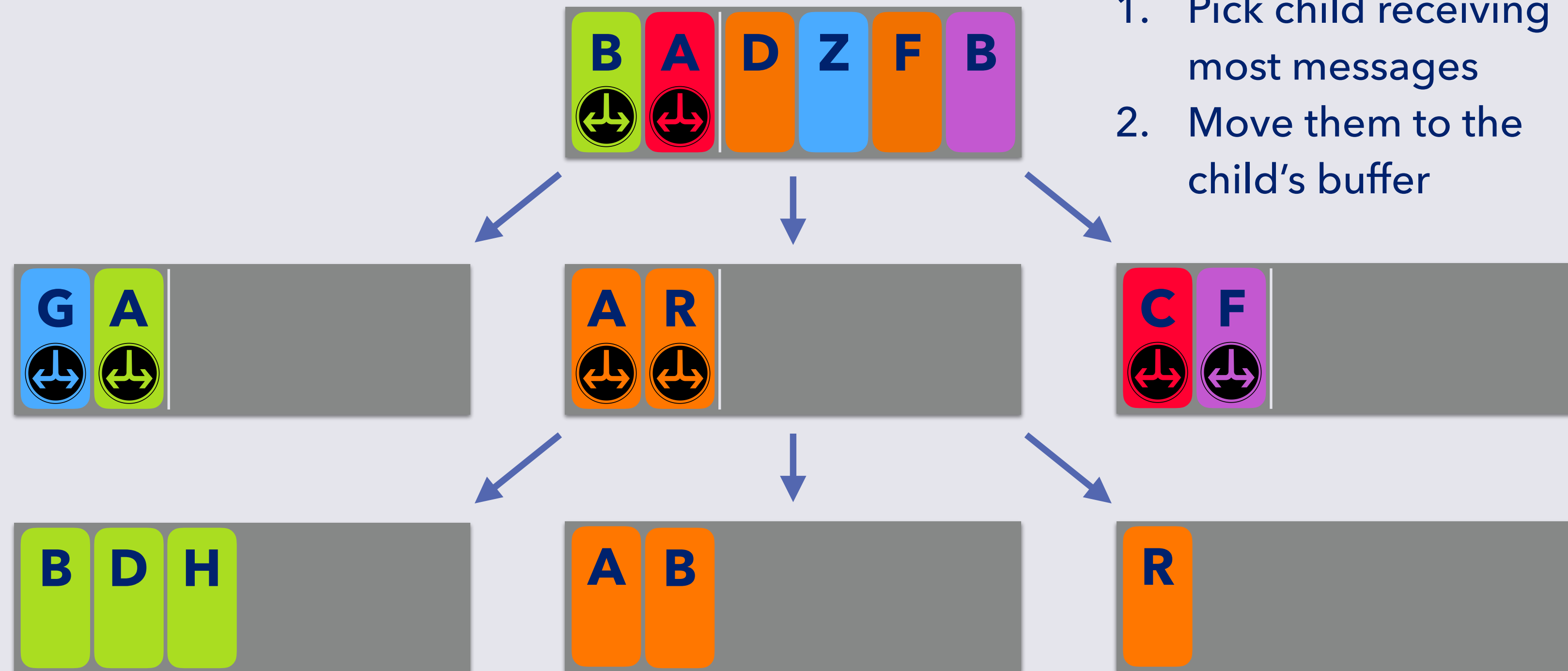
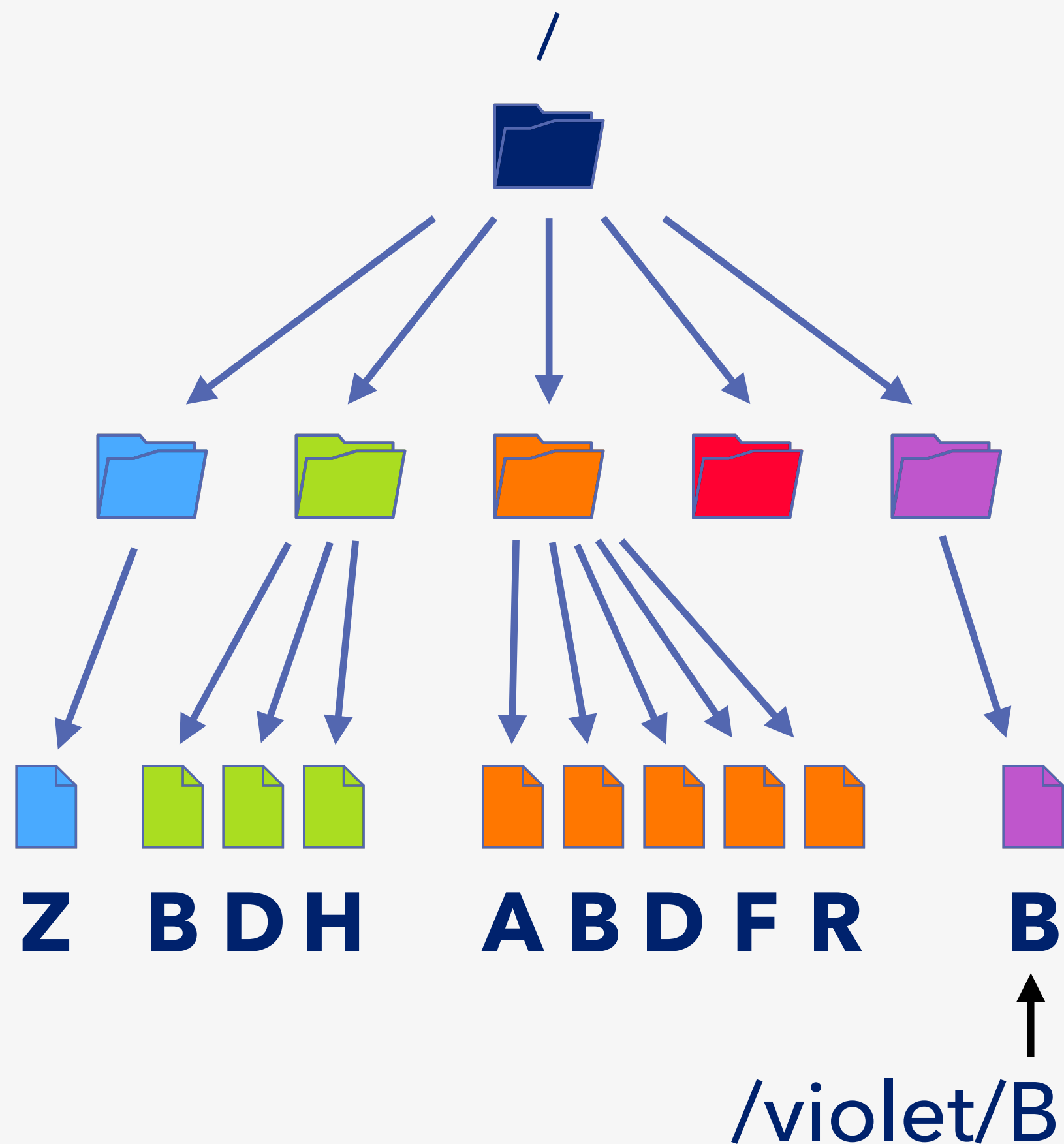
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

directory tree



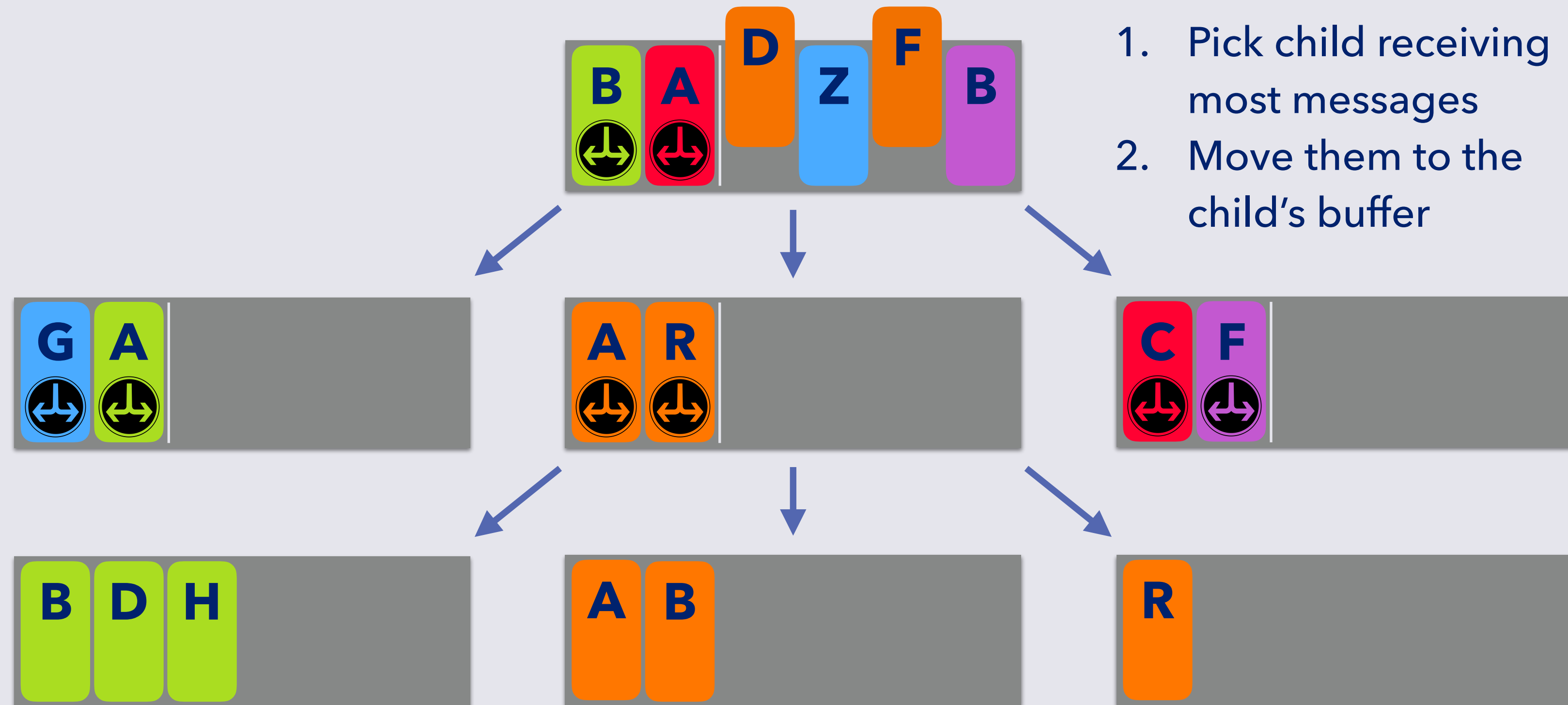
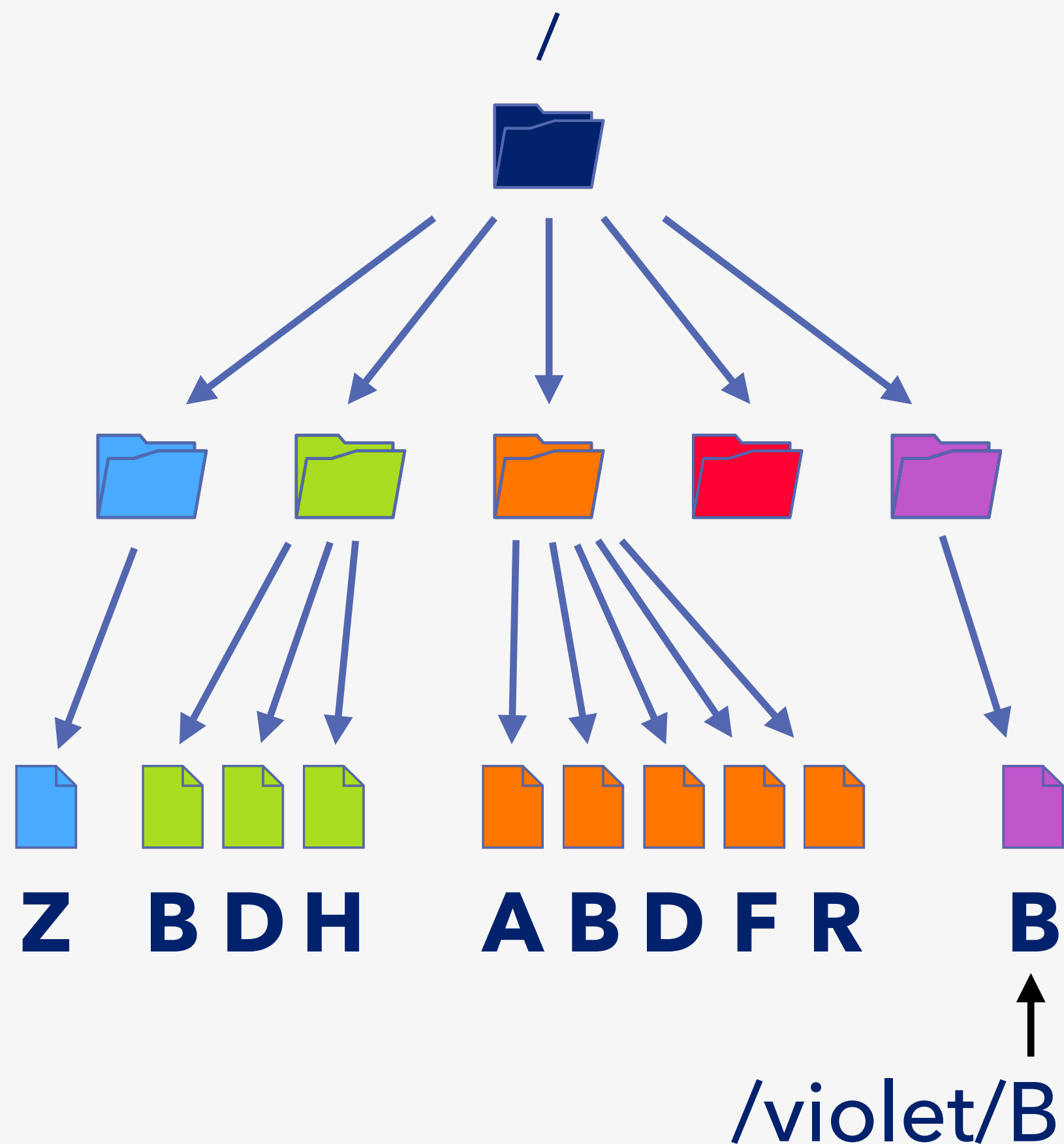
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

directory tree

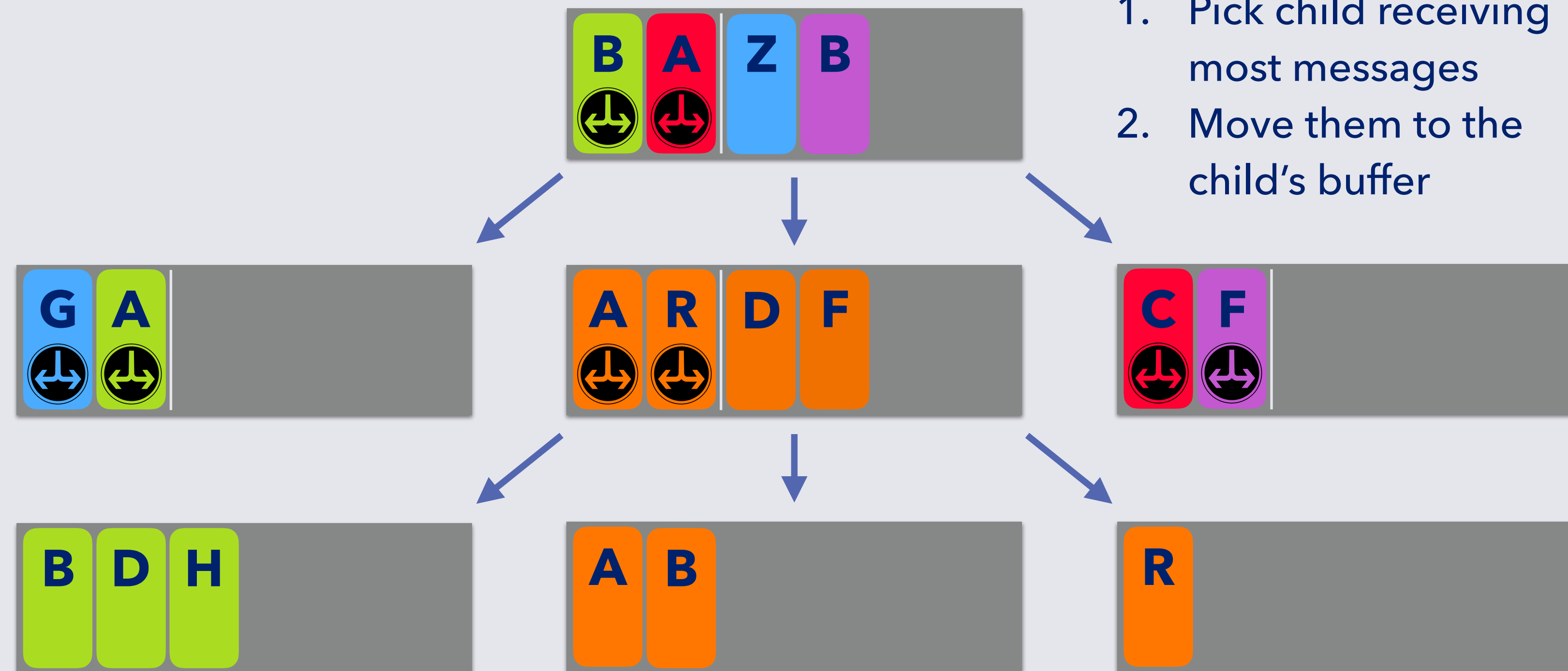
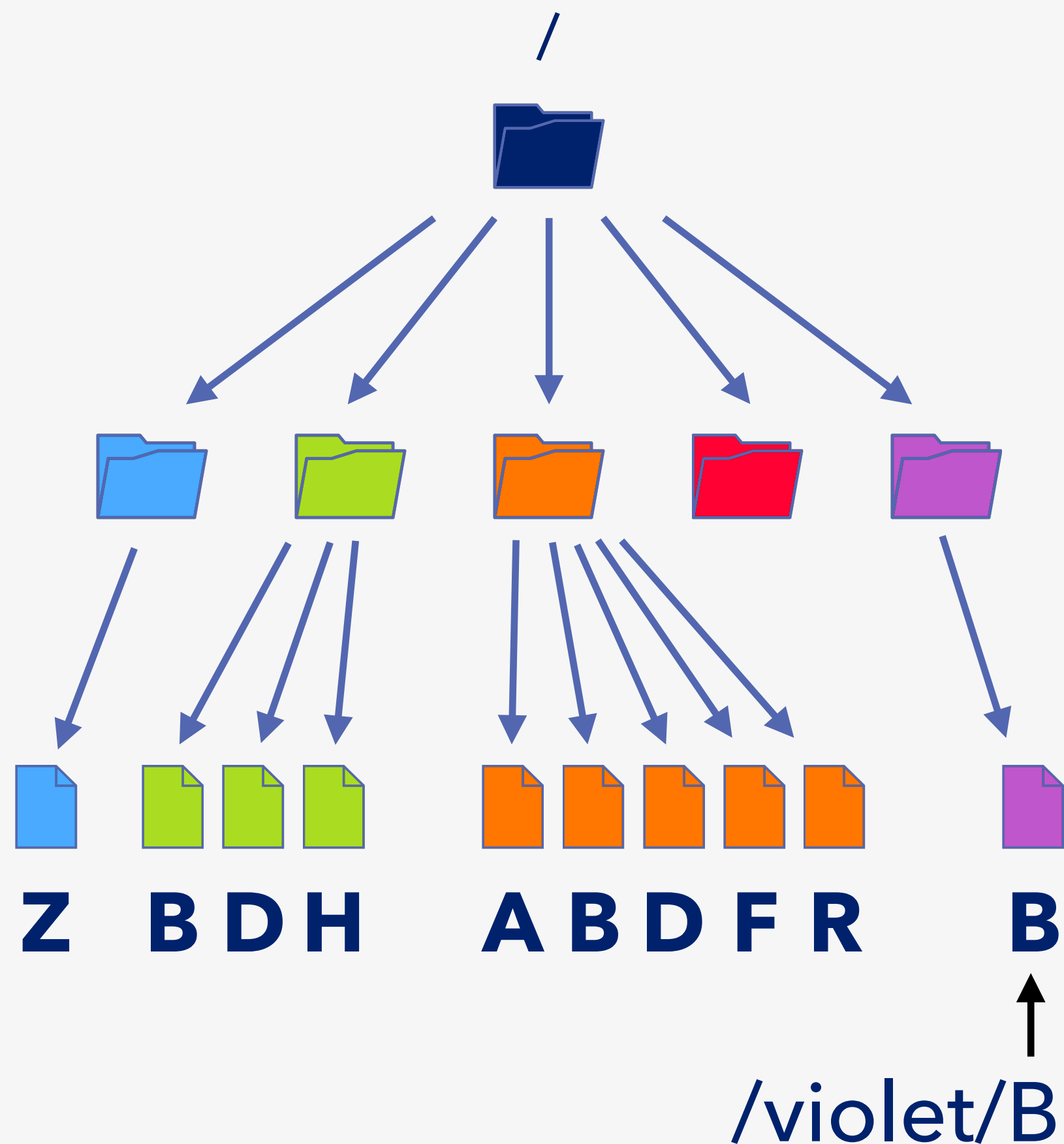


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree

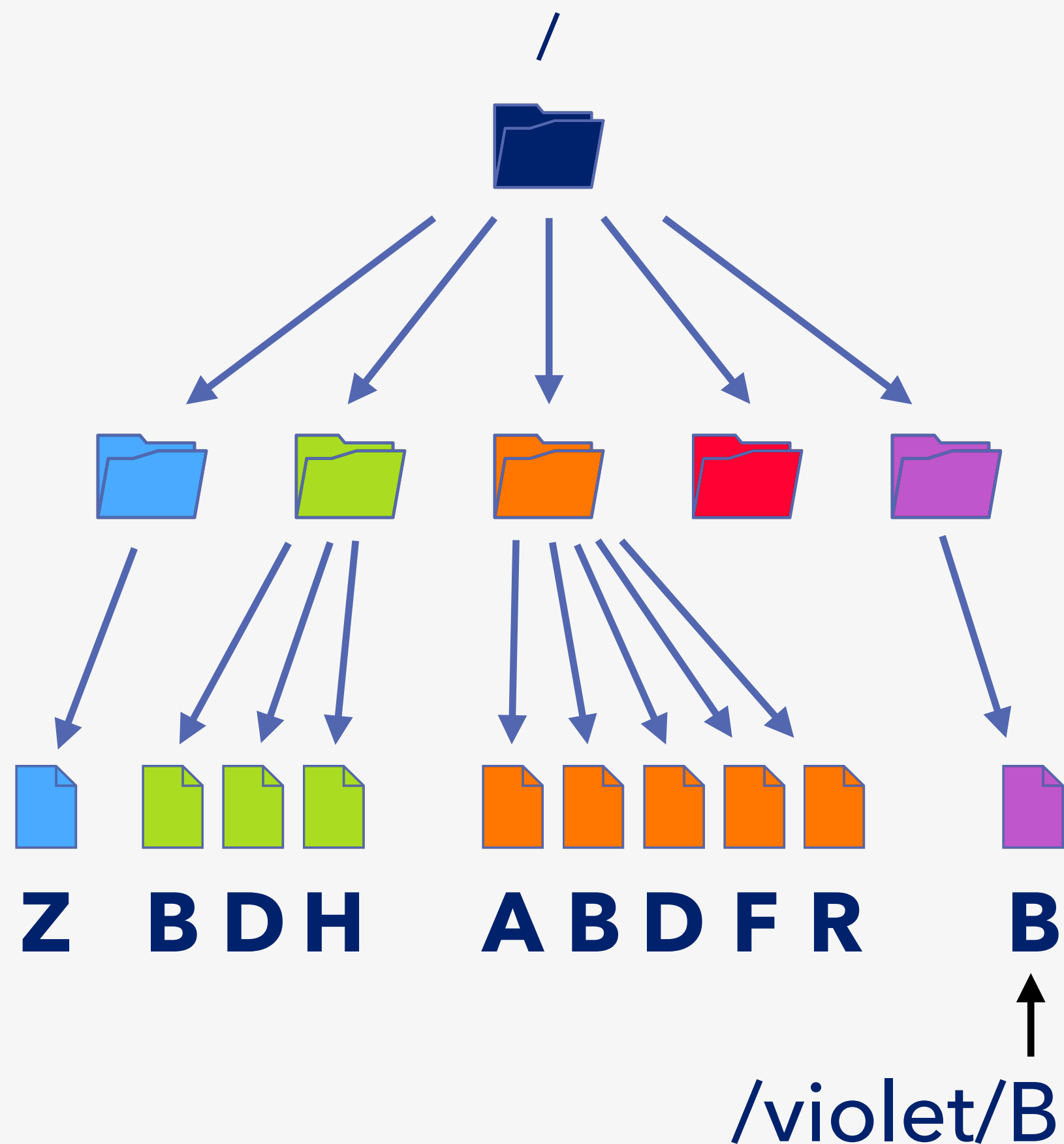


When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

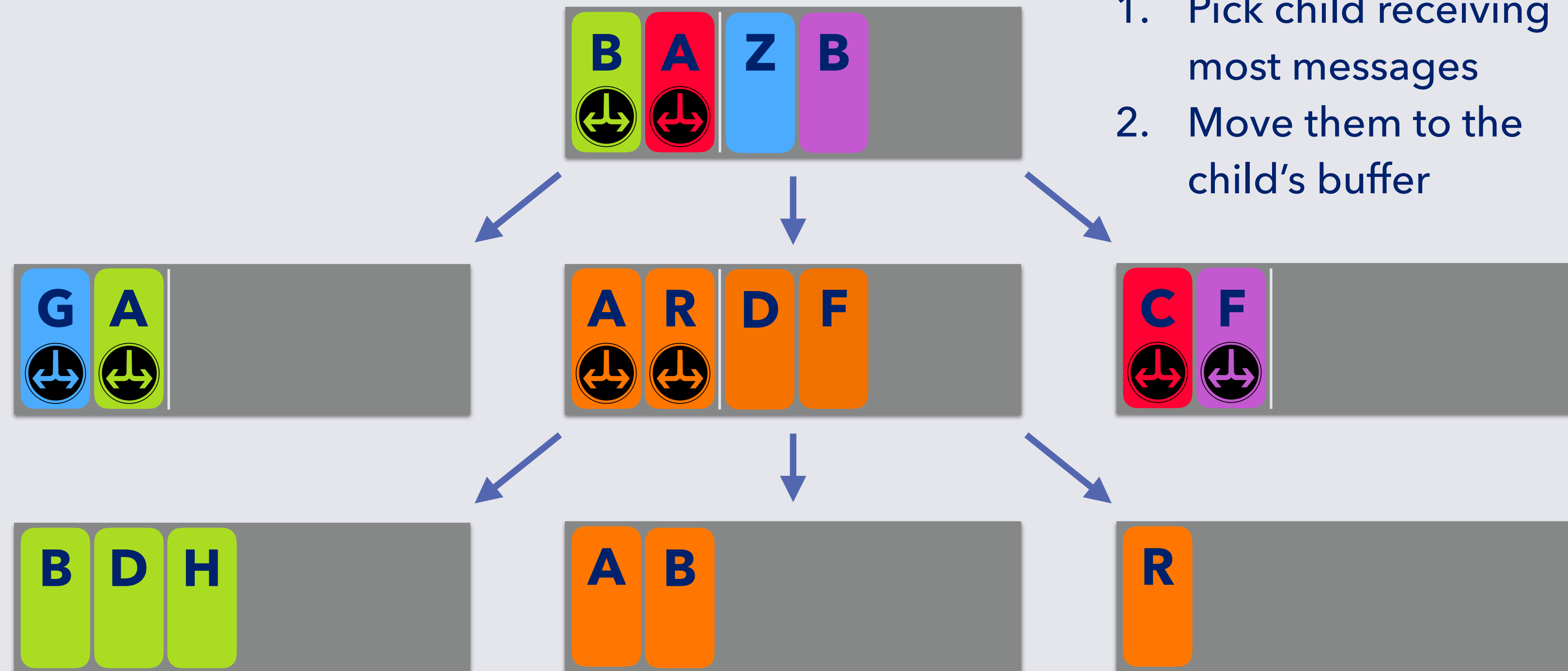
directory tree



L

Inserts get put in
the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer



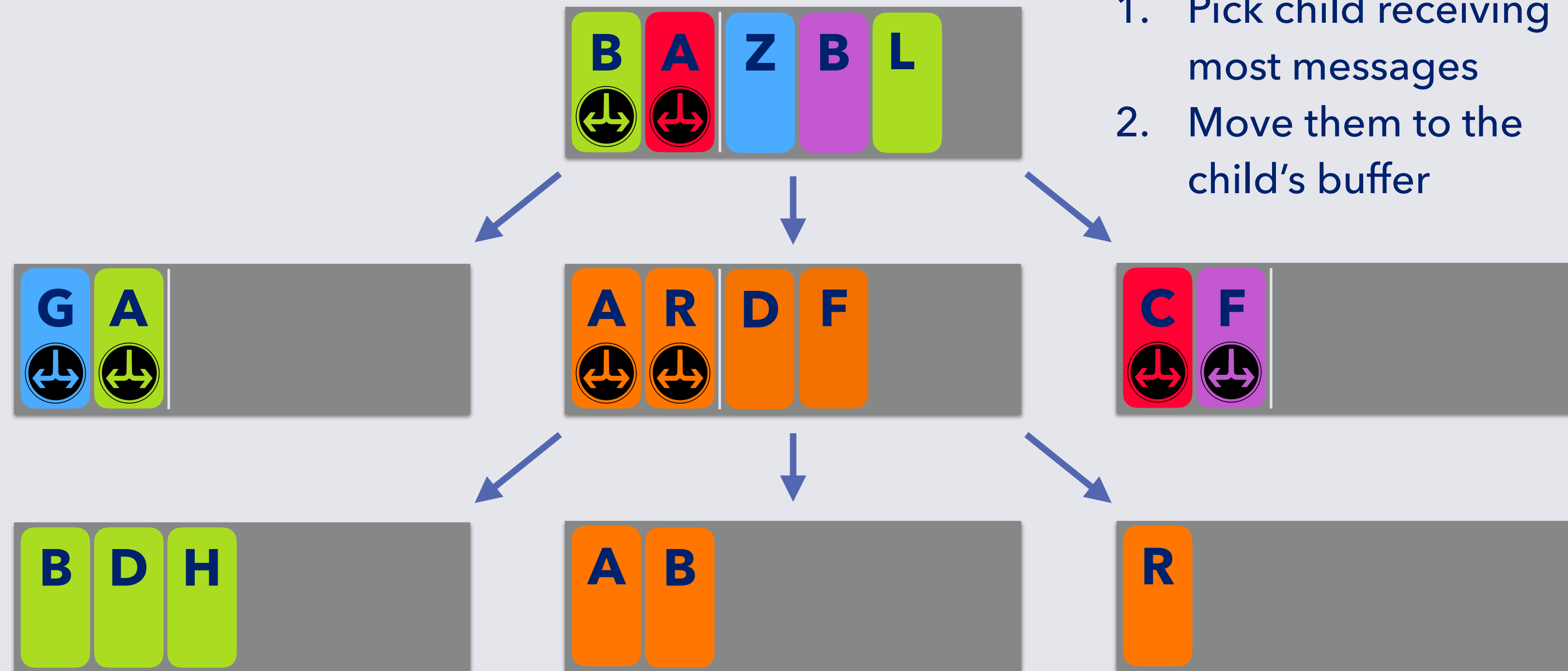
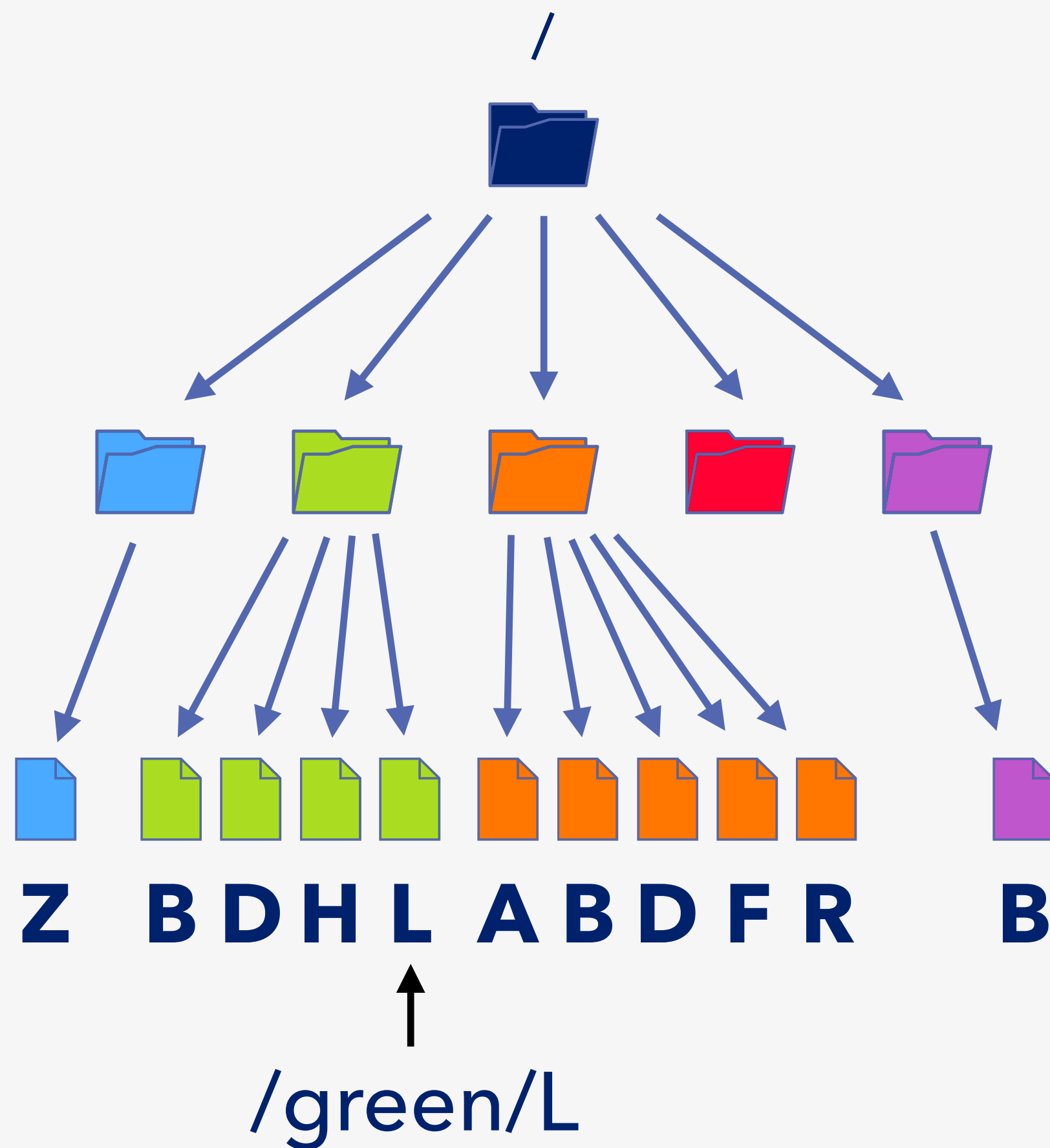
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

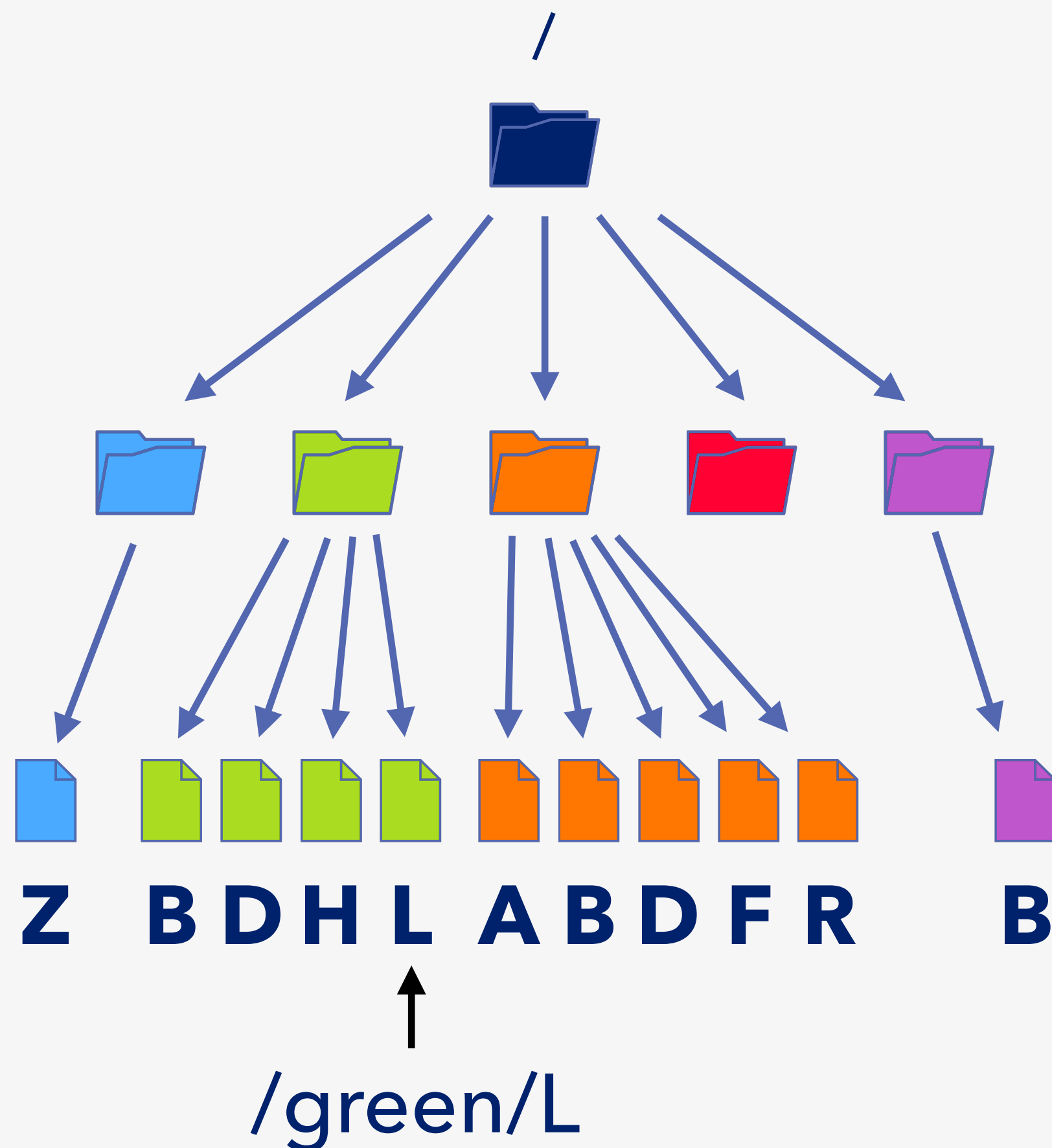
directory tree



B ϵ -Trees

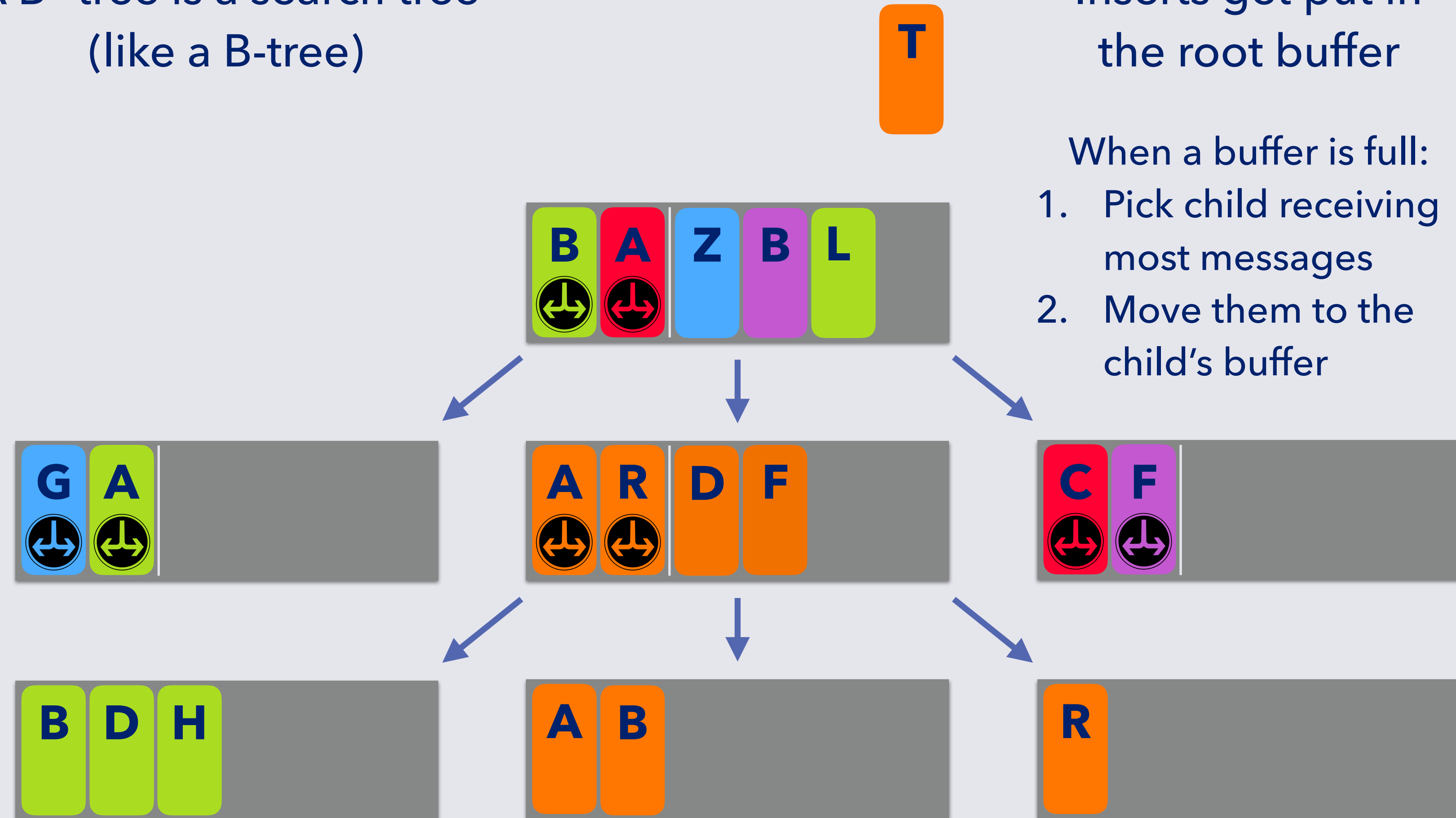
A B ϵ -tree is a search tree
(like a B-tree)

directory tree



Inserts get put in
the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer

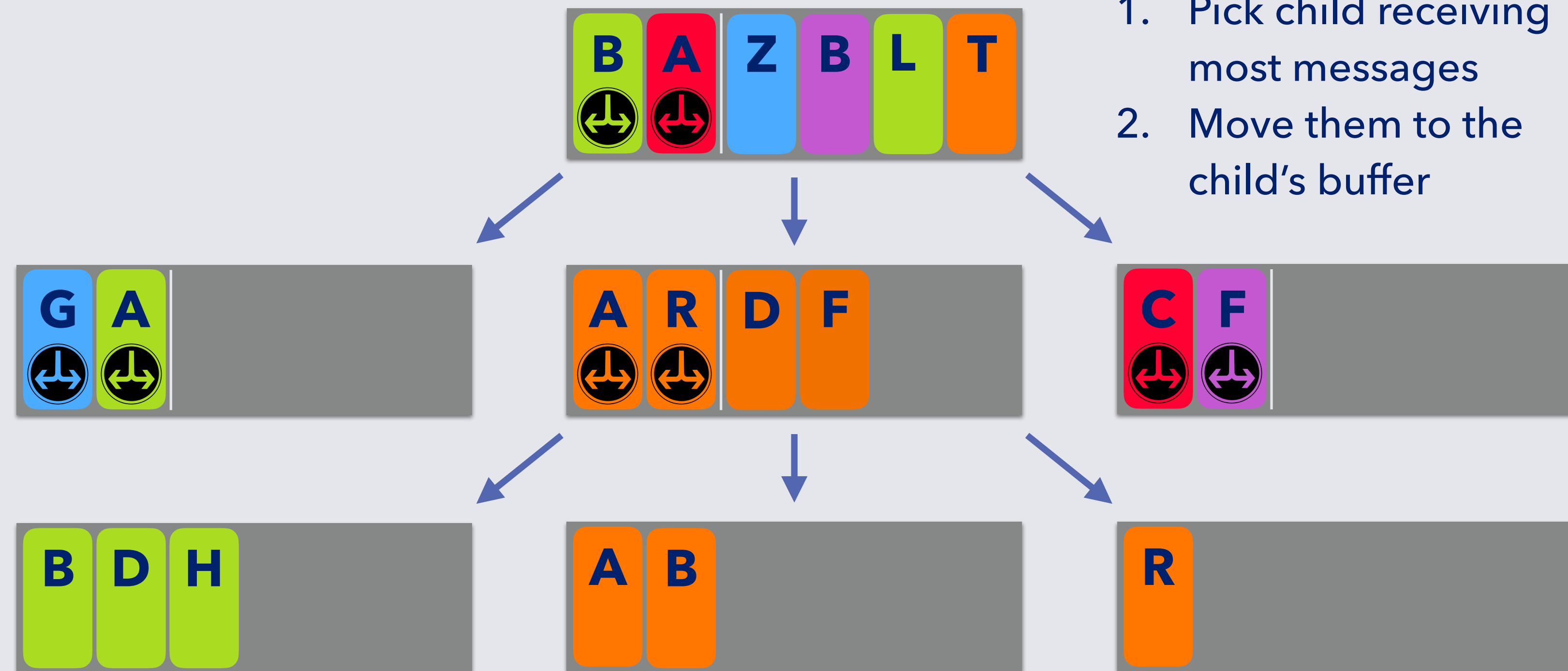
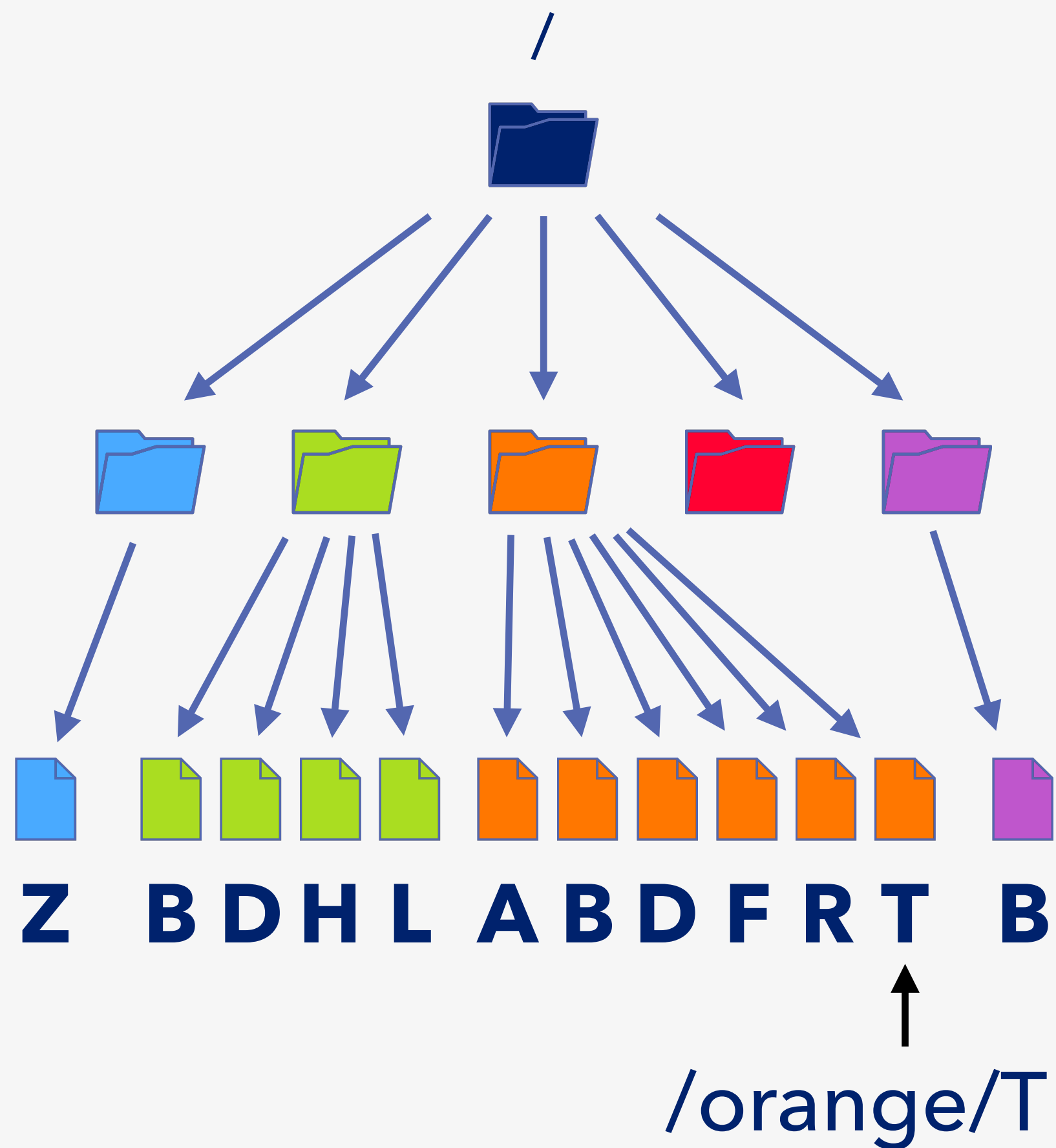


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree



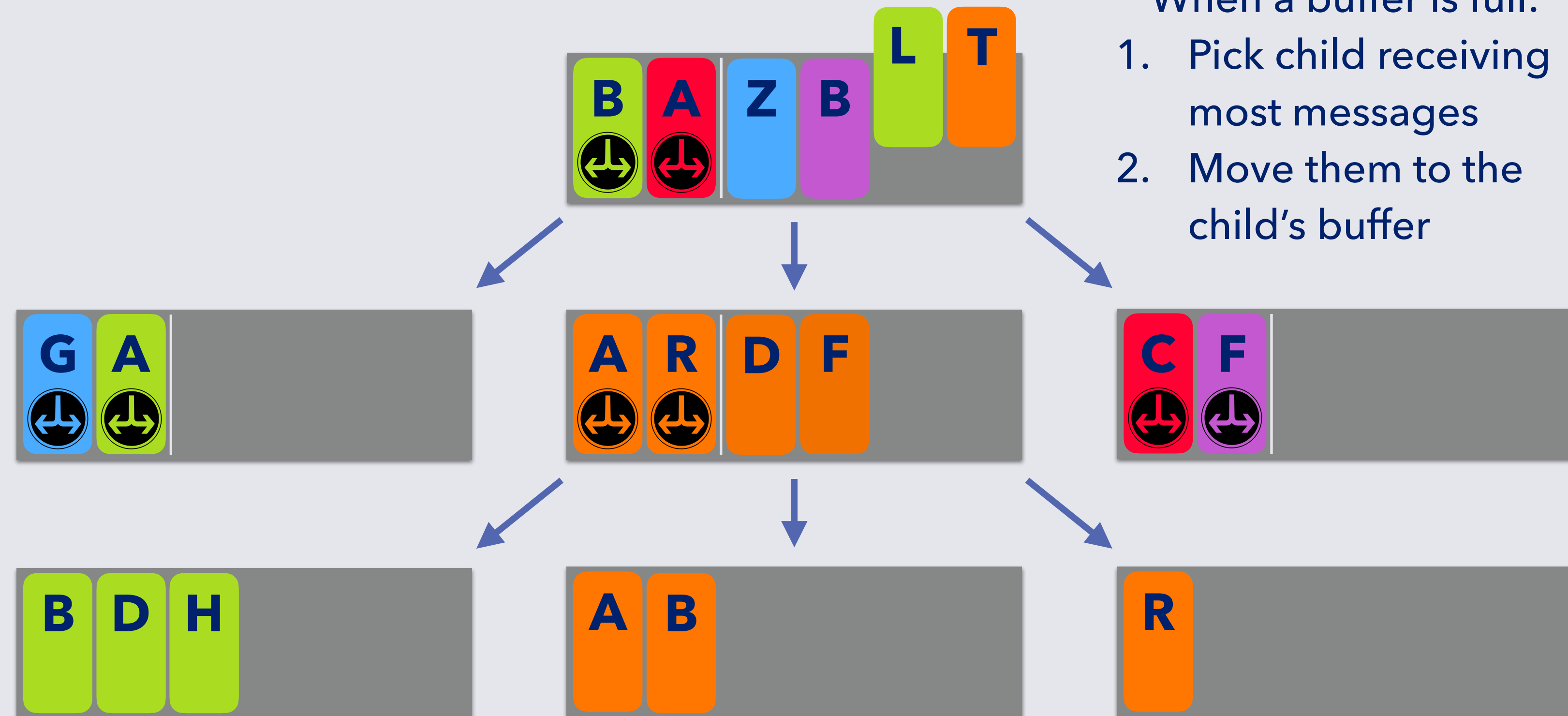
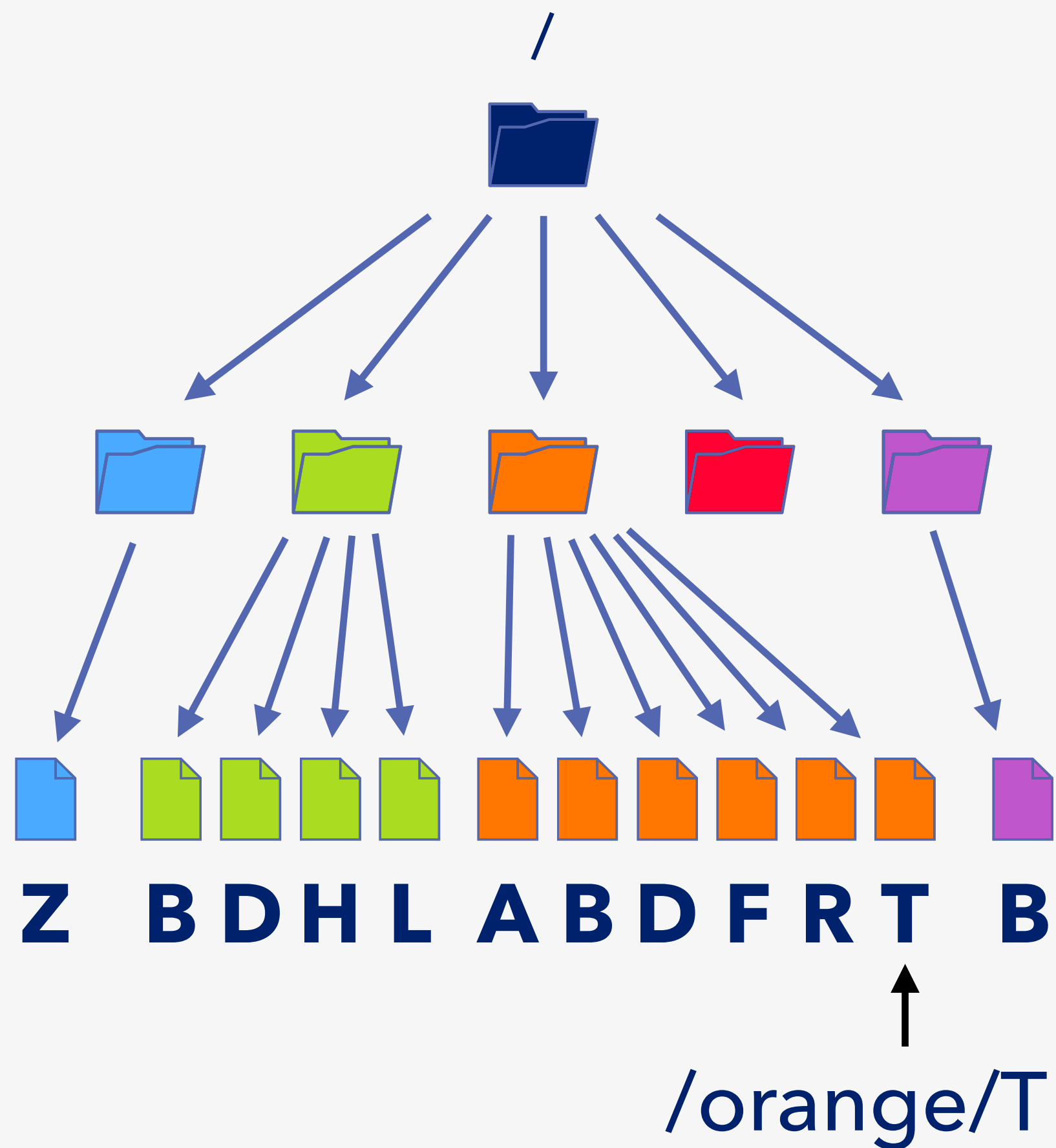
When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree



When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

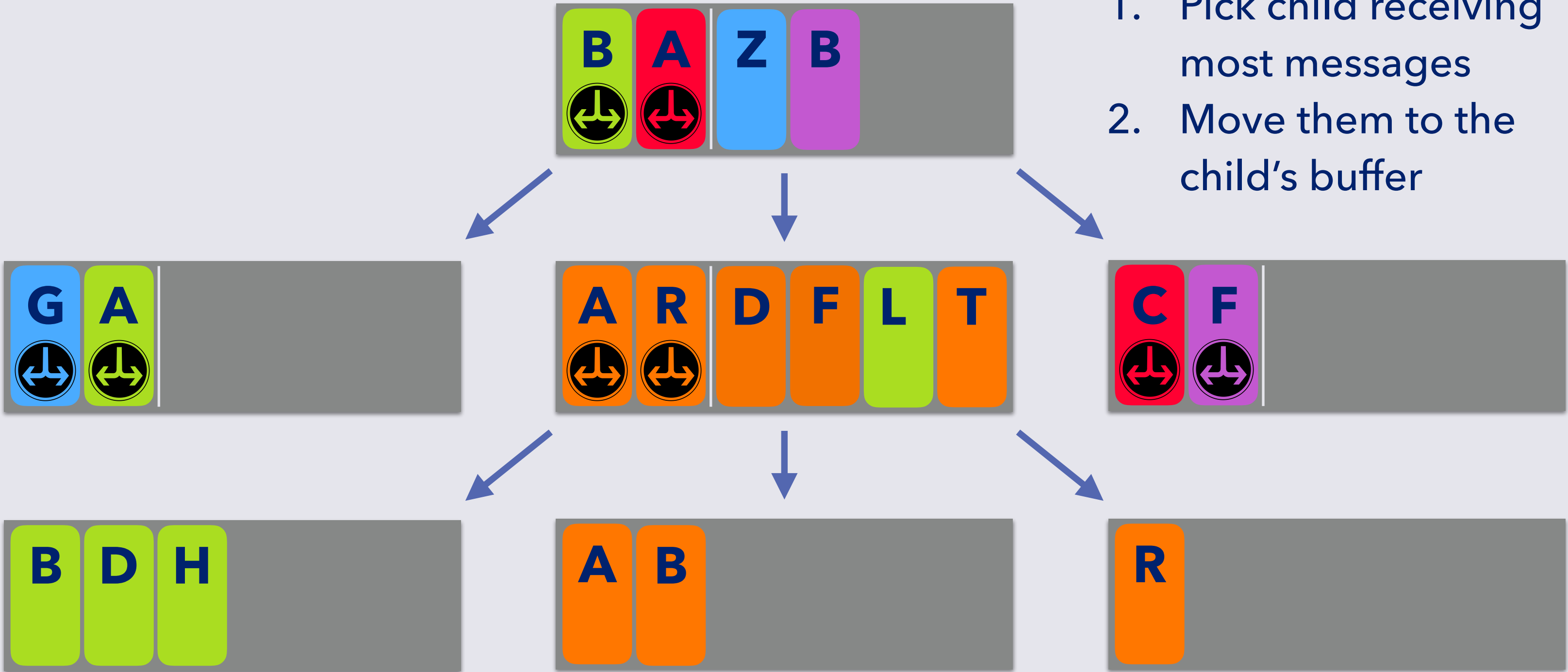
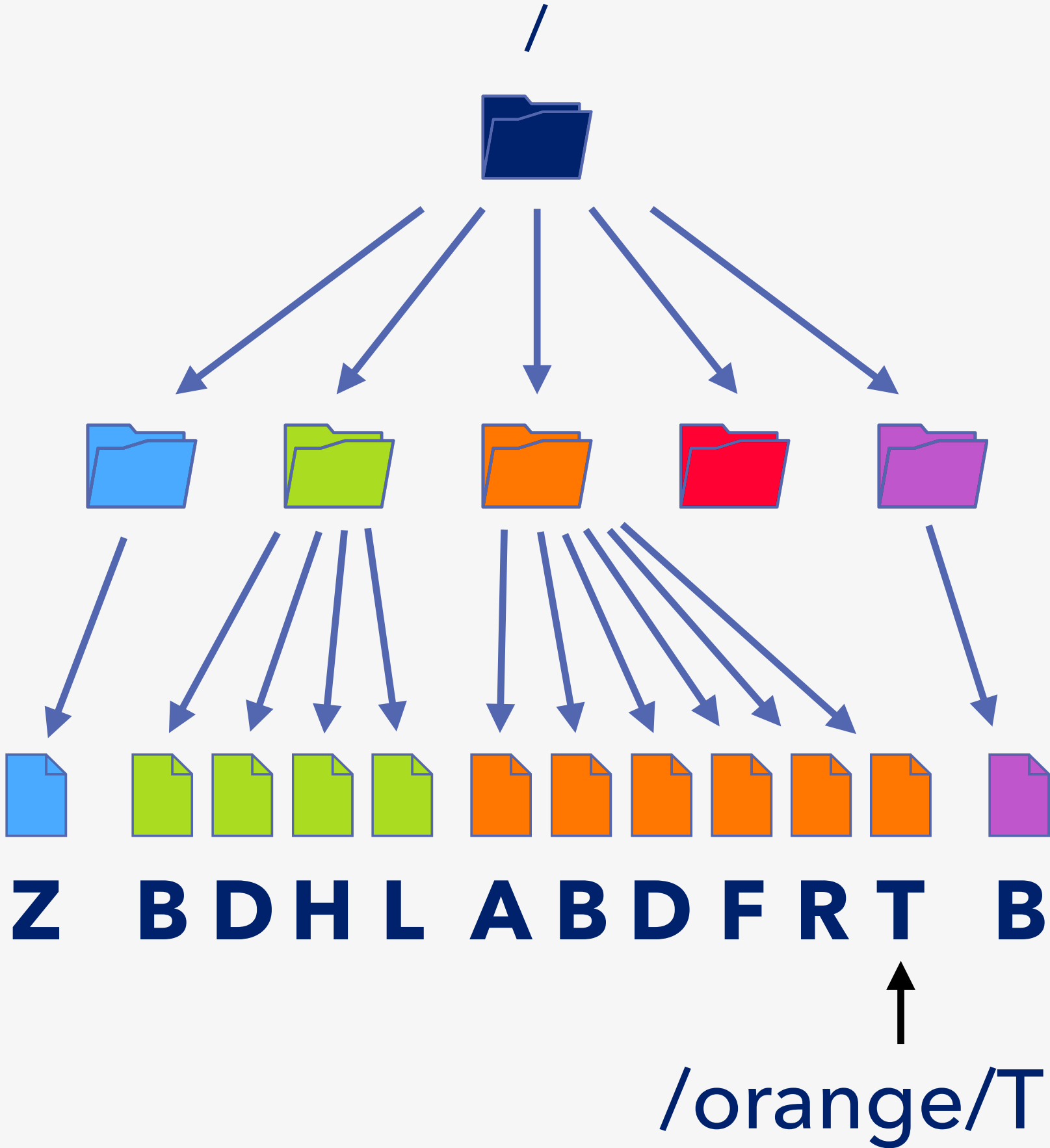
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

- When a buffer is full:
1. Pick child receiving most messages
 2. Move them to the child's buffer

directory tree

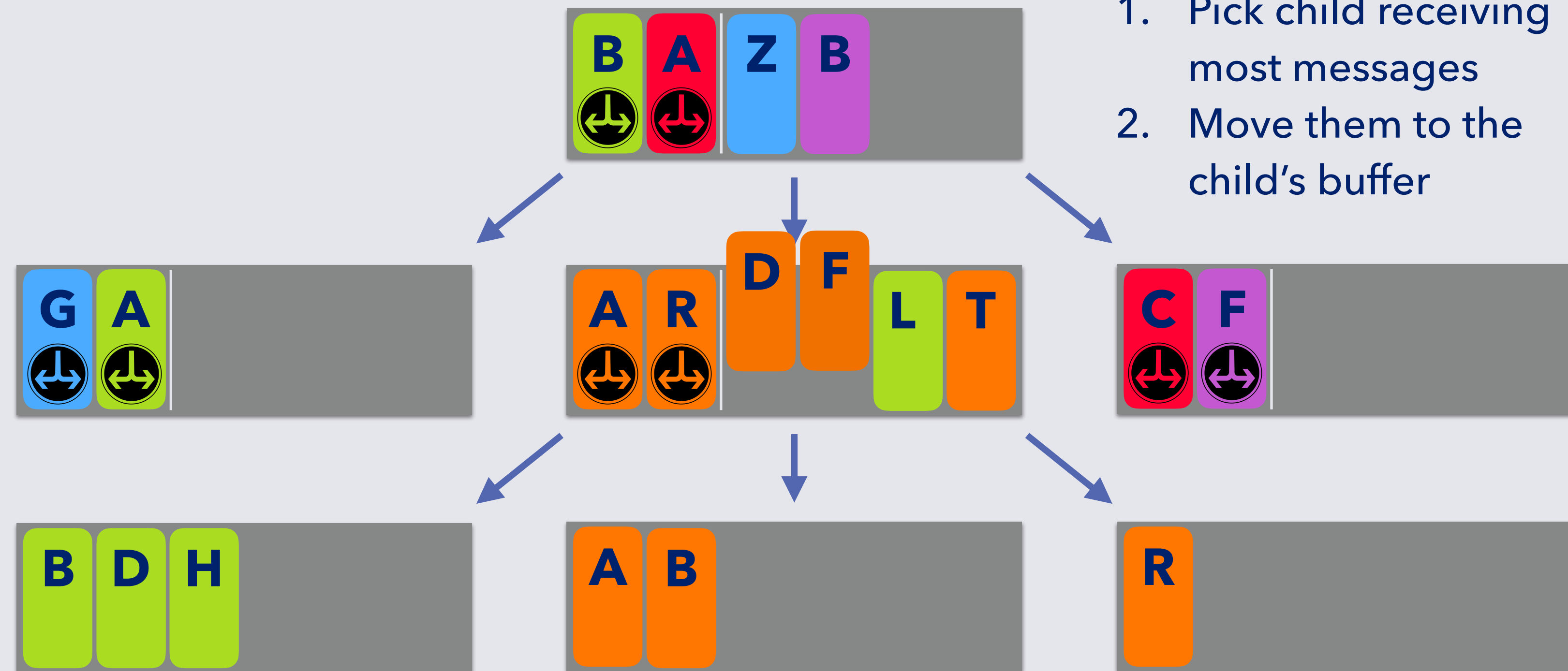
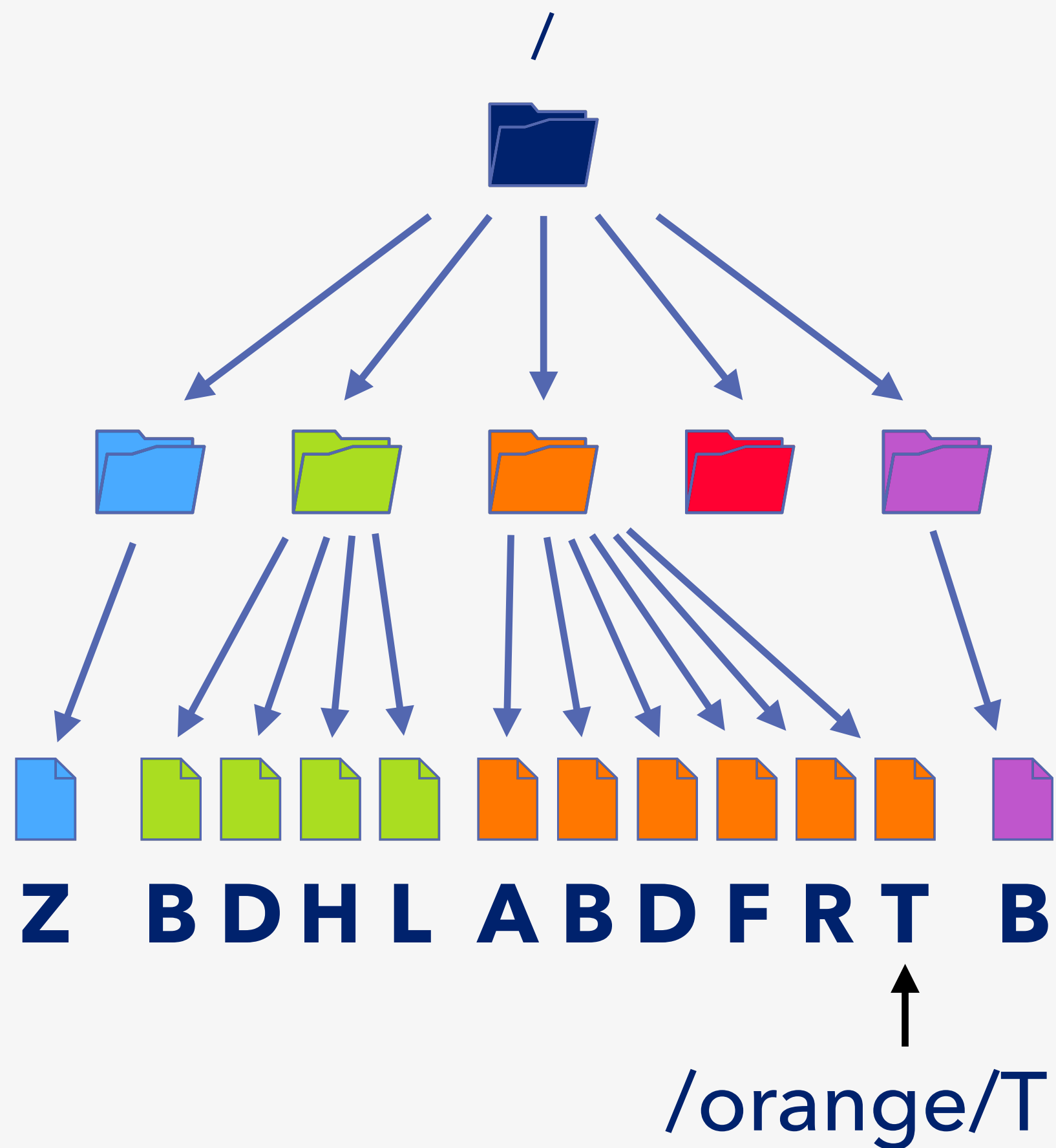


B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree



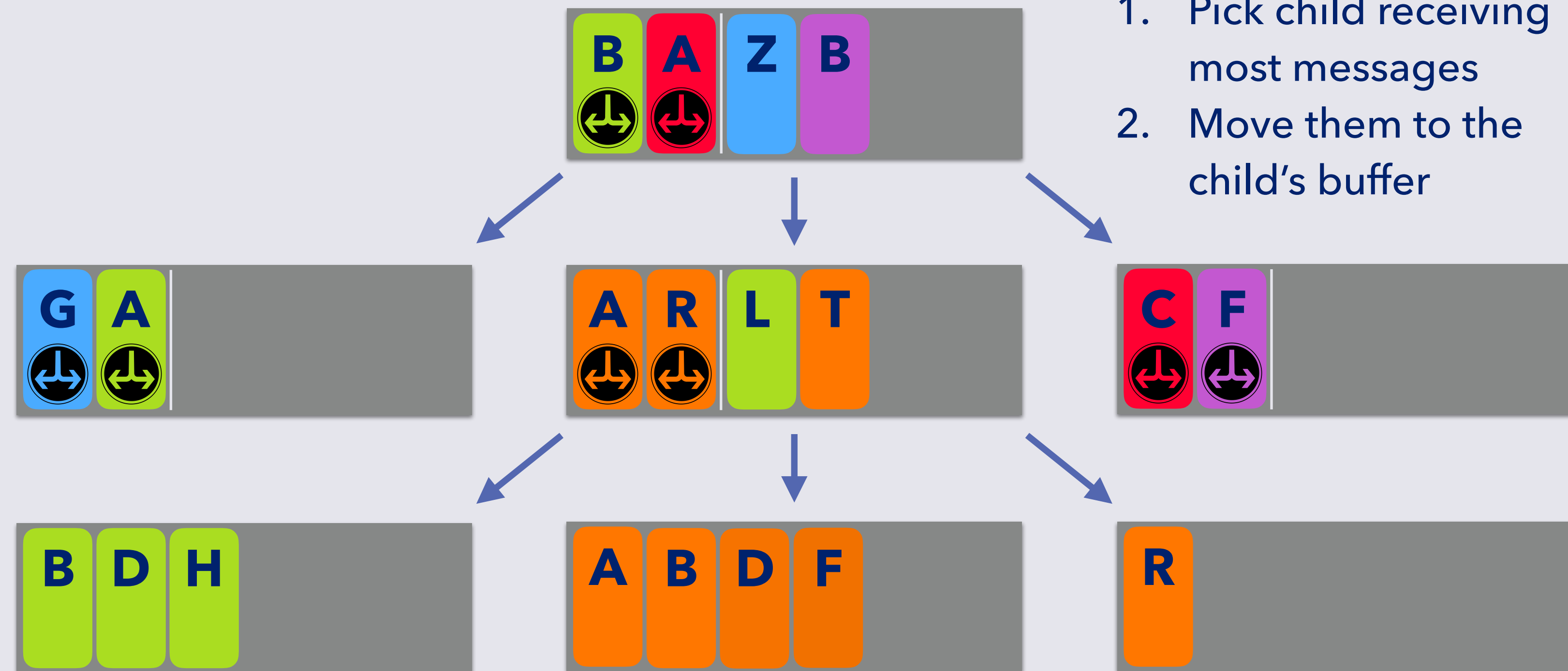
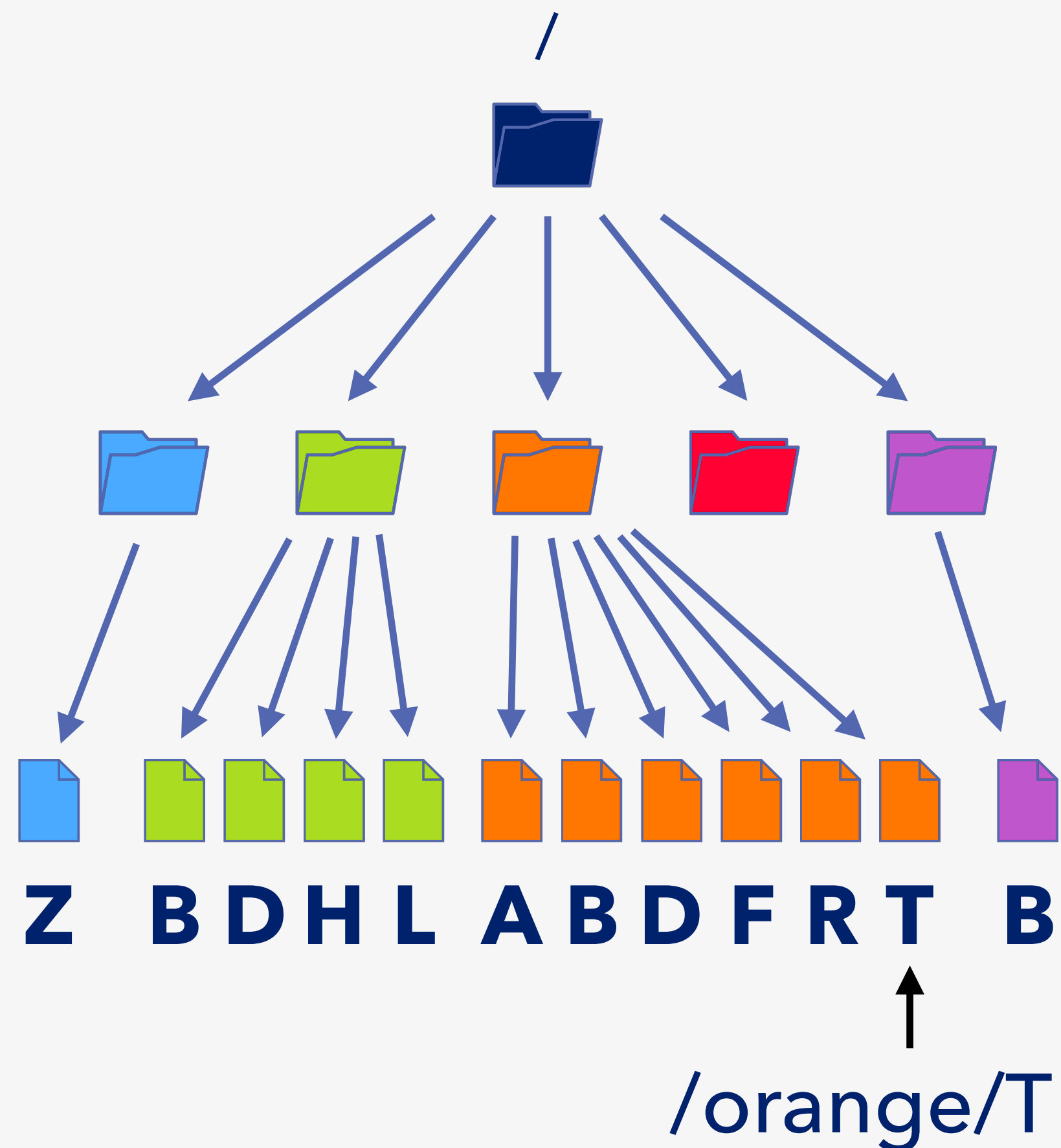
When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

directory tree



When a buffer is full:
1. Pick child receiving
most messages
2. Move them to the
child's buffer

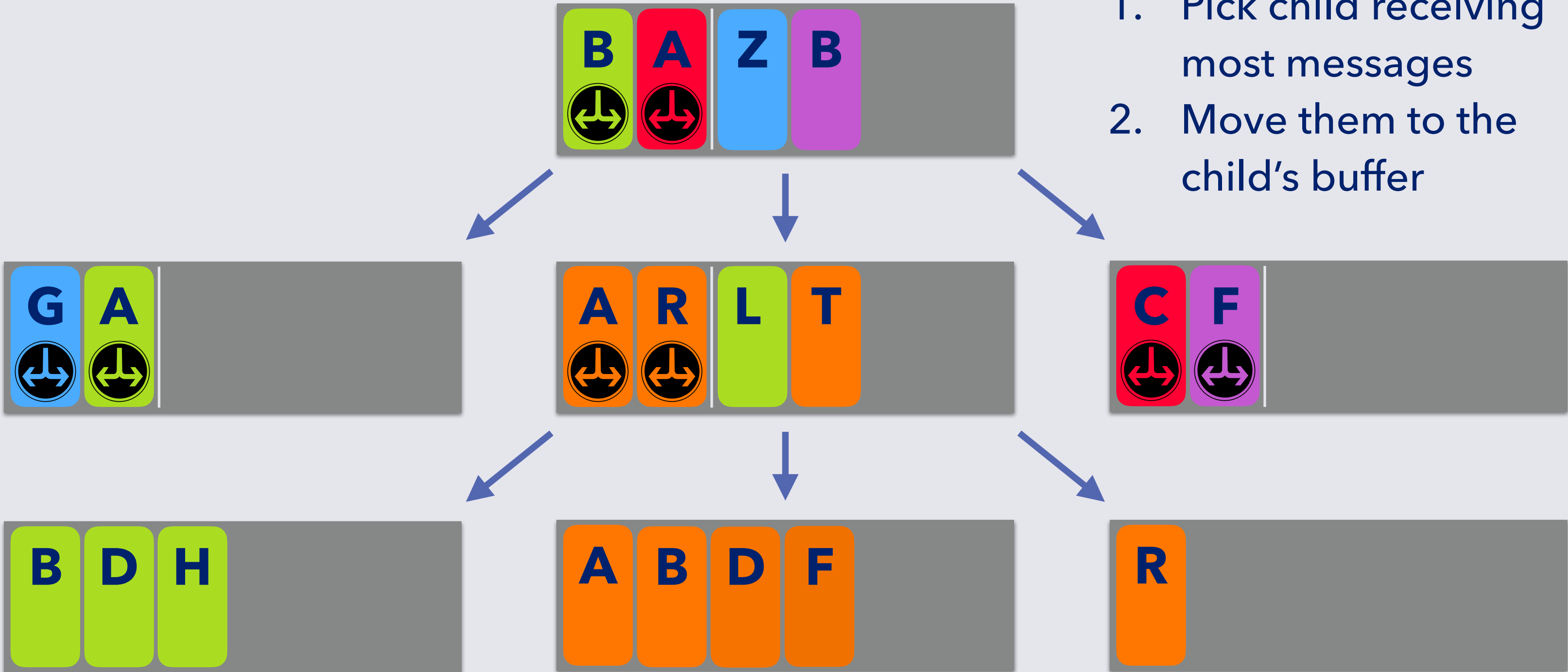
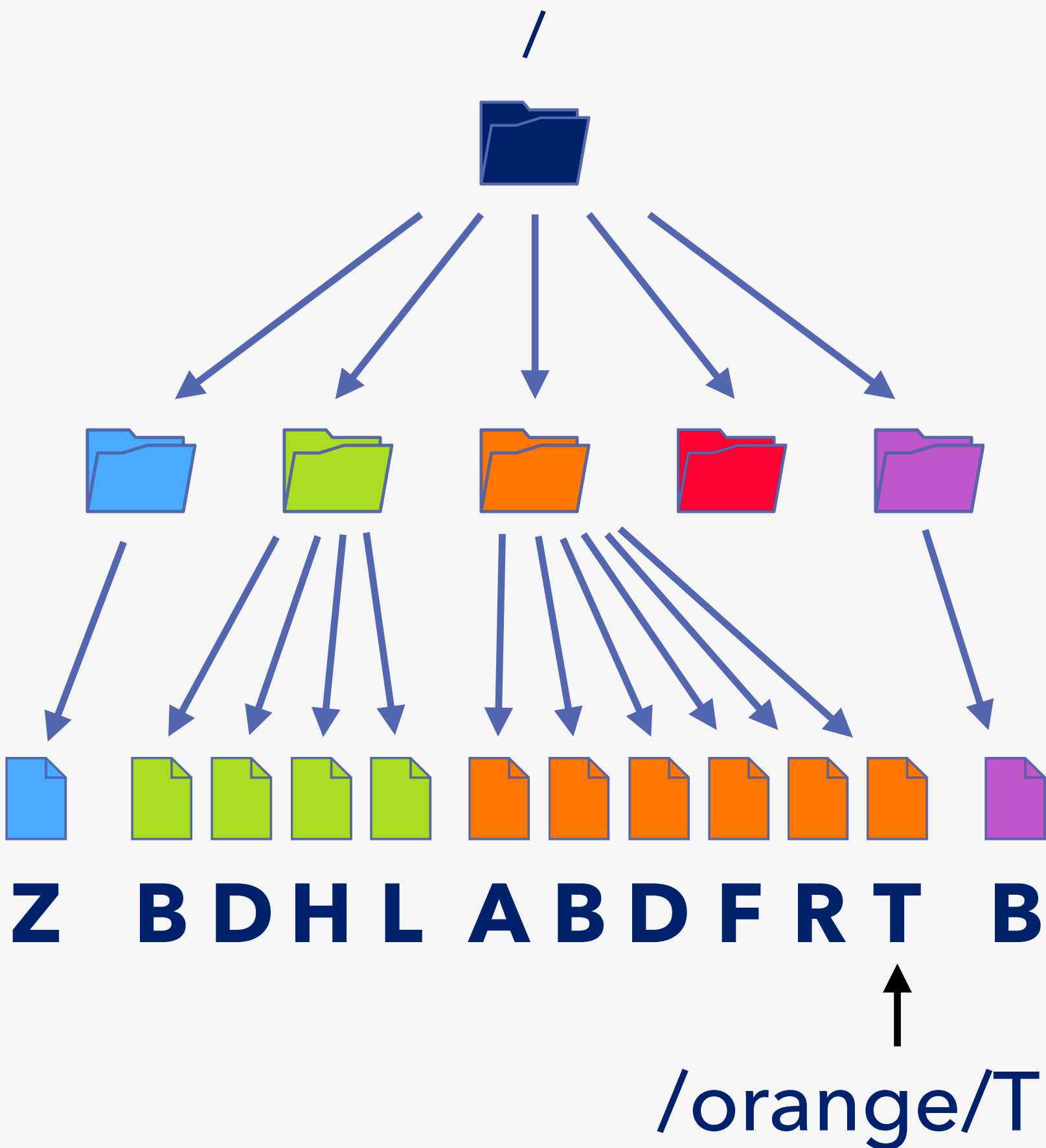
B ϵ -Trees

A B ϵ -tree is a search tree
(like a B-tree)

Inserts get put in
the root buffer

- When a buffer is full:
- 1. Pick child receiving most messages
 - 2. Move them to the child's buffer

directory tree

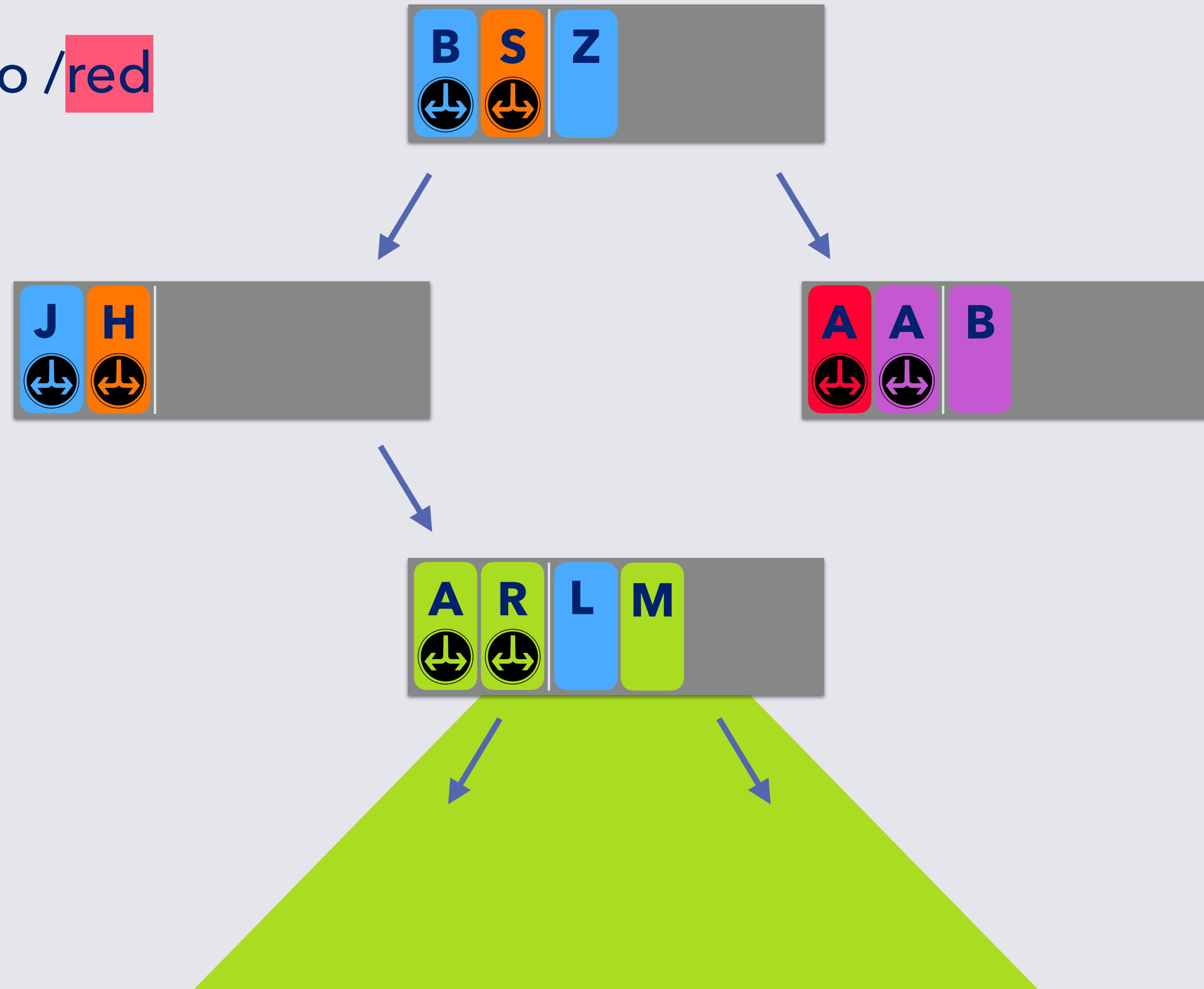


Key Insight: Each flush applies many small changes

Logical Copy with B^ϵ -DAGs

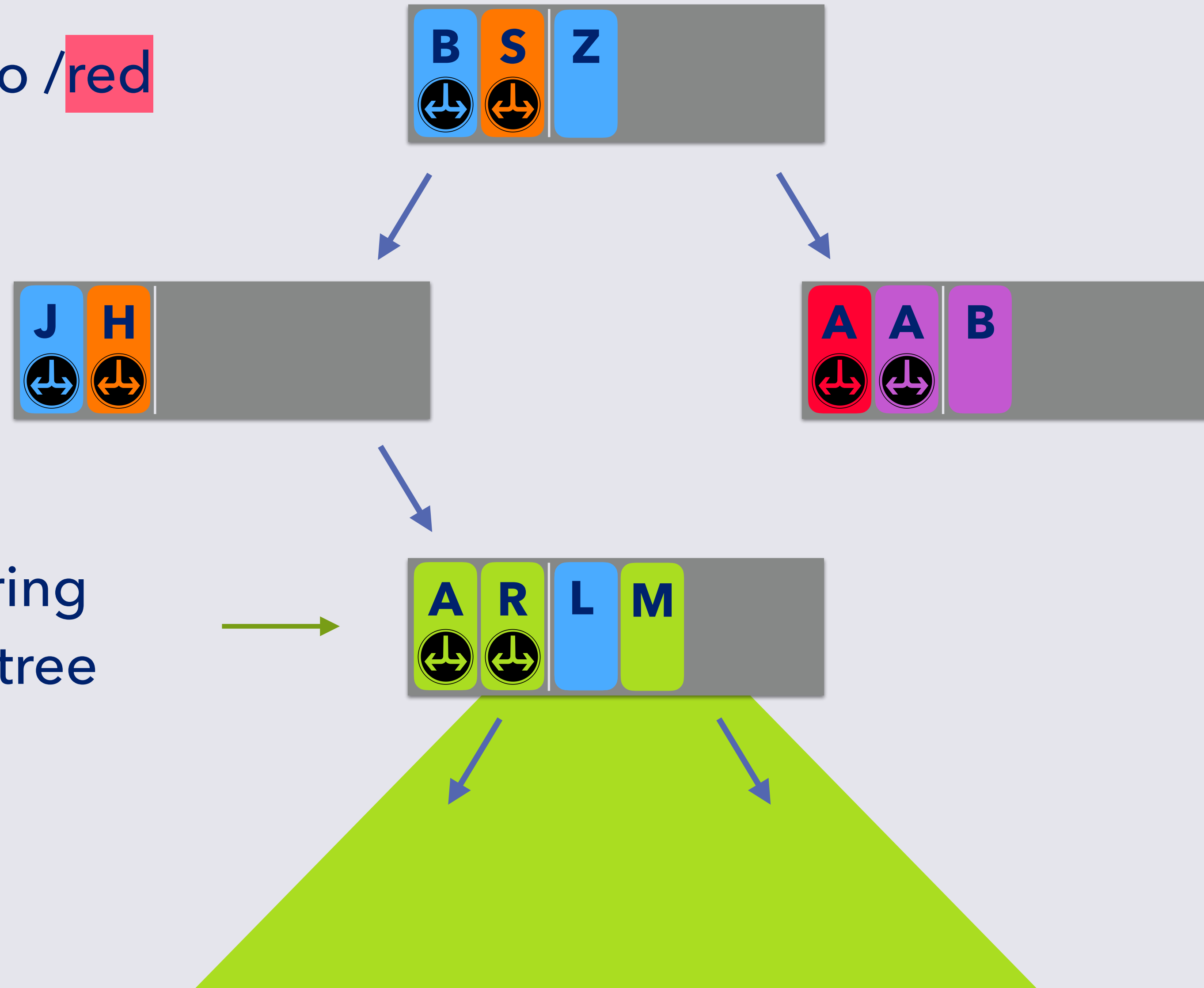
Logical Copy with B^ϵ -DAGs

Copy /green to /red



Logical Copy with B^ϵ -DAGs

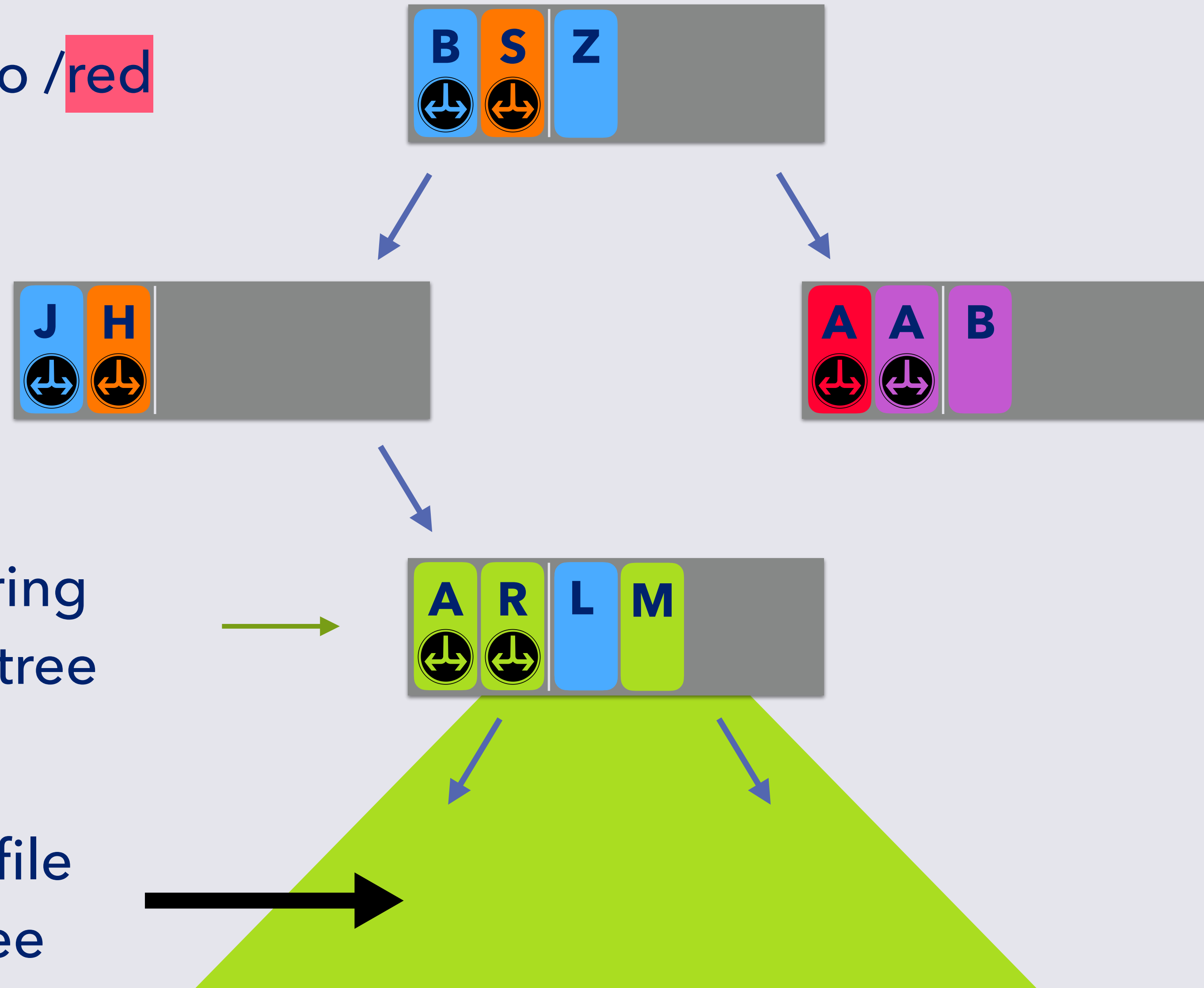
Copy /green to /red



Node covering
/green/ subtree

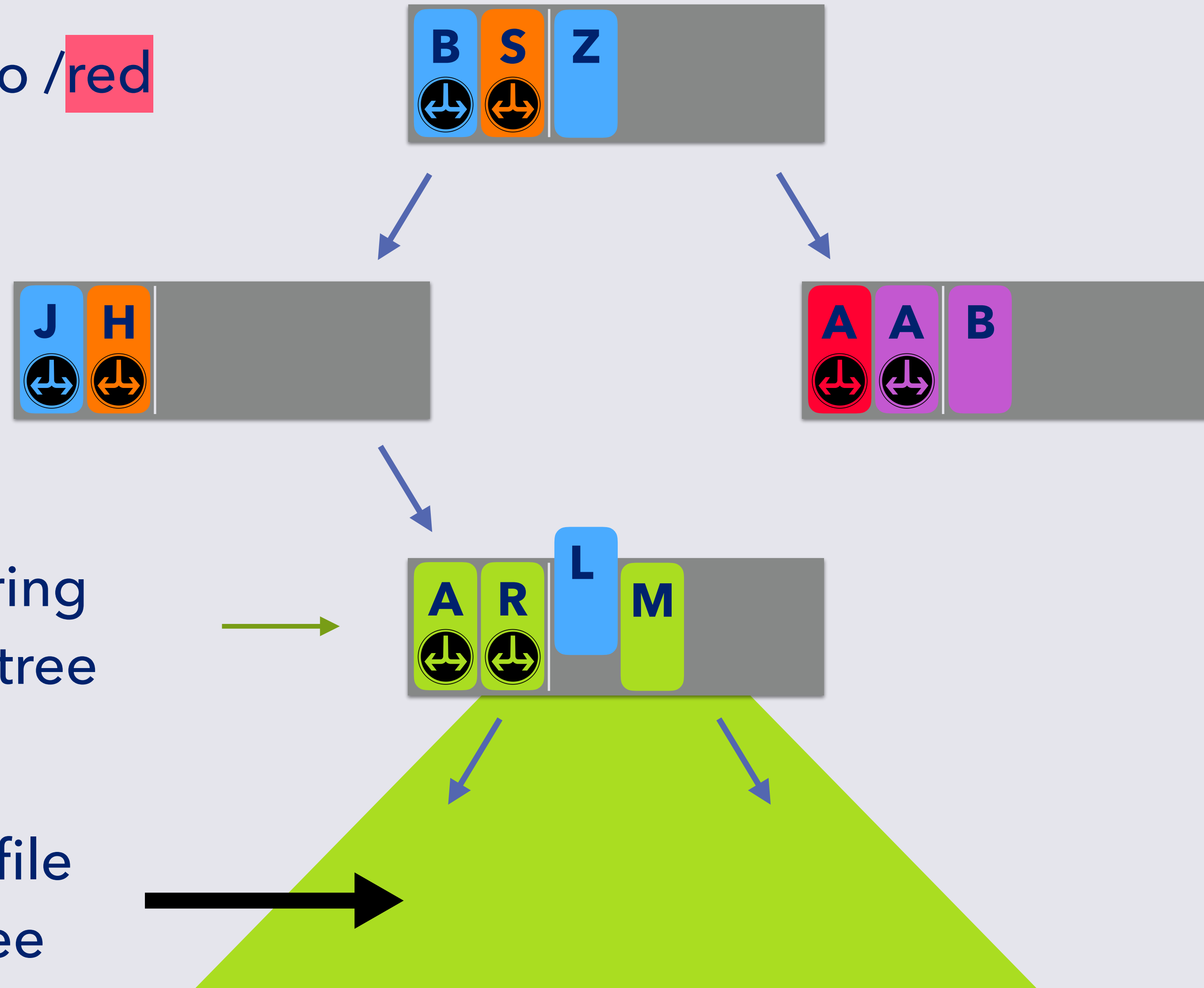
Logical Copy with B^ϵ -DAGs

Copy /green to /red



Logical Copy with B^ϵ -DAGs

Copy /green to /red

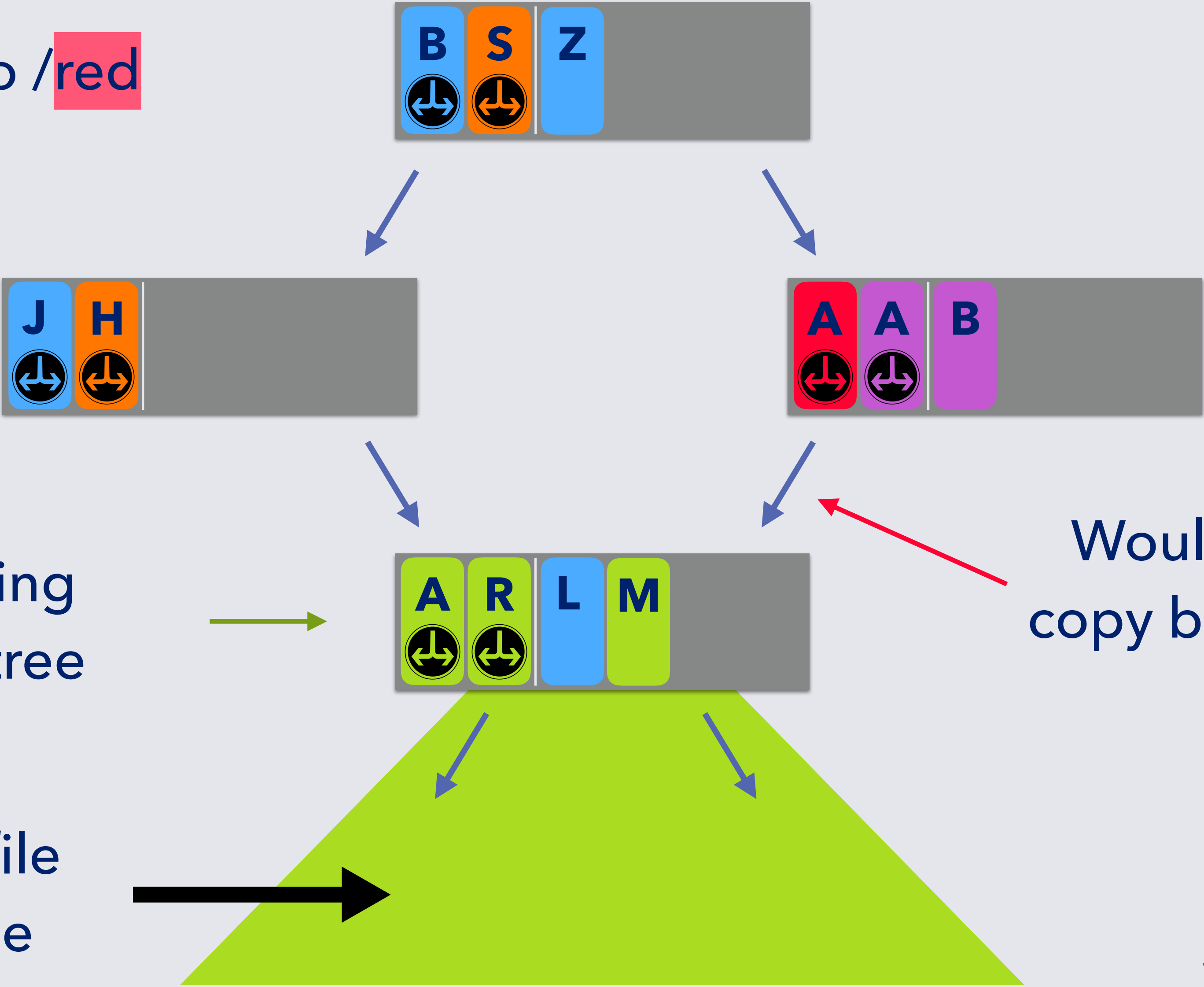


Node covering
/green/ subtree

Every /green/* file
is in this subtree

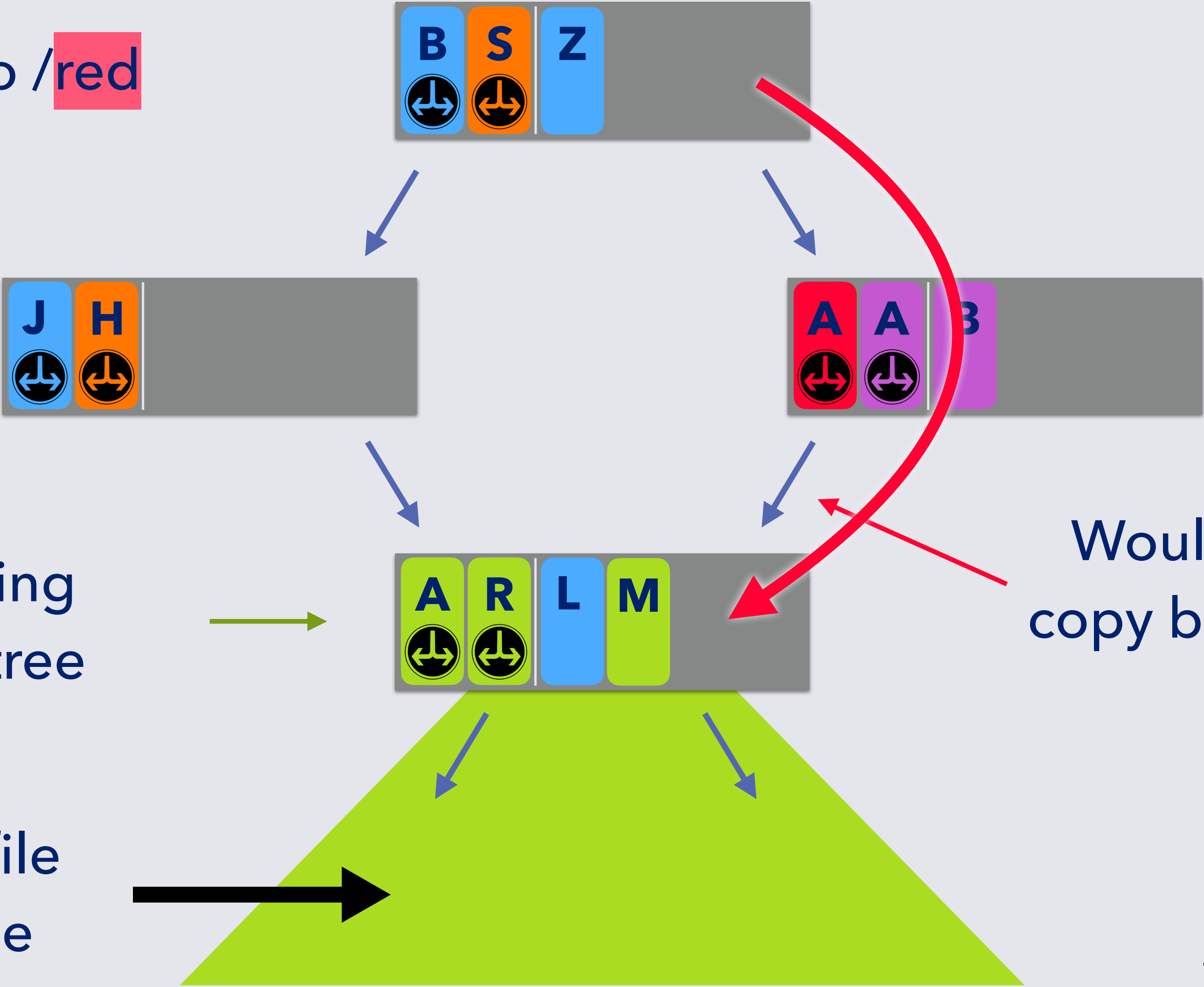
Logical Copy with B^ϵ -DAGs

Copy /green to /red



Logical Copy with B^ϵ -DAGs

Copy /green to /red

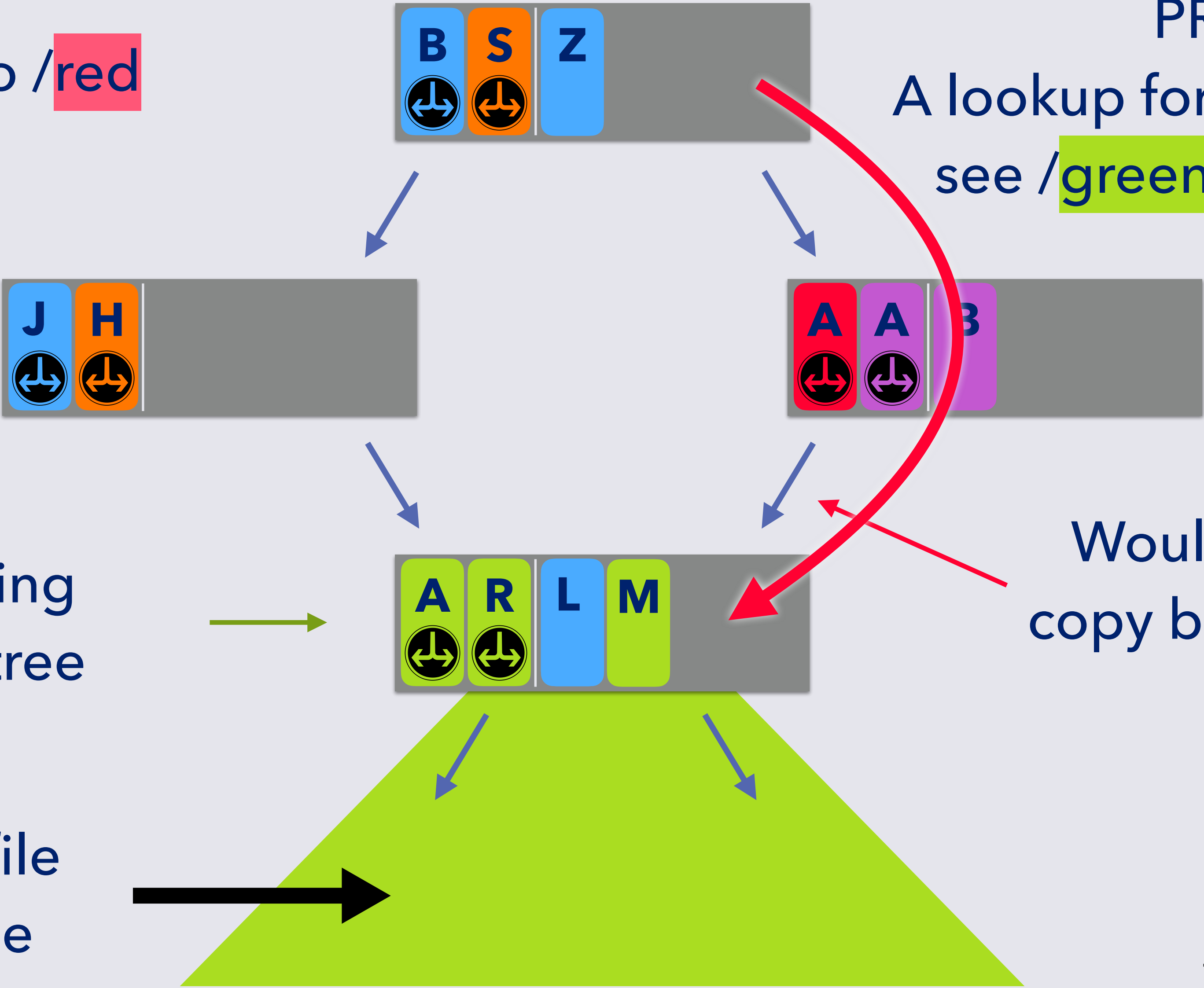


Logical Copy with B^ϵ -DAGs

Copy /green to /red

PROBLEM:

A lookup for /red/M is going to see /green/M in the subtree

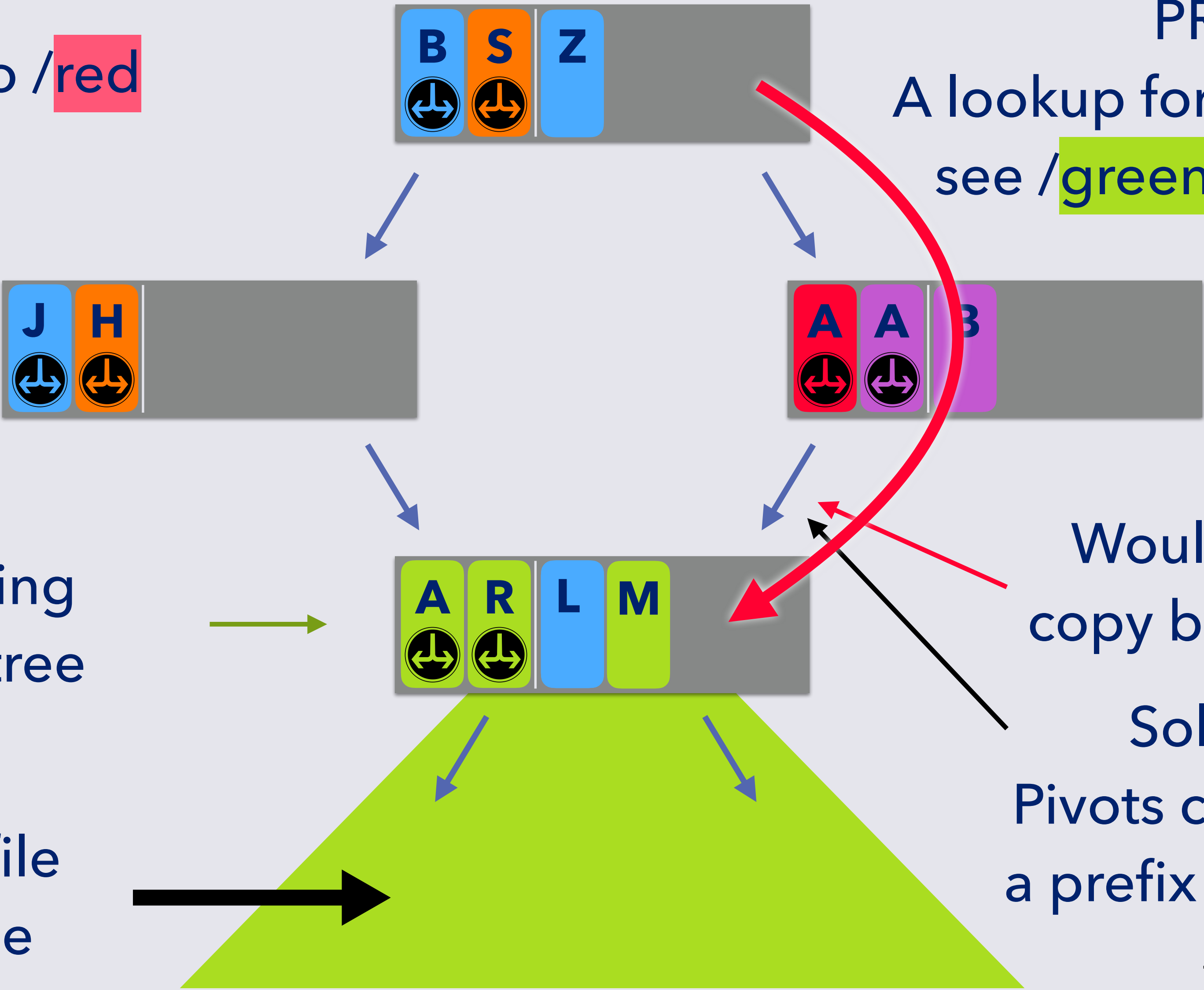


Logical Copy with B^ϵ -DAGs

Copy /green to /red

PROBLEM:

A lookup for /red/M is going to see /green/M in the subtree



Node covering /green/ subtree

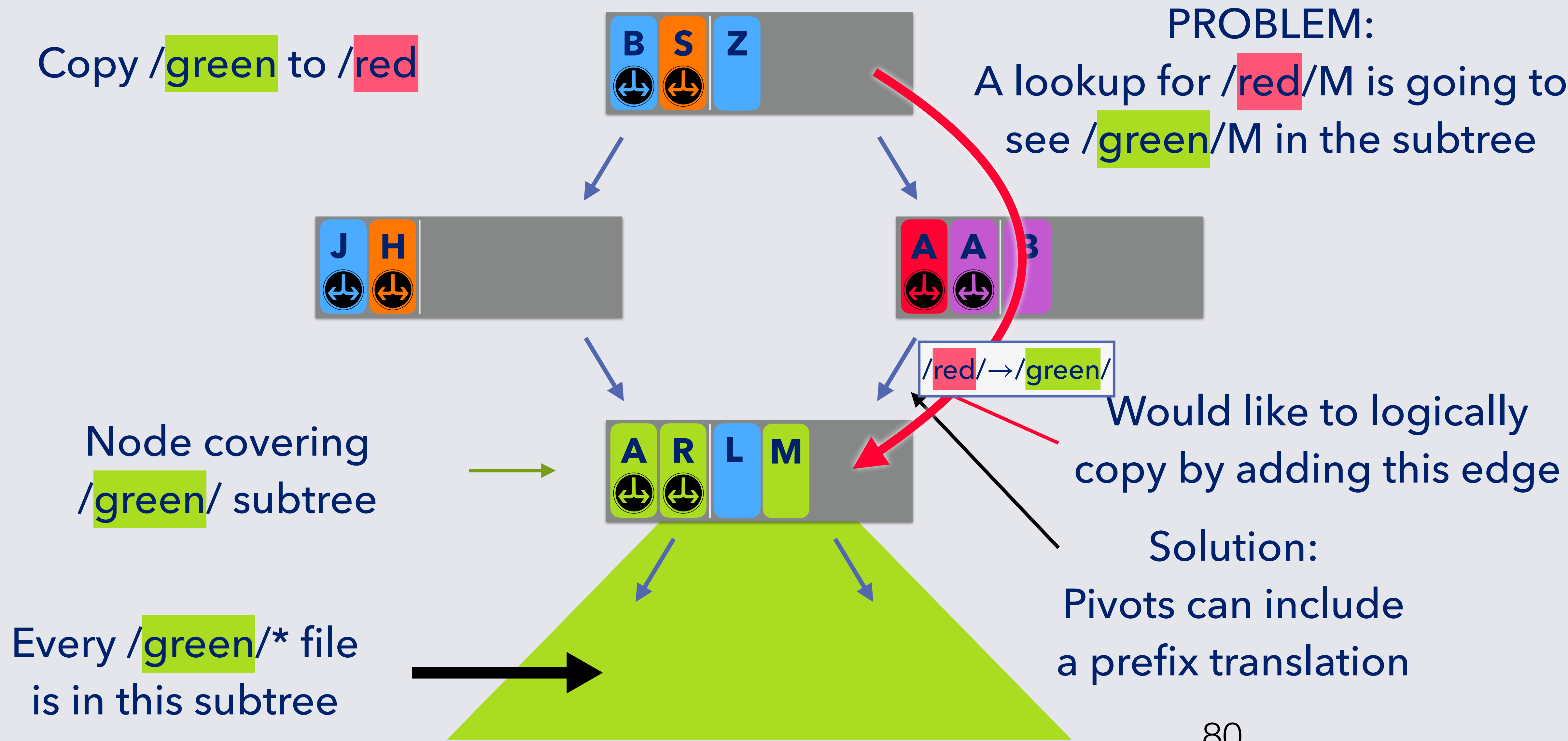
Would like to logically copy by adding this edge

Solution:

Pivots can include a prefix translation

Every /green/* file is in this subtree

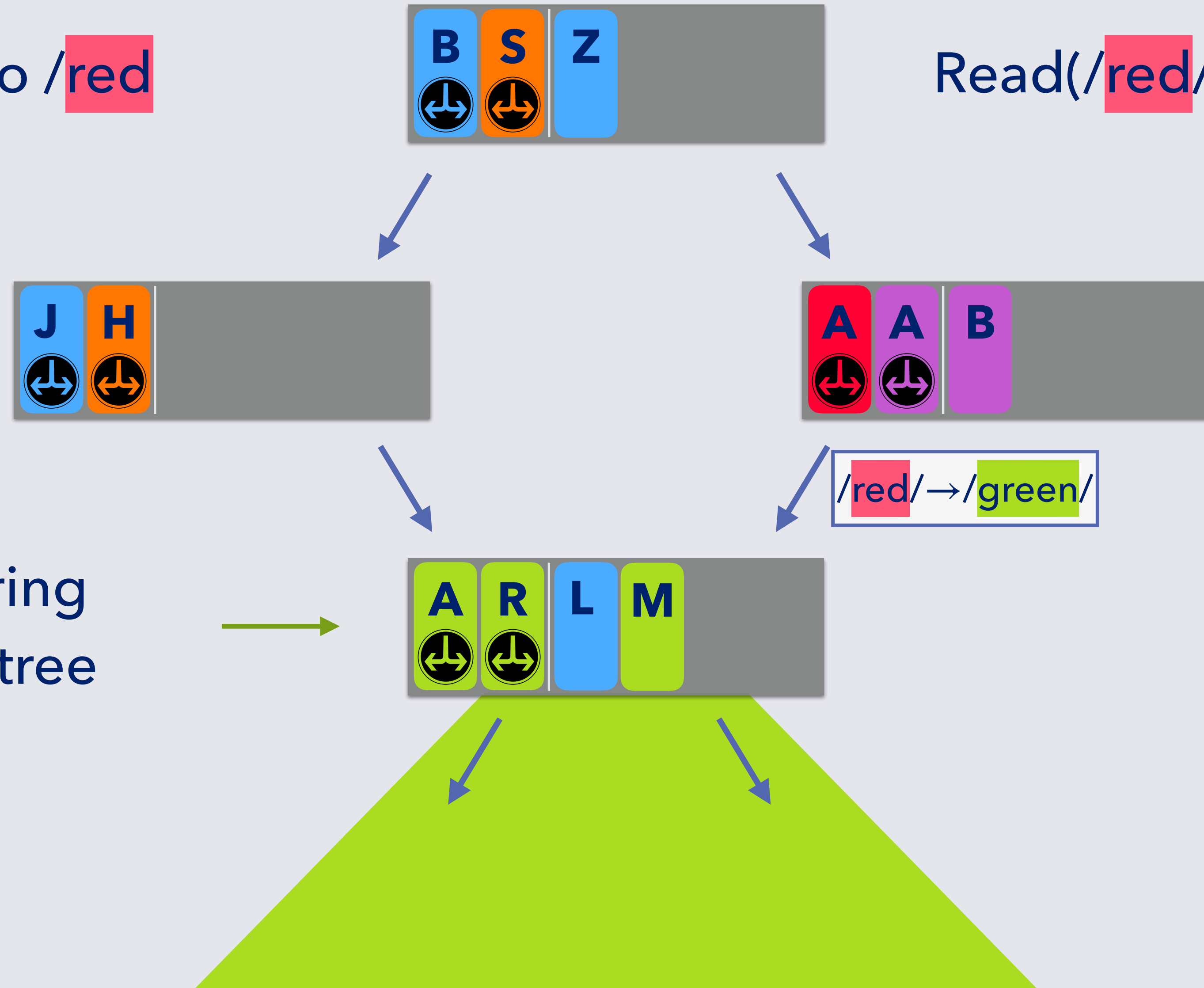
Logical Copy with B^ϵ -DAGs



Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/

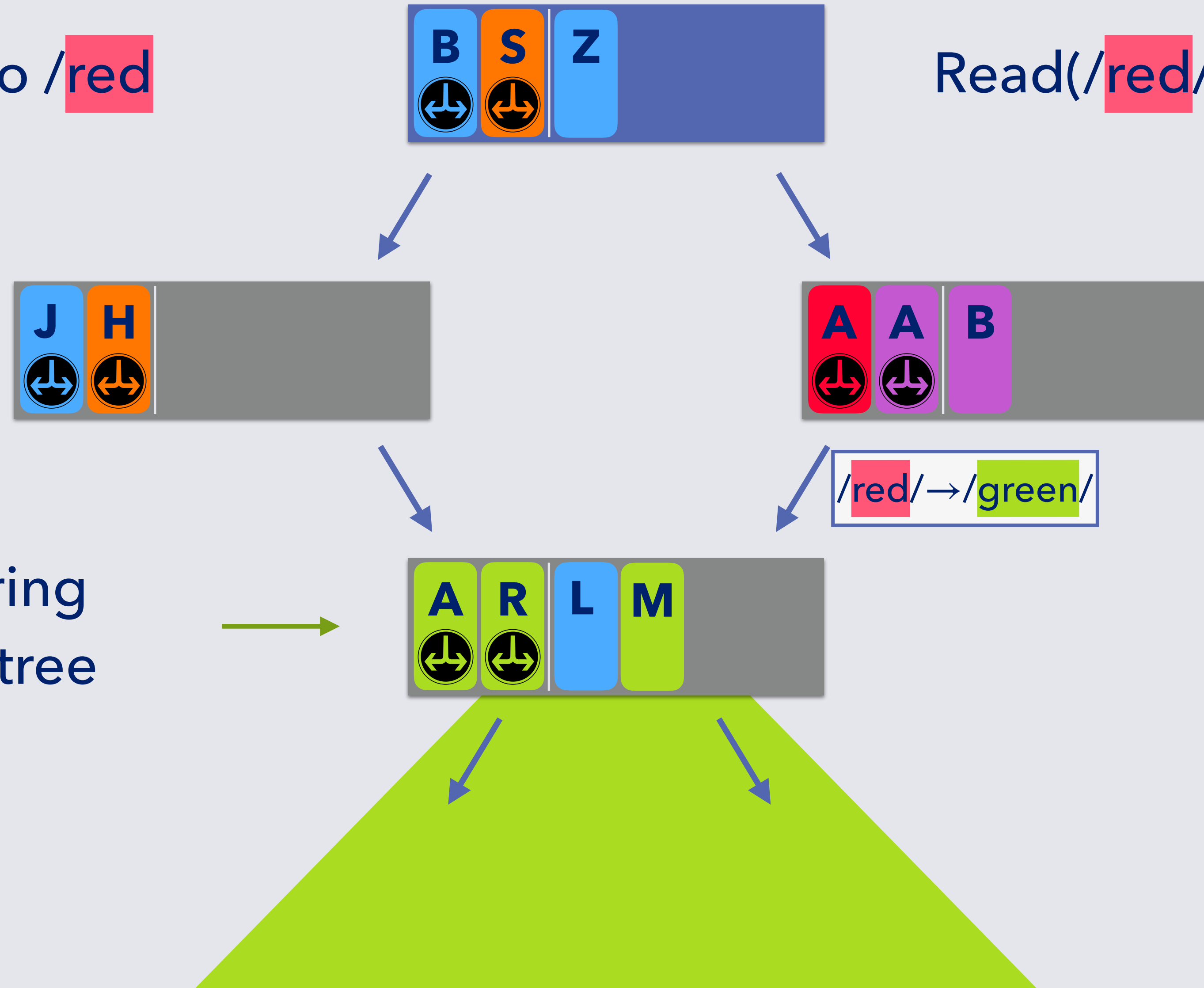
Read(/red/M)



Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/

Read(/red/M)

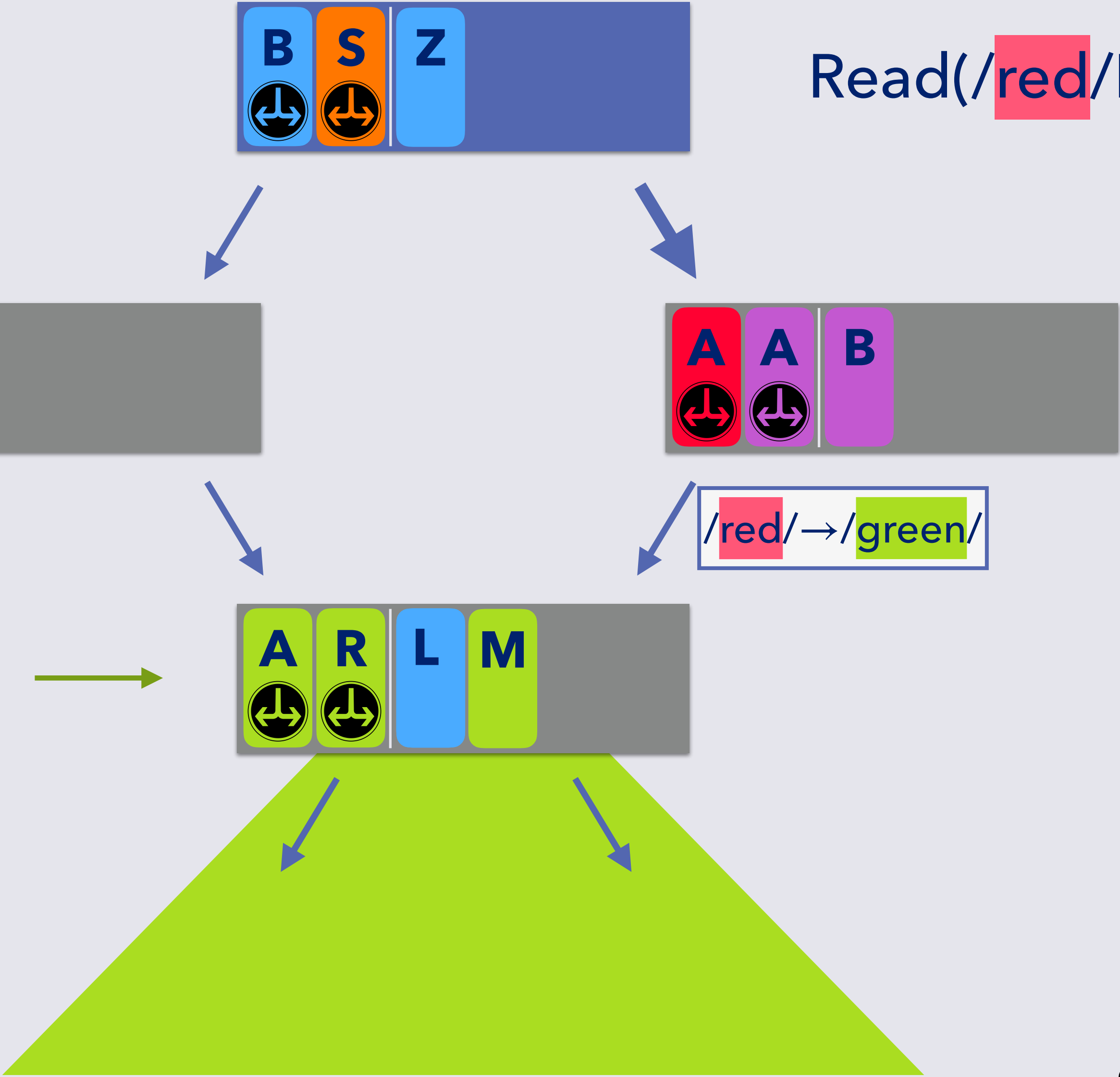


Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/

Read(/red/M)

Node covering
/green/ subtree

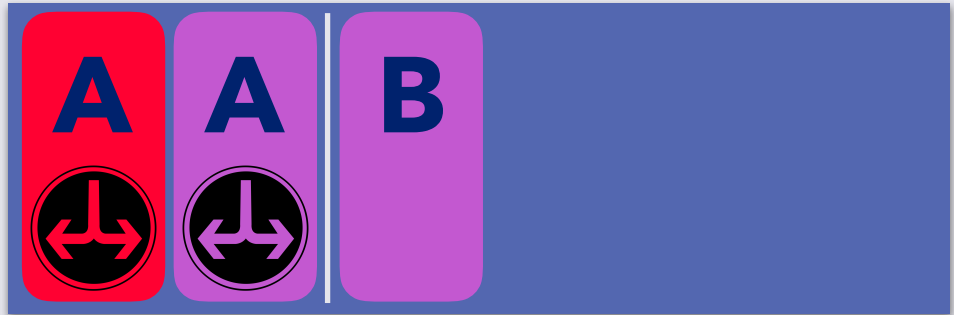


Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/



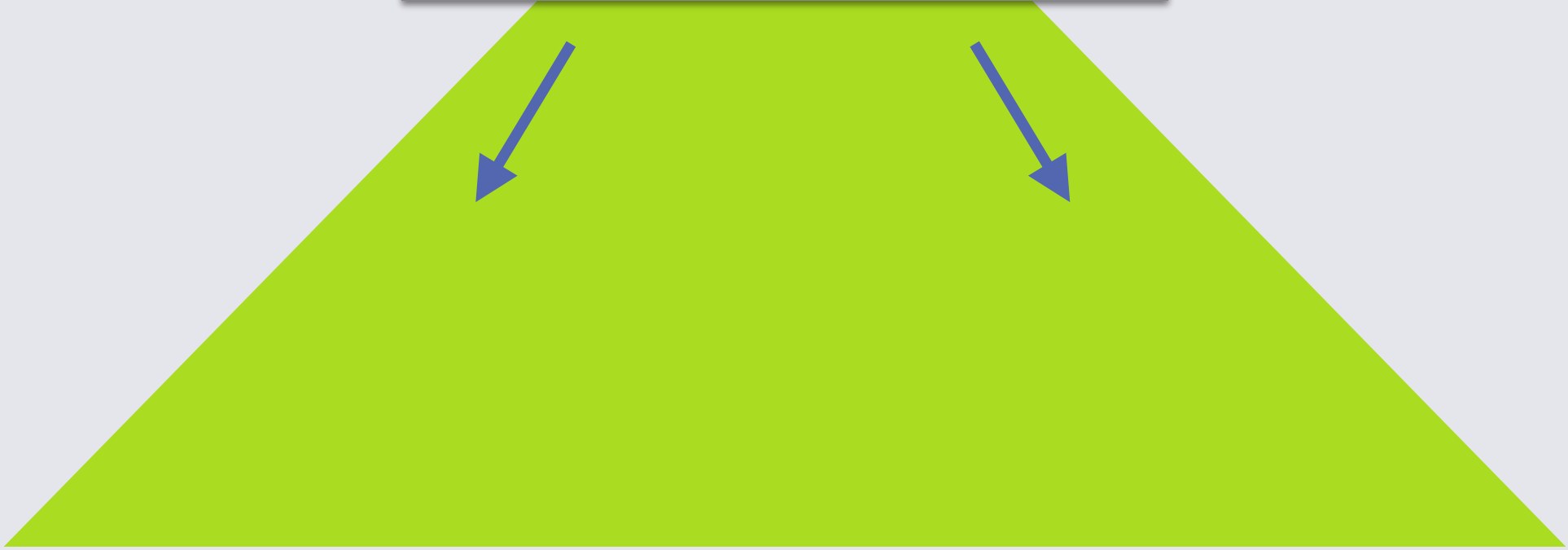
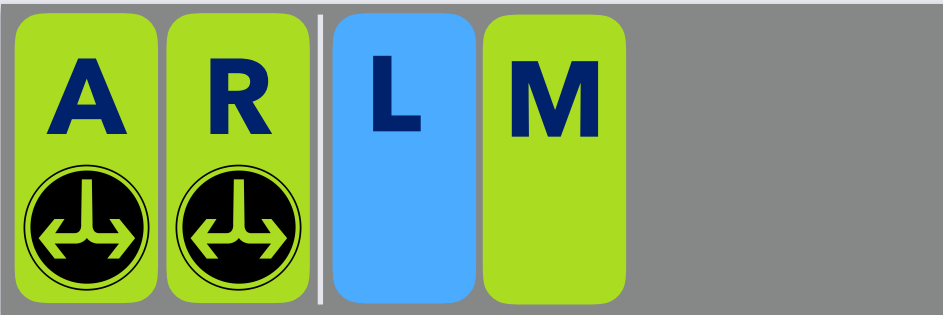
Read(/red/M)



Read(/red/M)

/red/ → /green/

Node covering
/green/ subtree

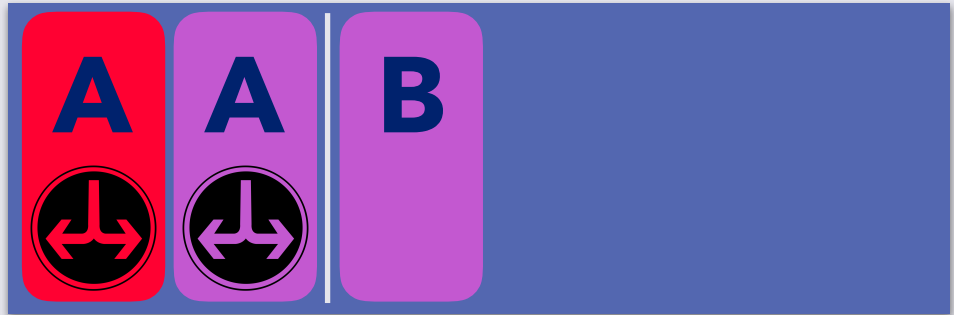


Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/



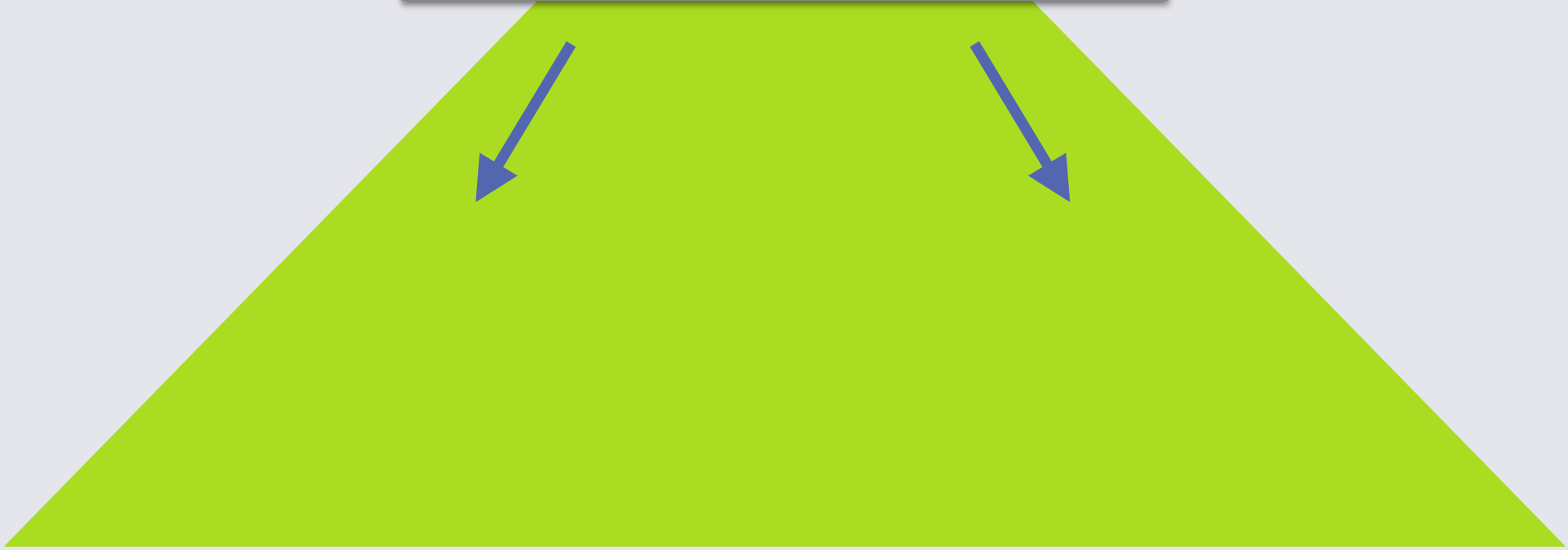
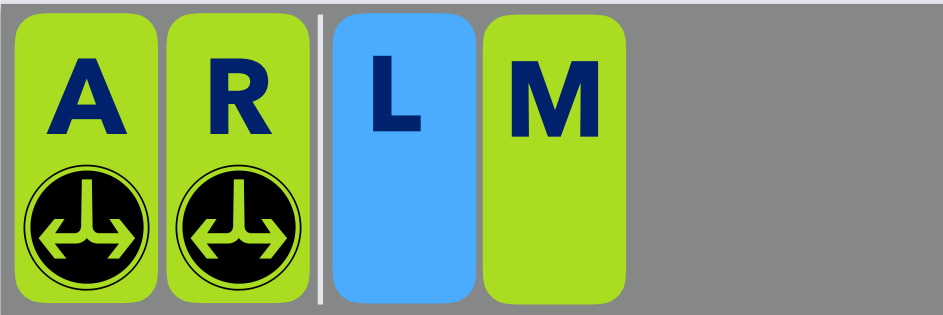
Read(/red/M)



Read(/red/M)

/red/ → /green/

Node covering
/green/ subtree

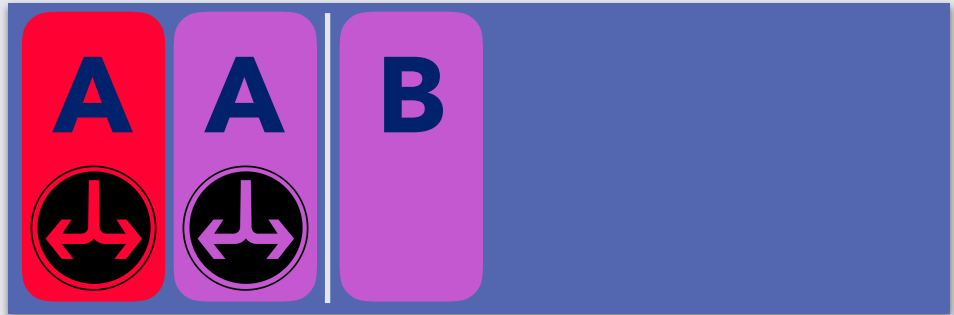


Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/



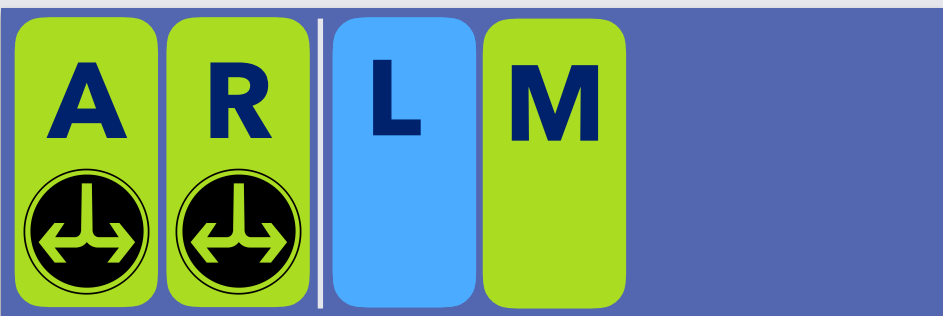
Read(/red/M)



Read(/red/M)



Node covering /green/ subtree



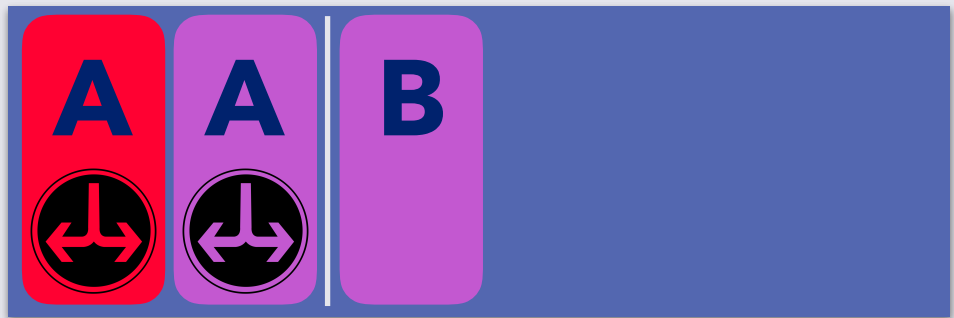
Read(/green/M)

Logical Copy with B^ϵ -DAGs

Copy /green/ to /red/



Read(/red/M)



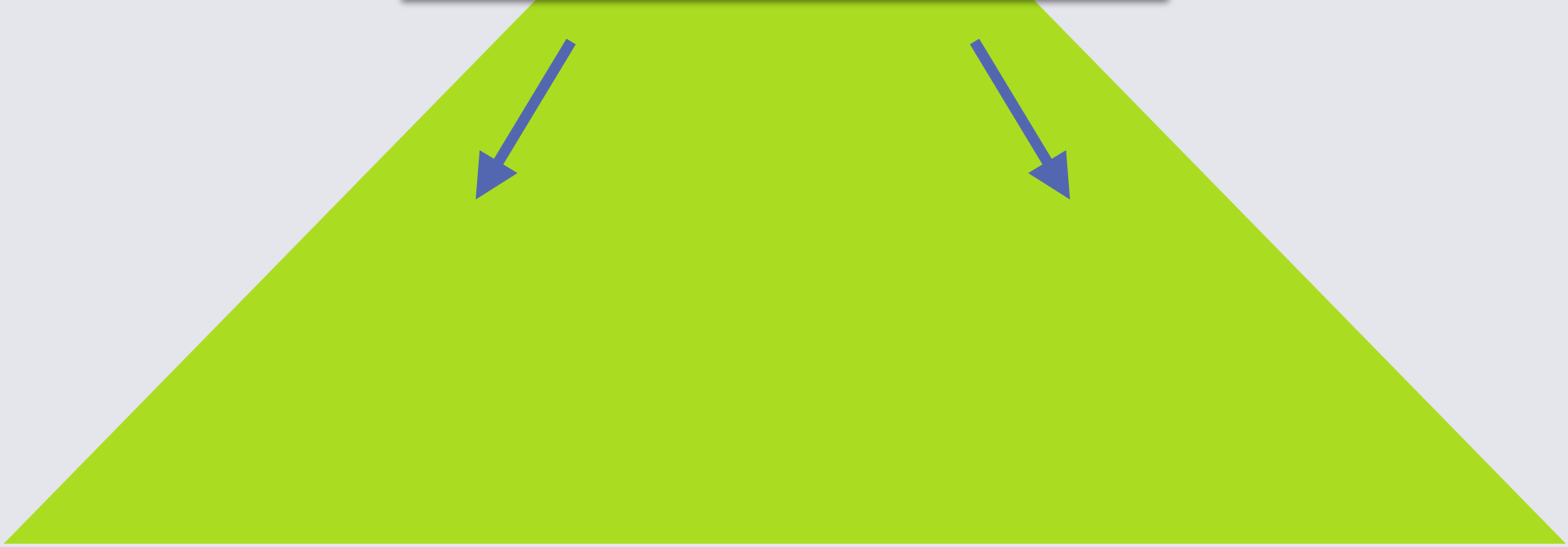
Read(/red/M)



Node covering /green/ subtree



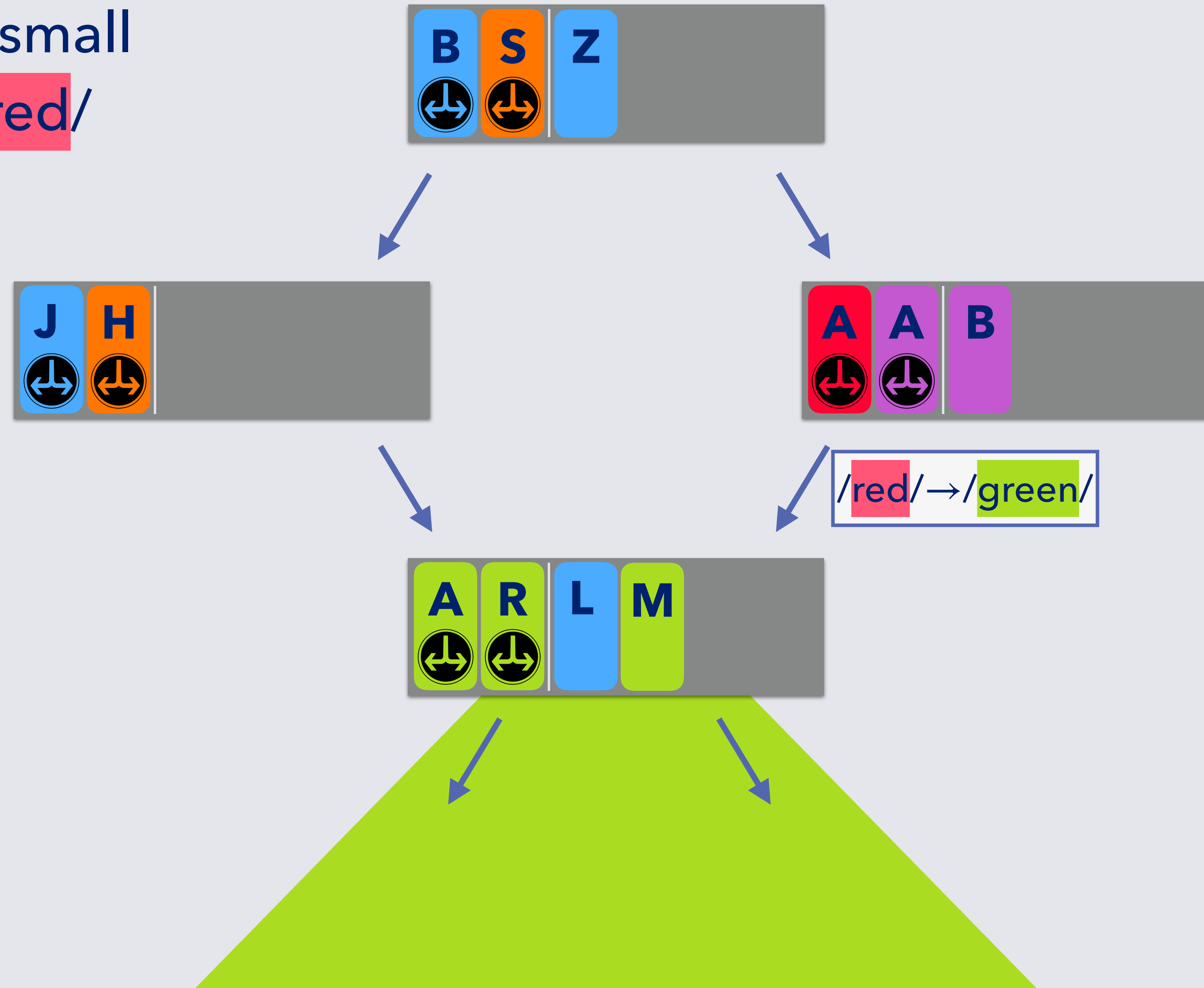
Read(/green/M)



B_ε -DAGs and Small Writes

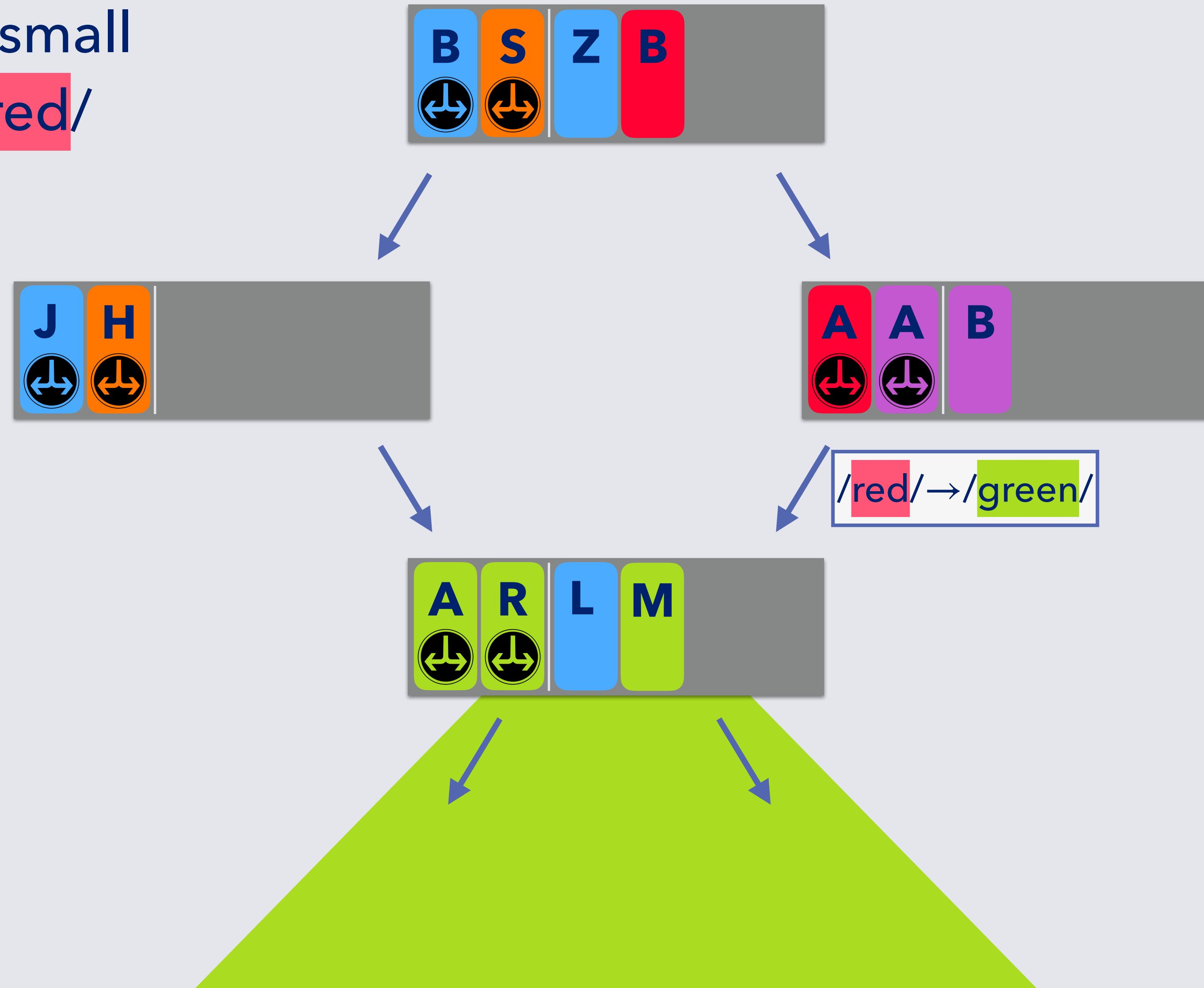
B ϵ -DAGs and Small Writes

Make some small
writes to /red/



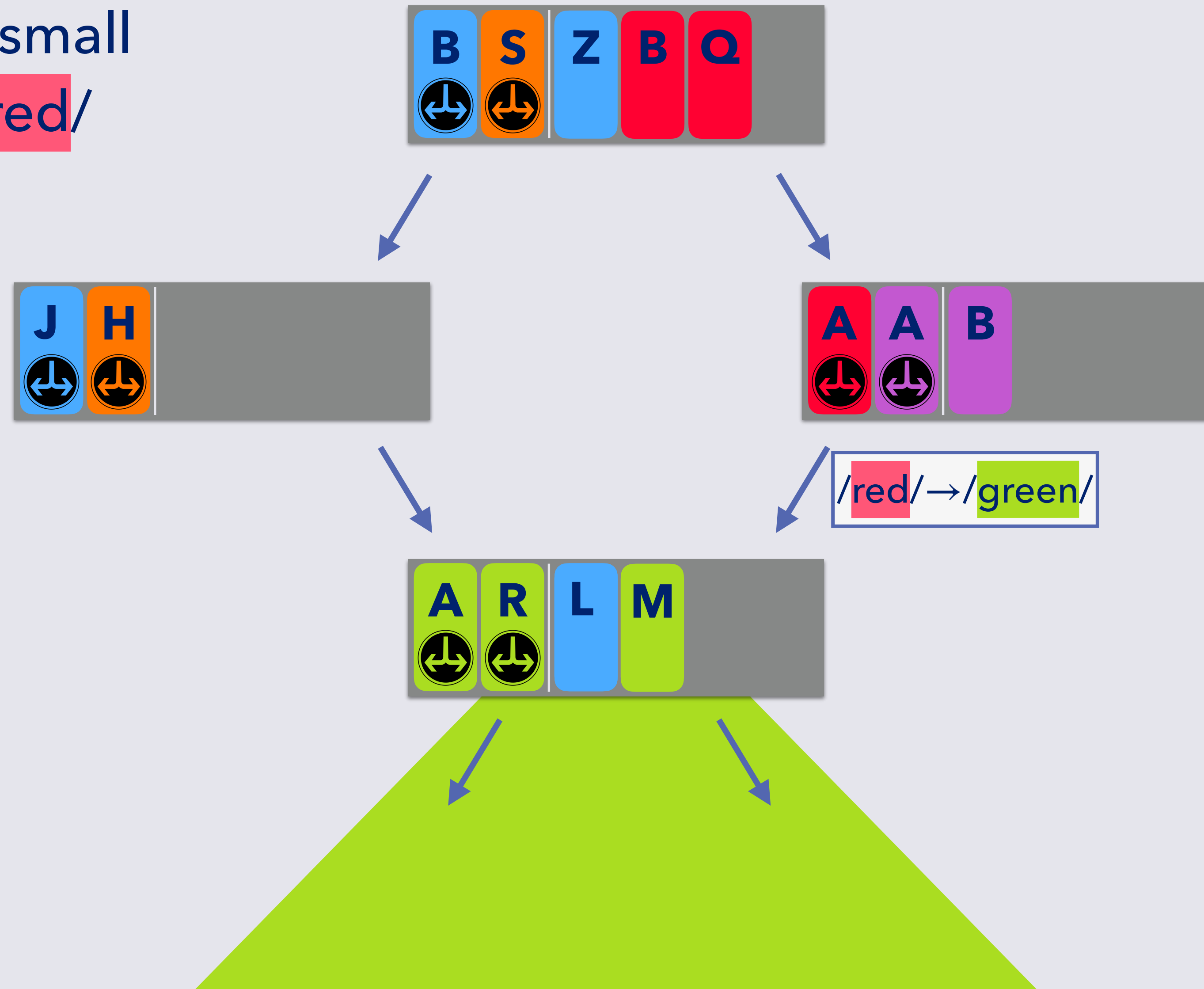
B ϵ -DAGs and Small Writes

Make some small
writes to /red/



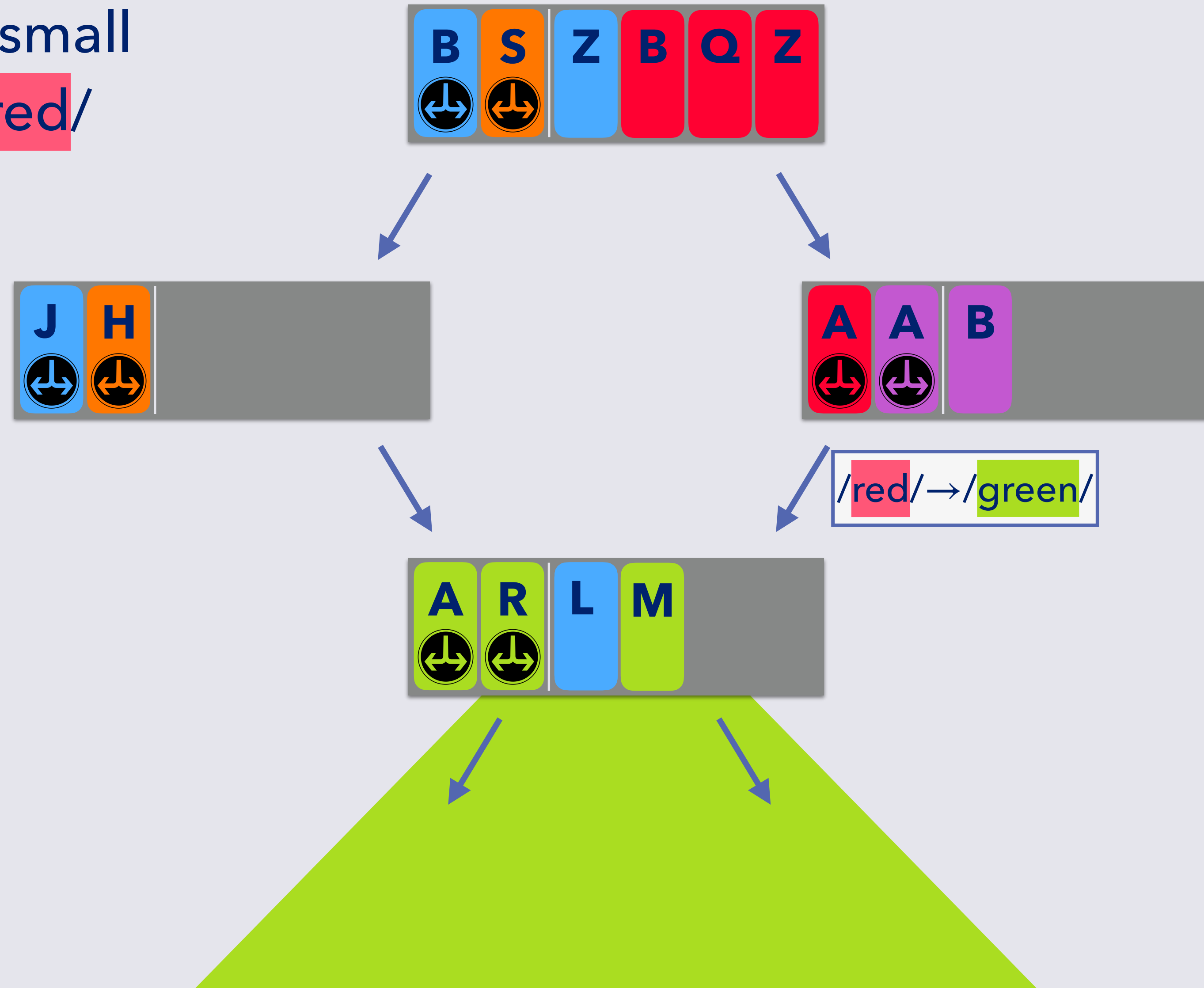
B $^\epsilon$ -DAGs and Small Writes

Make some small
writes to /red/



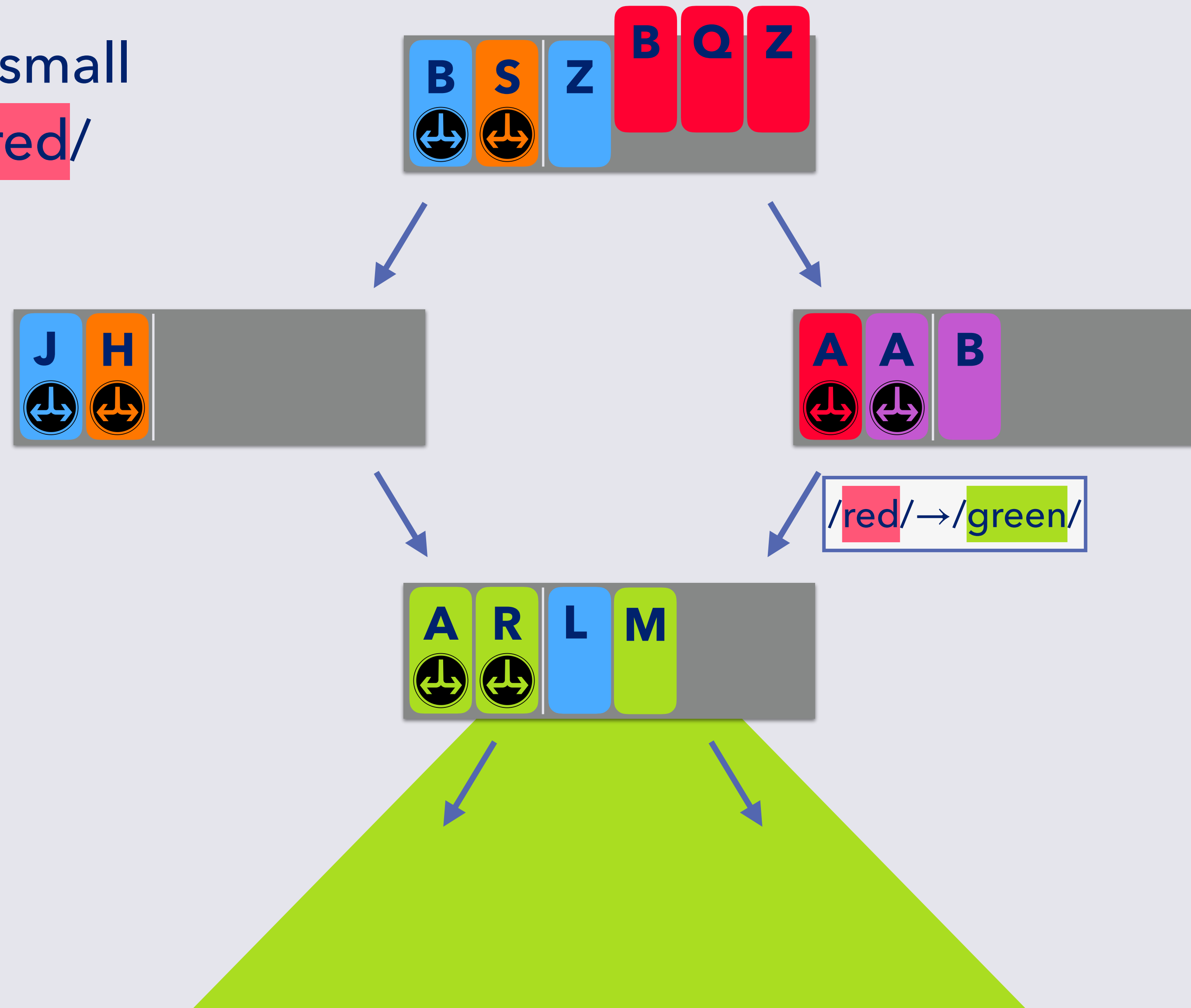
B ϵ -DAGs and Small Writes

Make some small
writes to /red/



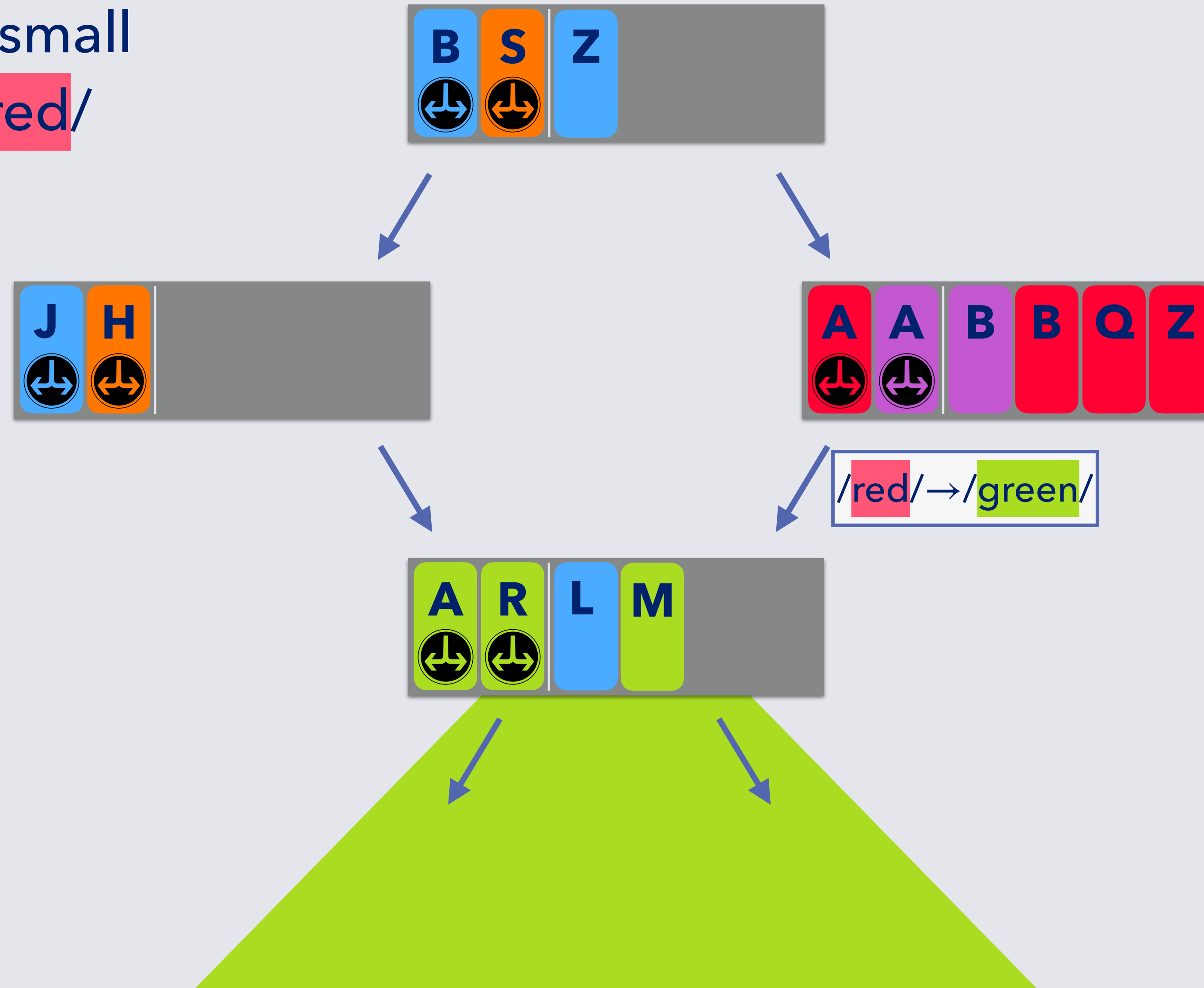
B ϵ -DAGs and Small Writes

Make some small
writes to /red/



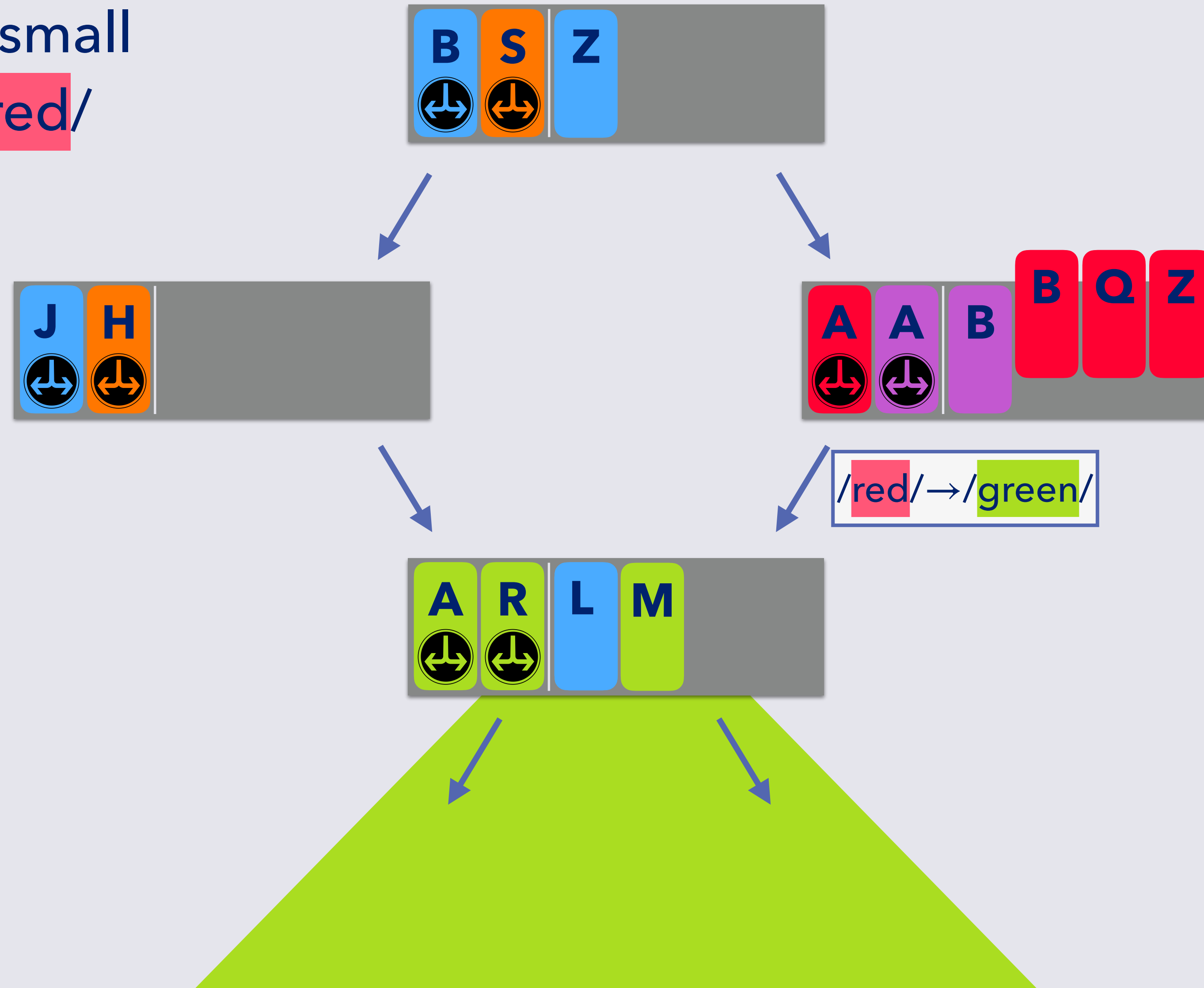
B ϵ -DAGs and Small Writes

Make some small
writes to /red/



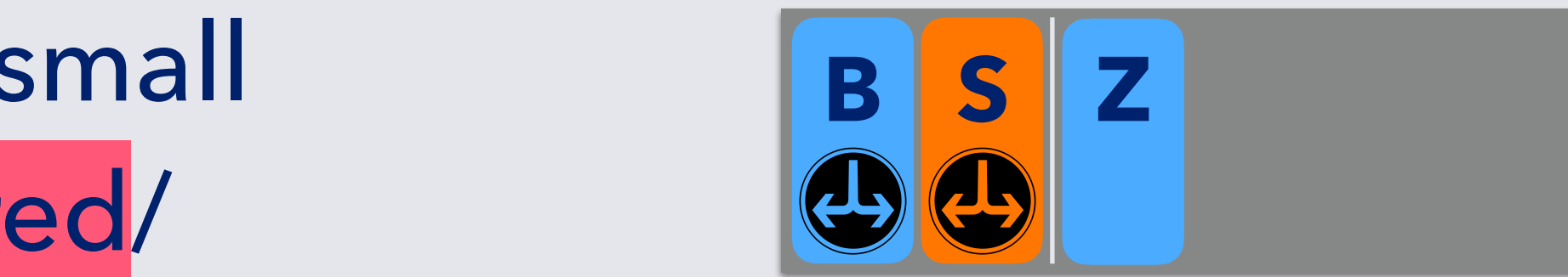
B ϵ -DAGs and Small Writes

Make some small
writes to /red/

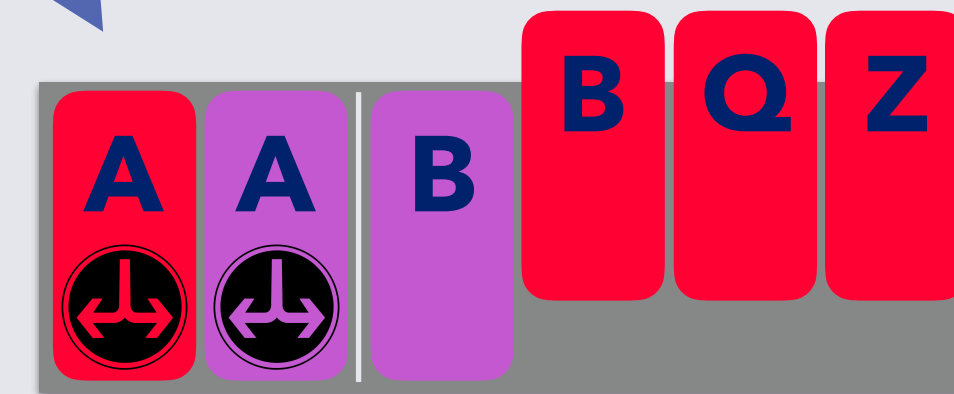


B ϵ -DAGs and Small Writes

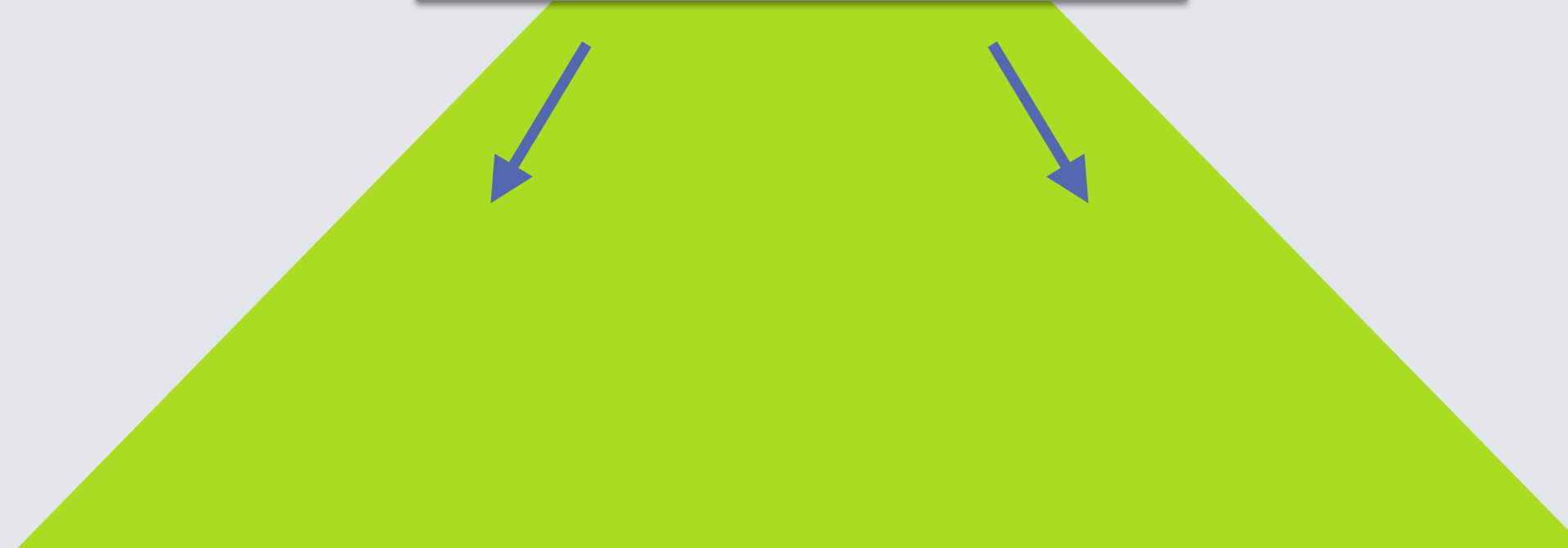
Make some small
writes to /red/



Now flush with
copy-on-write

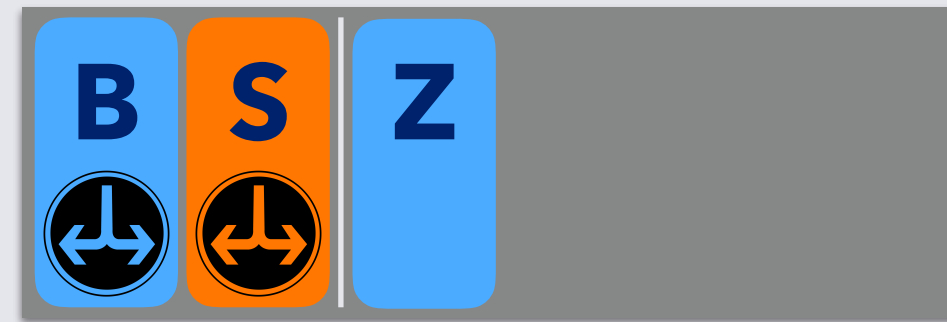


/red/ → /green/

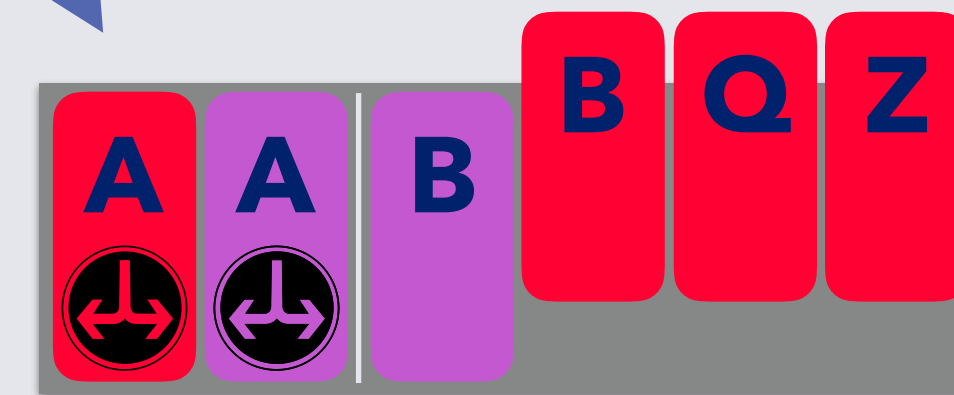


B ϵ -DAGs and Small Writes

Make some small
writes to /red/



Now flush with
copy-on-write



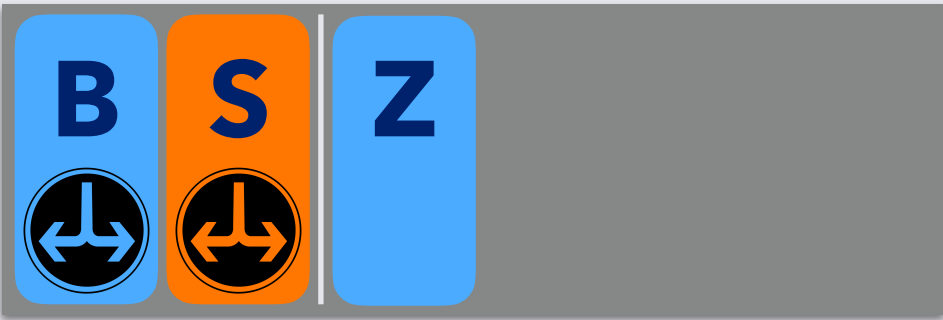
/red/ → /green/

1. copy

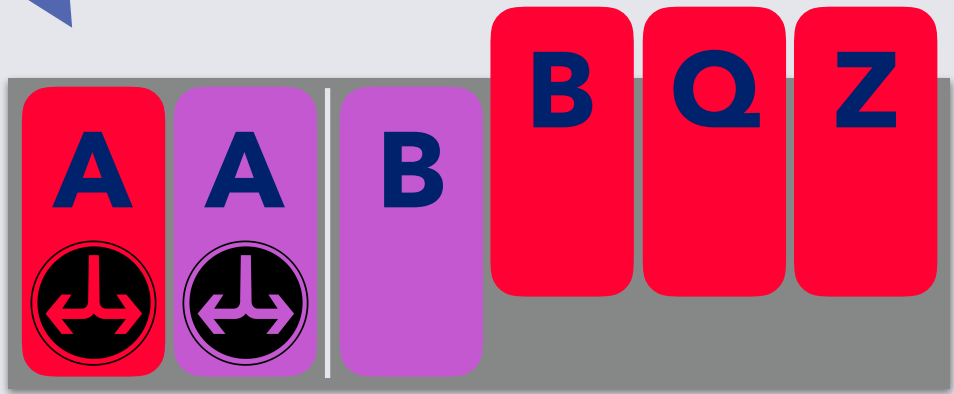


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

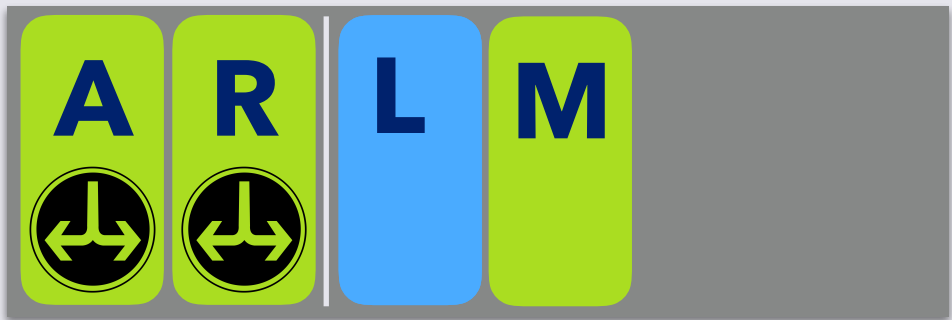
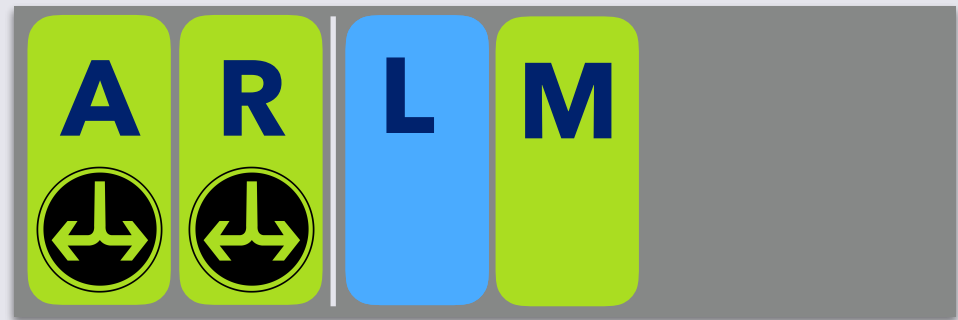


Now flush with
copy-on-write



/red/ → /green/

1. copy

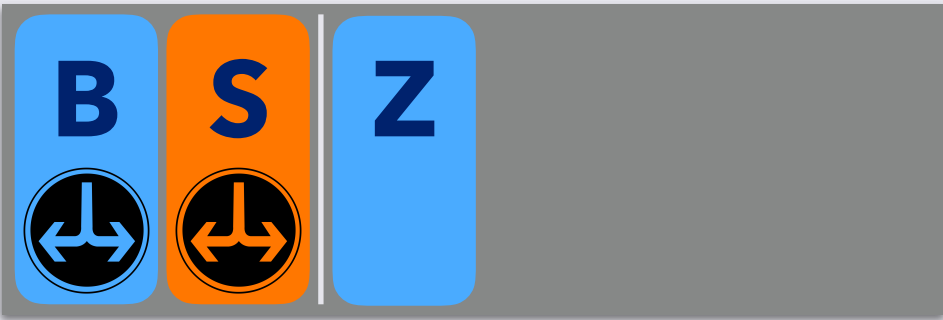


2. Translate
/green/ → /red/

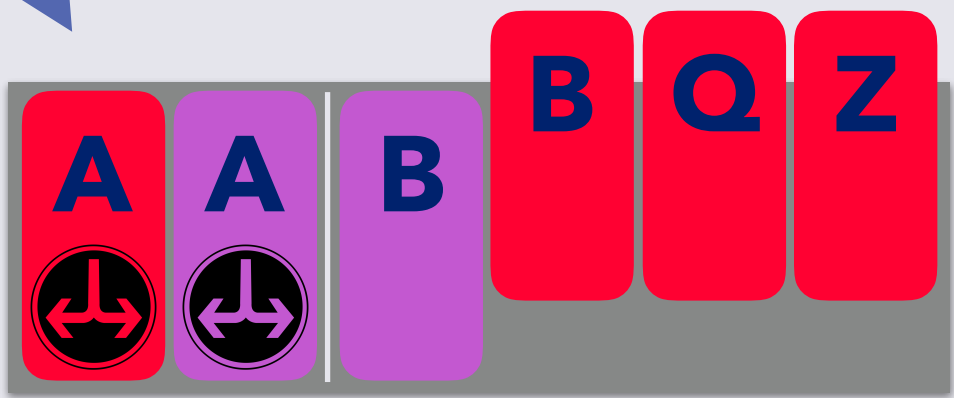


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

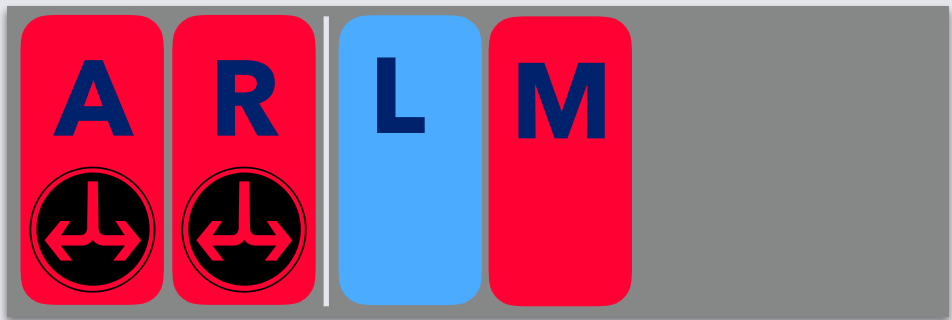
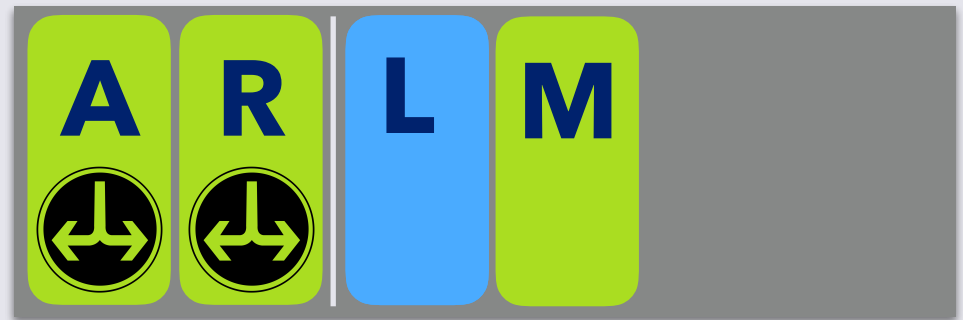


Now flush with
copy-on-write



/red/ → /green/

1. copy

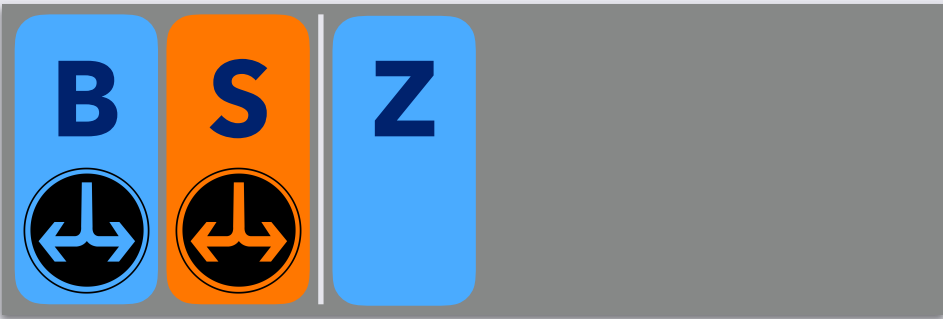


2. Translate
/green/ → /red/

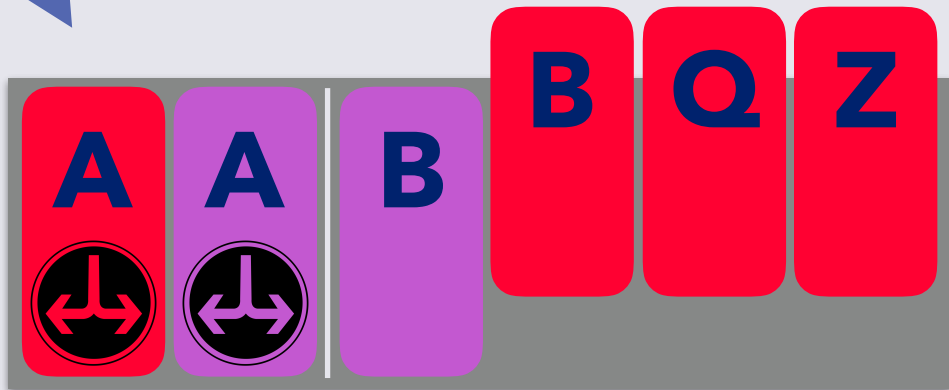


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

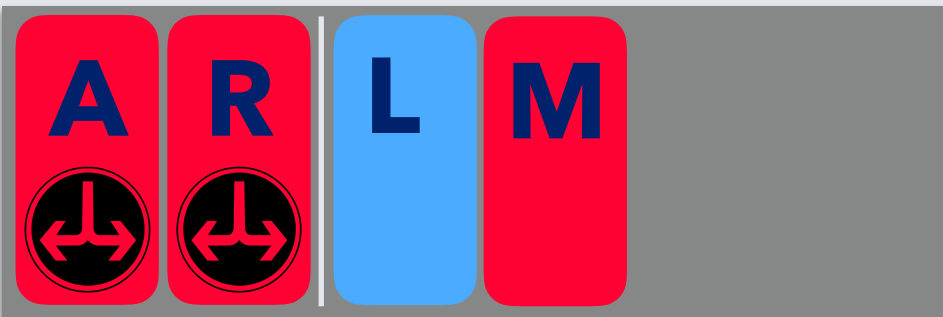


Now flush with
copy-on-write



/red/ → /green/

1. copy



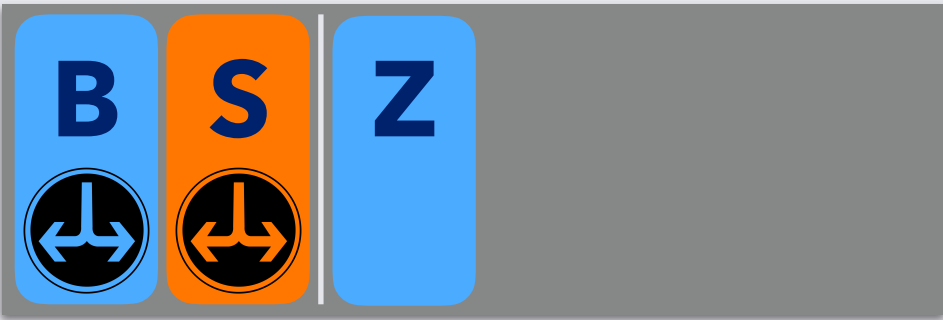
2. Translate
/green/ → /red/

3. Delete
unreachable data

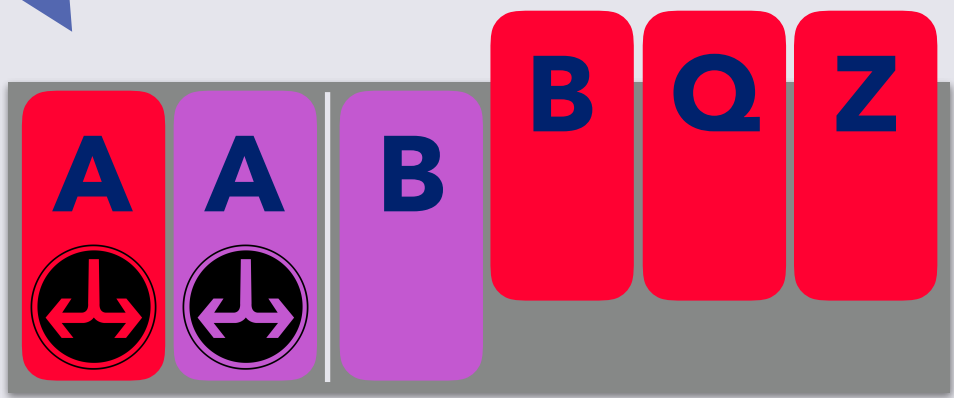


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

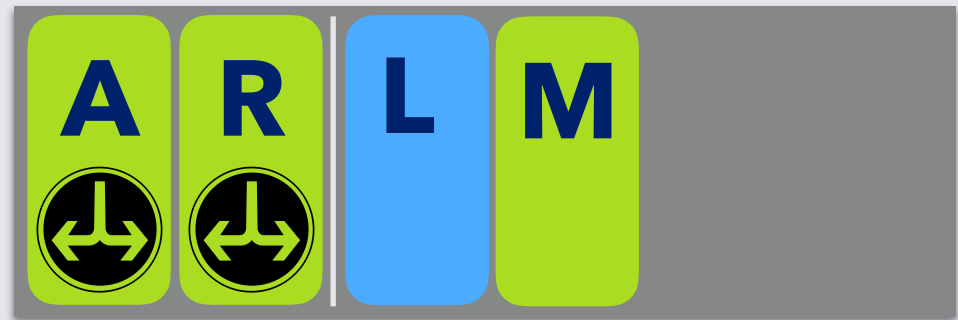


Now flush with
copy-on-write



/red/ → /green/

1. copy



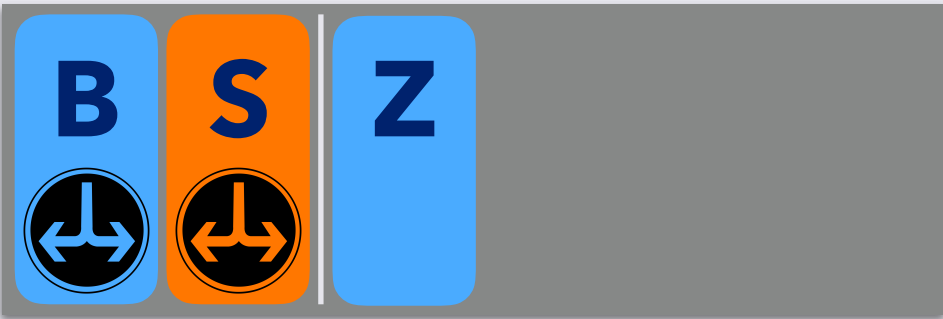
2. Translate
/green/ → /red/

3. Delete
unreachable data

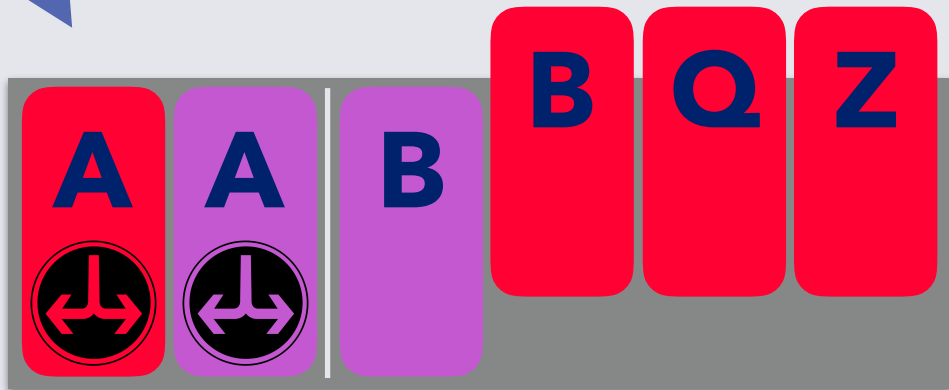


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

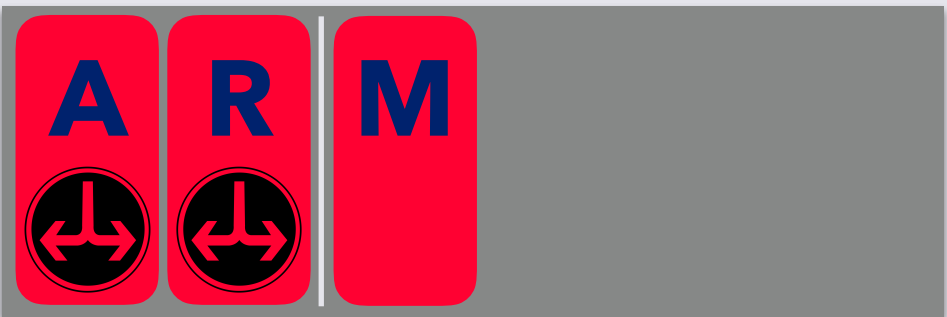


Now flush with
copy-on-write



/red/ → /green/

1. copy



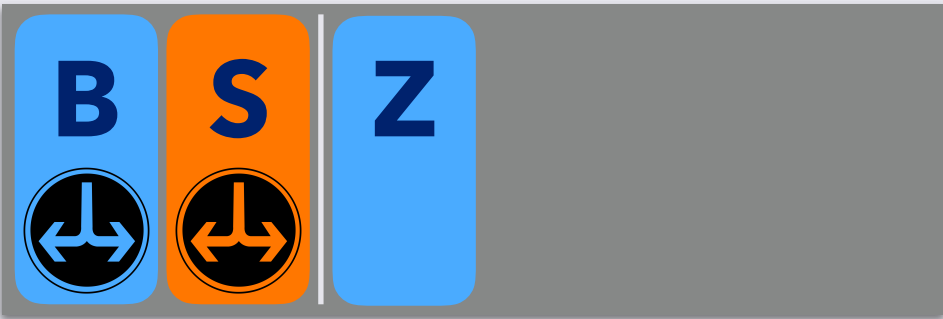
2. Translate
/green/ → /red/

- 3. Delete unreachable data
- 4. Move pivot translation

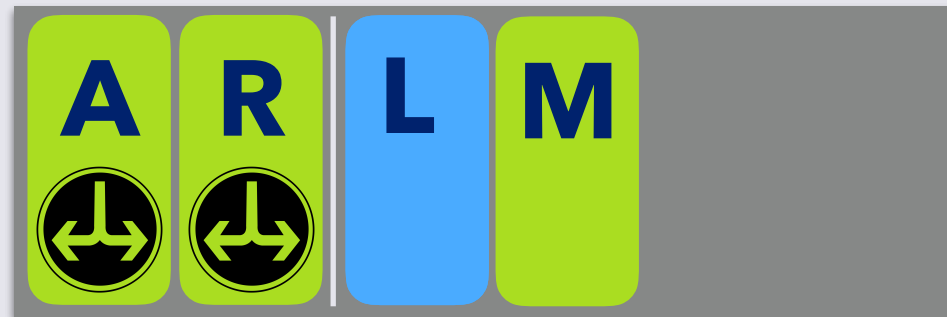
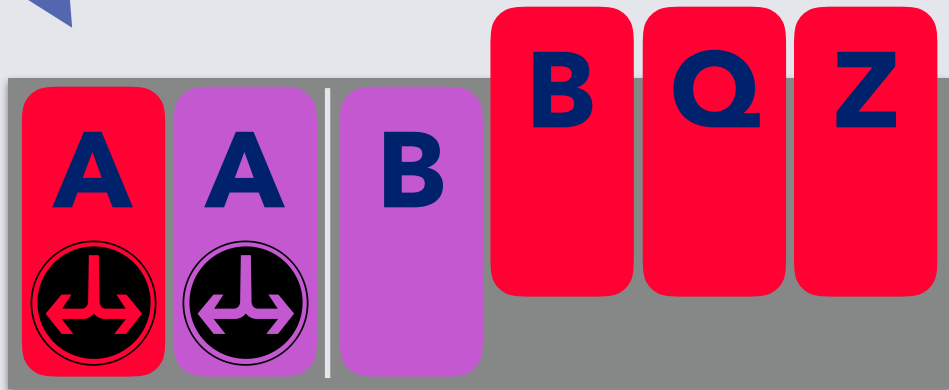


B ϵ -DAGs and Small Writes

Make some small
writes to /red/



Now flush with
copy-on-write



1. copy



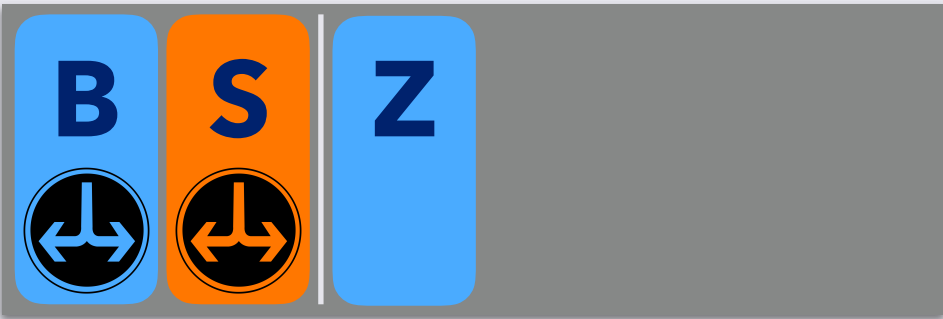
2. Translate
/green/ → /red/

- 3. Delete unreachable data
- 4. Move pivot translation

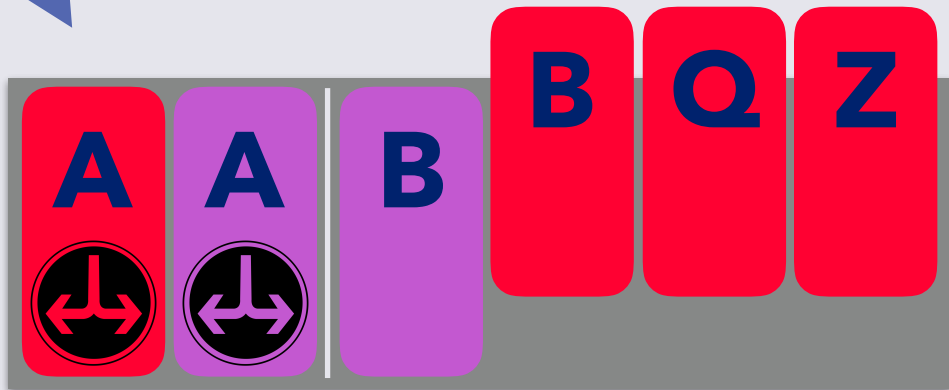


B ϵ -DAGs and Small Writes

Make some small
writes to /red/



Now flush with
copy-on-write



1. copy

5. Flush



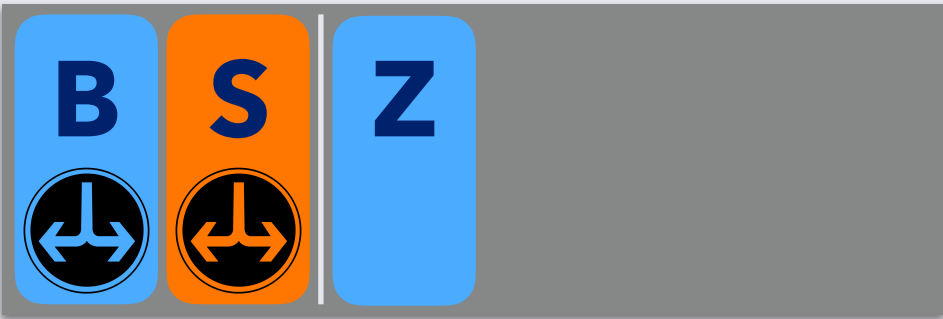
2. Translate
/green/ → /red/

- 3. Delete unreachable data
- 4. Move pivot translation

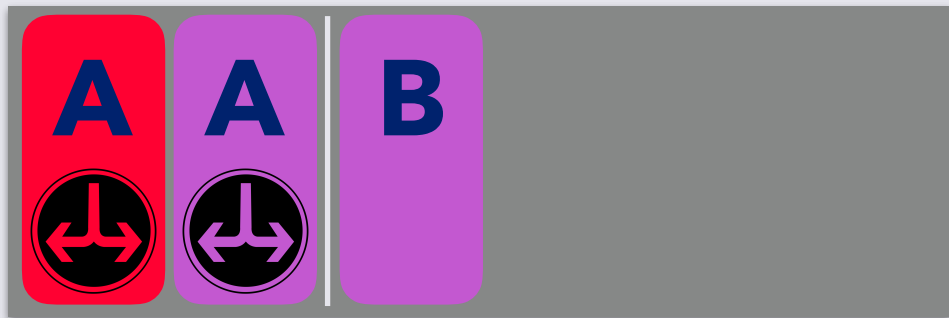


B ϵ -DAGs and Small Writes

Make some small
writes to /red/

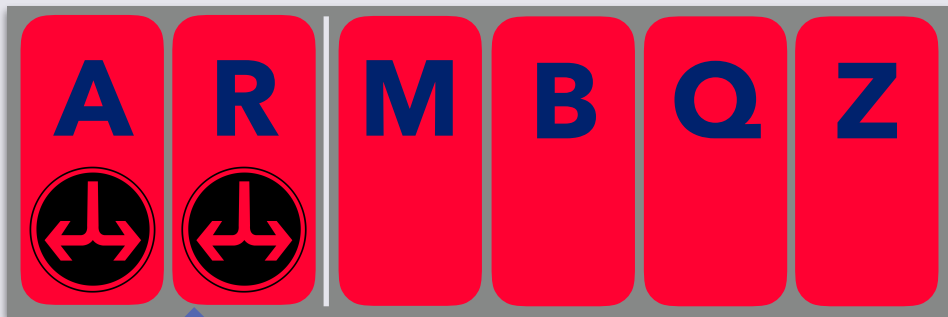


Now flush with
copy-on-write



1. copy

5. Flush



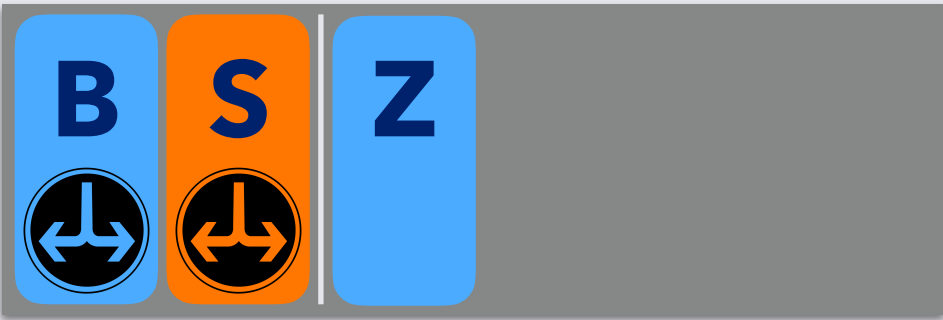
2. Translate
/green/ → /red/

3. Delete
unreachable data
4. Move pivot
translation



B ϵ -DAGs and Small Writes

Make some small
writes to /red/



Now flush with
copy-on-write

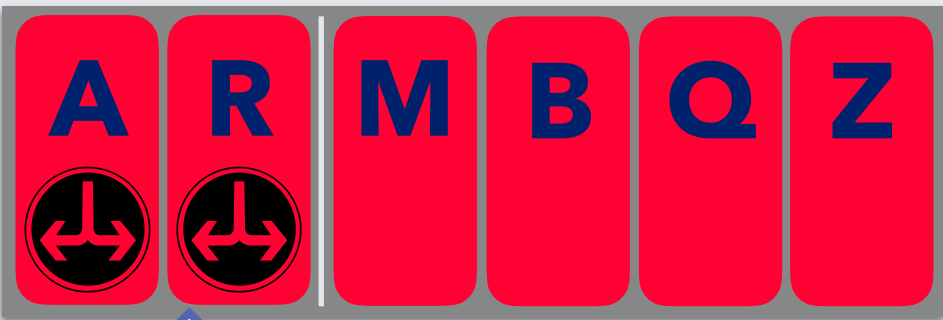
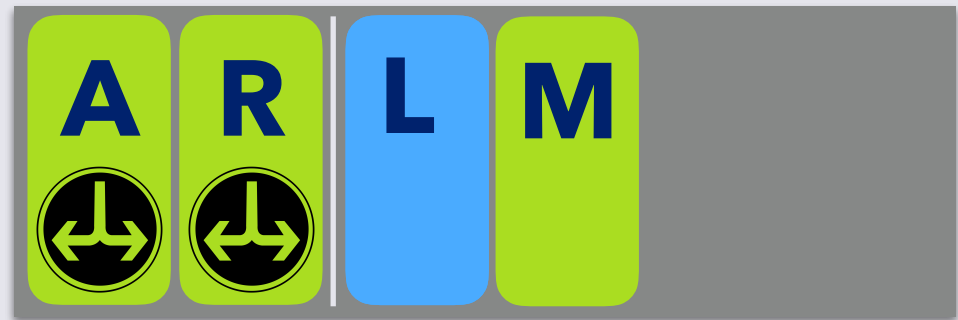


Applied multiple
small changes

5. Flush

2. Translate
/green/ → /red/

1. copy

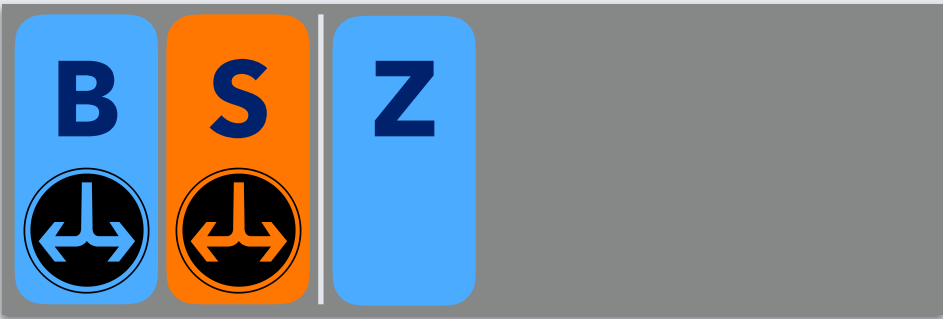


3. Delete
unreachable data
4. Move pivot
translation



B ϵ -DAGs and Small Writes

Make some small
writes to /red/



Now flush with
copy-on-write



Applied multiple
small changes

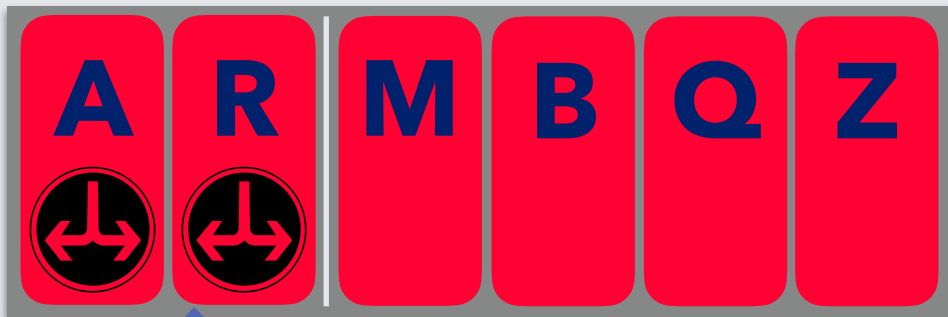
5. Flush

2. Translate
/green/ → /red/

Broke sharing of
one node



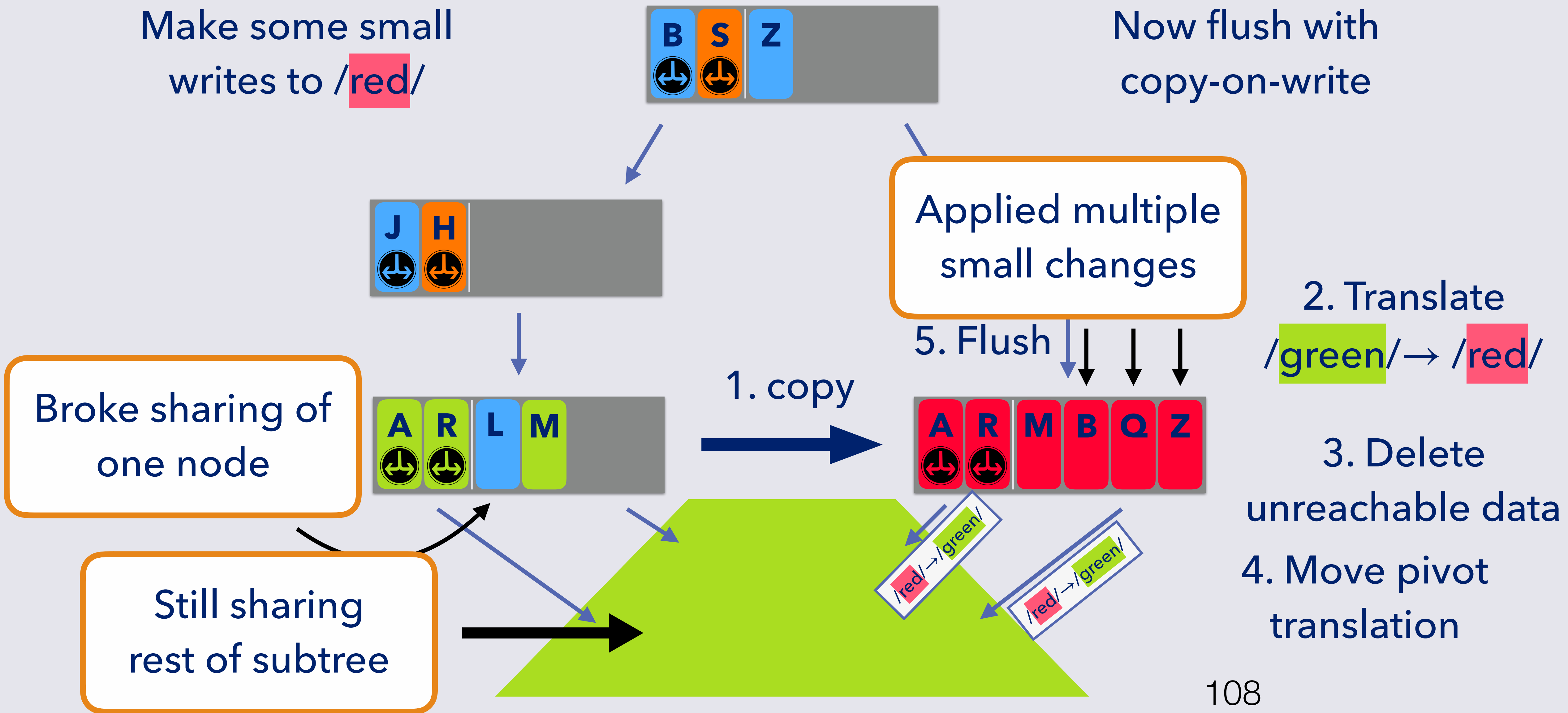
1. copy



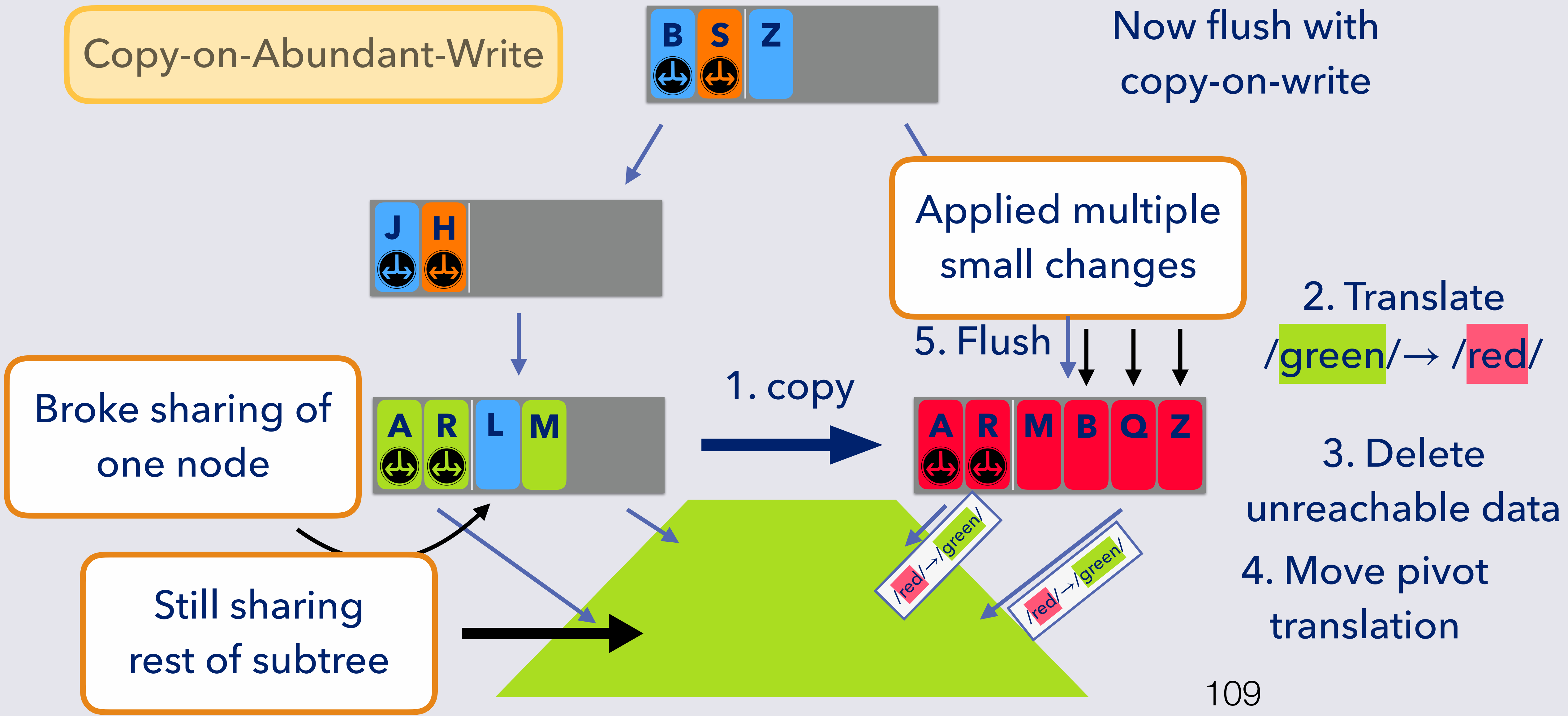
3. Delete
unreachable data
4. Move pivot
translation



B ϵ -DAGs and Small Writes



B ϵ -DAGs and Small Writes



Logical Copy with B^ϵ -DAGs

Performance Goals

Space efficient

?

Fast writes

?

Low latency

?

Fast reads

?

Copy-specific

General file system

Logical Copy with B^ϵ -DAGs

Performance Goals

Space efficient

?

Low latency

?

Copy-specific

B^ϵ -trees have large nodes

Fast writes

✓

⇒ Locality*

Fast reads

✓

General file system

**File Systems Fated for Senescence? Nonsense, Says Science!,
Conway et al, FAST 2017*

Logical Copy with B^ϵ -DAGs

Performance Goals

Copy-on-Abundant-Write

Space efficient



Low latency



Copy-specific

B^ϵ -trees have large nodes

Fast writes



\Rightarrow Locality*

Fast reads



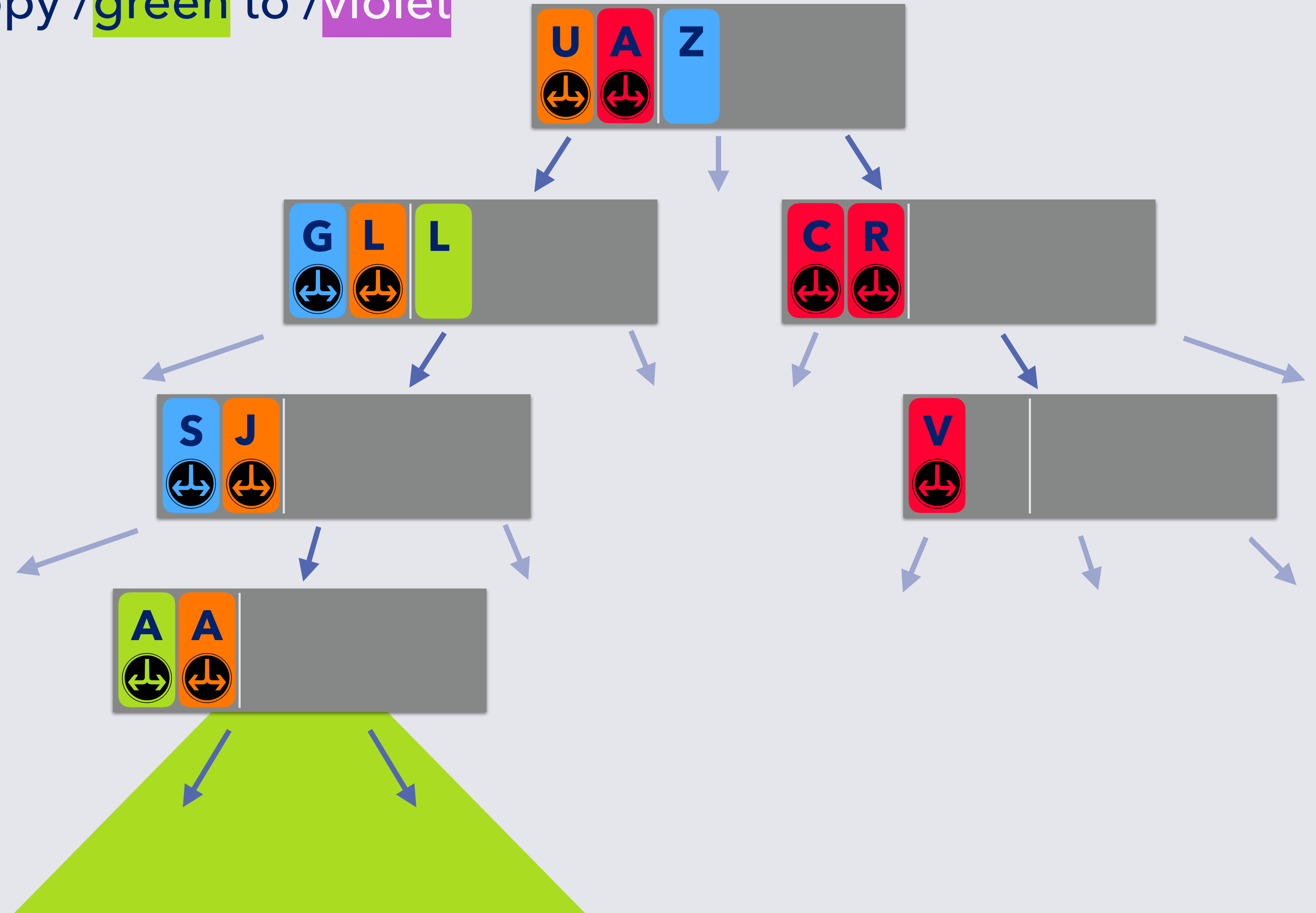
General file system

**File Systems Fated for Senescence? Nonsense, Says Science!,
Conway et al, FAST 2017*

Reducing Copy Latency in B^ϵ -DAGs

Reducing Copy Latency in B^ϵ -DAGs

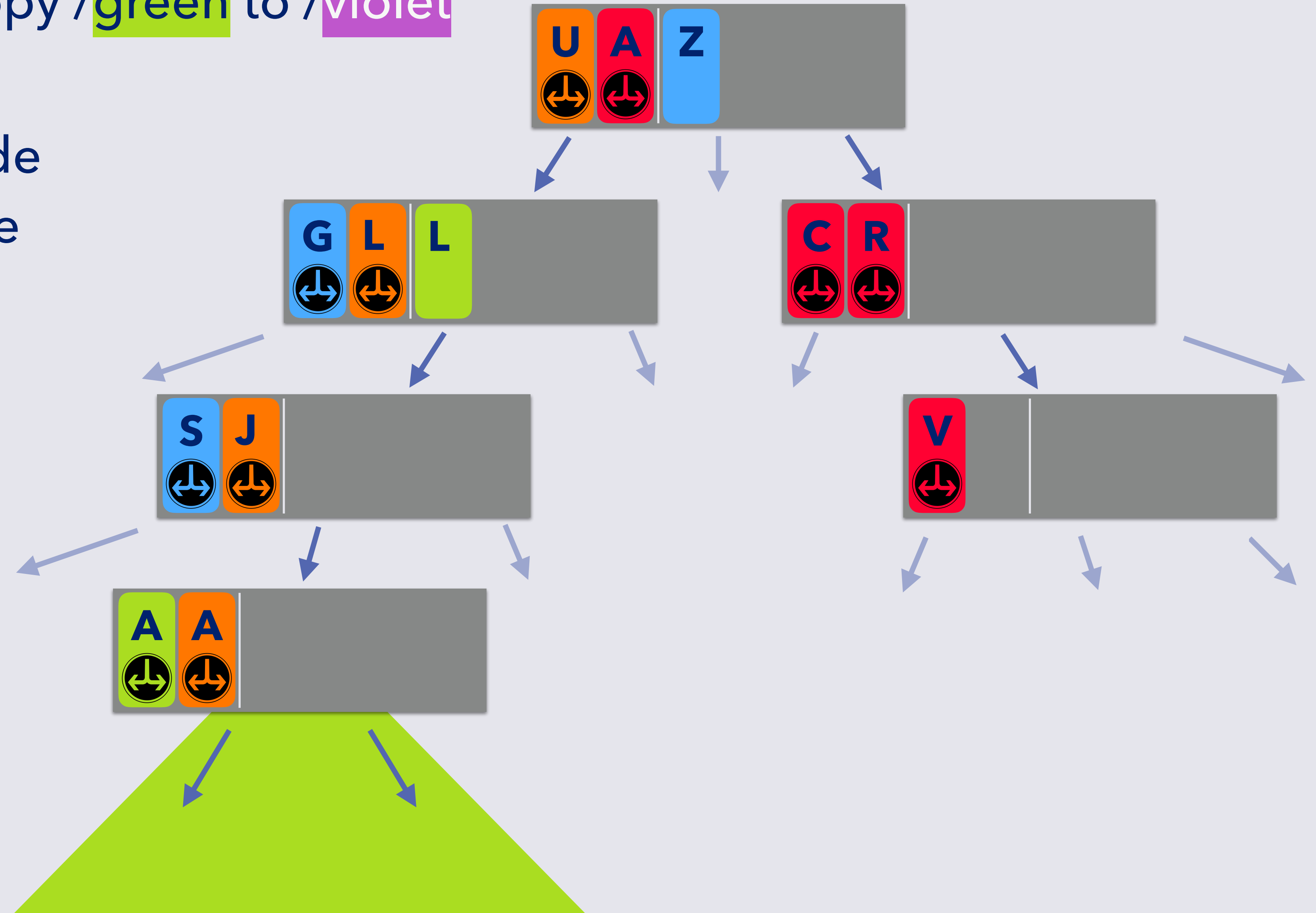
Copy /green to /violet



Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

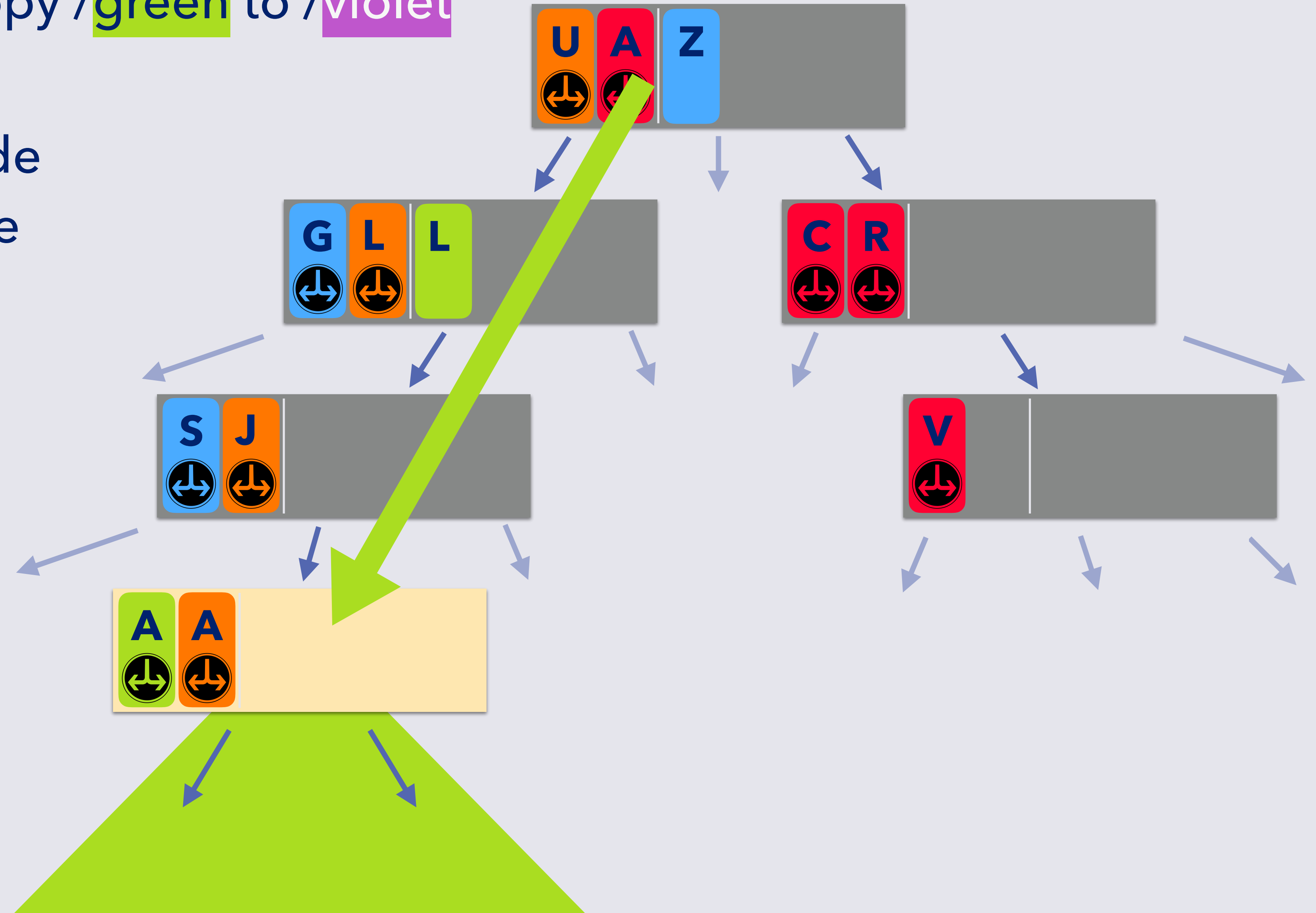
1. Flush messages to node covering /green subtree



Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

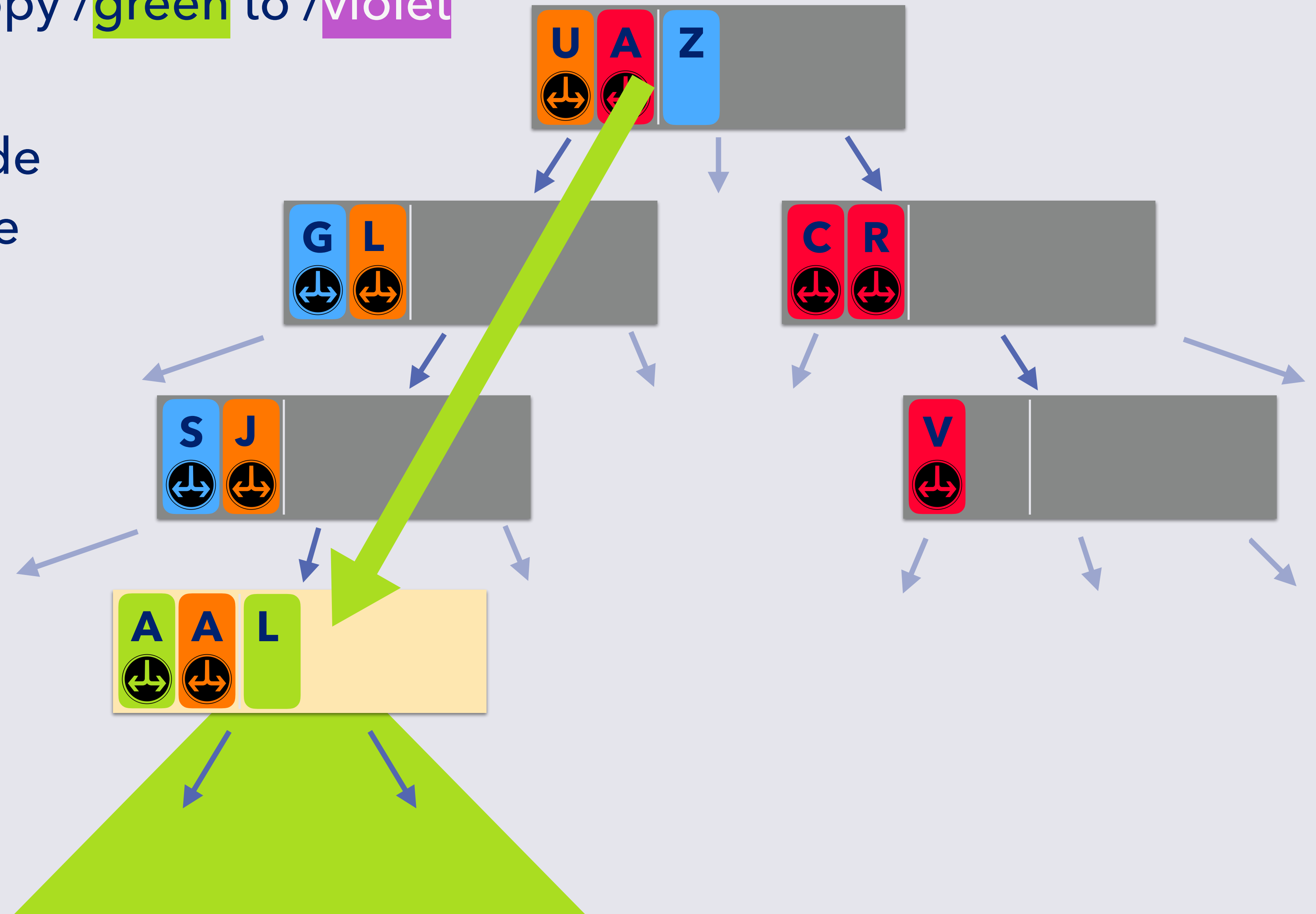
1. Flush messages to node covering /green subtree



Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

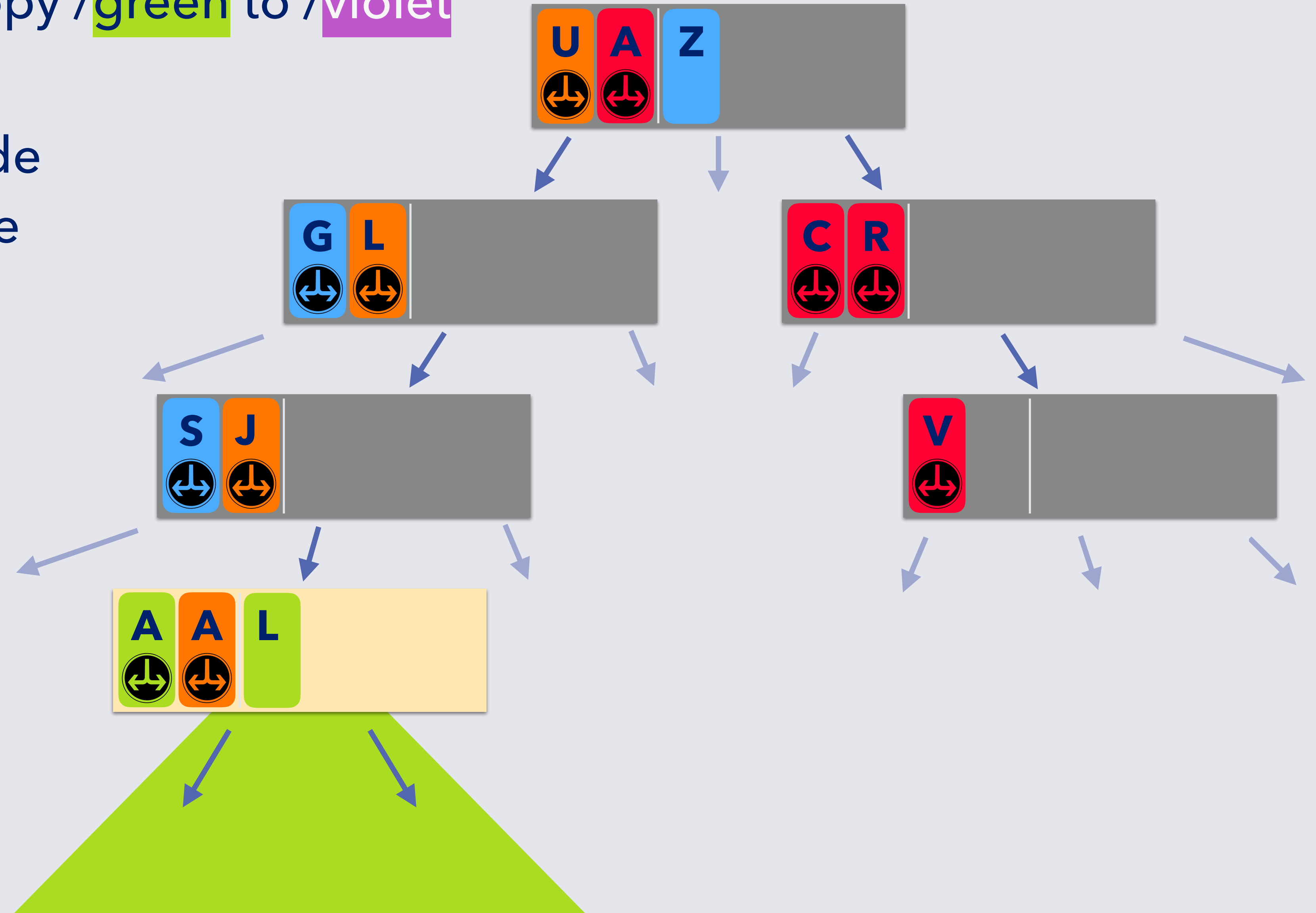
1. Flush messages to node covering /green subtree



Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

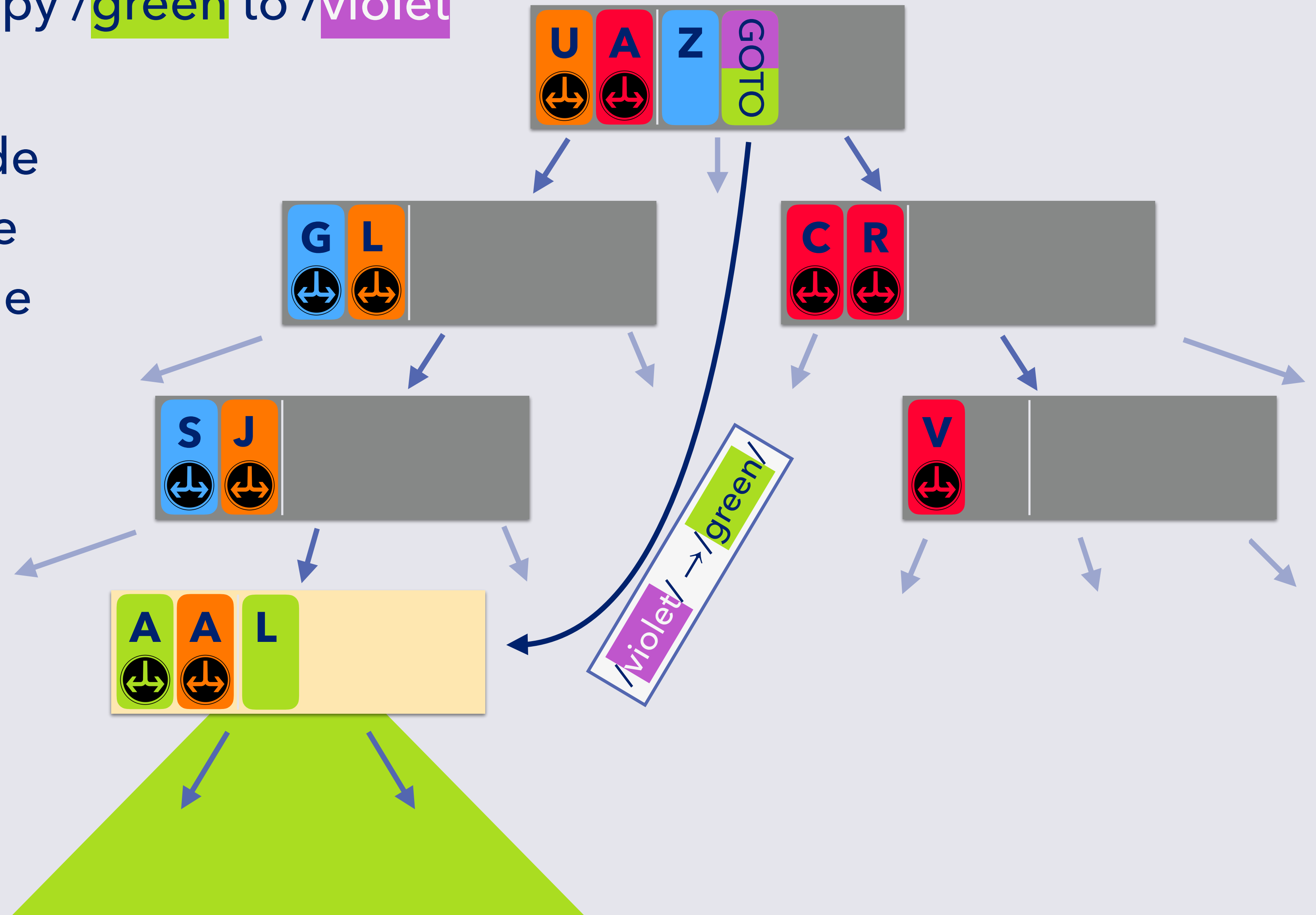
- # 1. Flush messages to node covering /green subtree



Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message



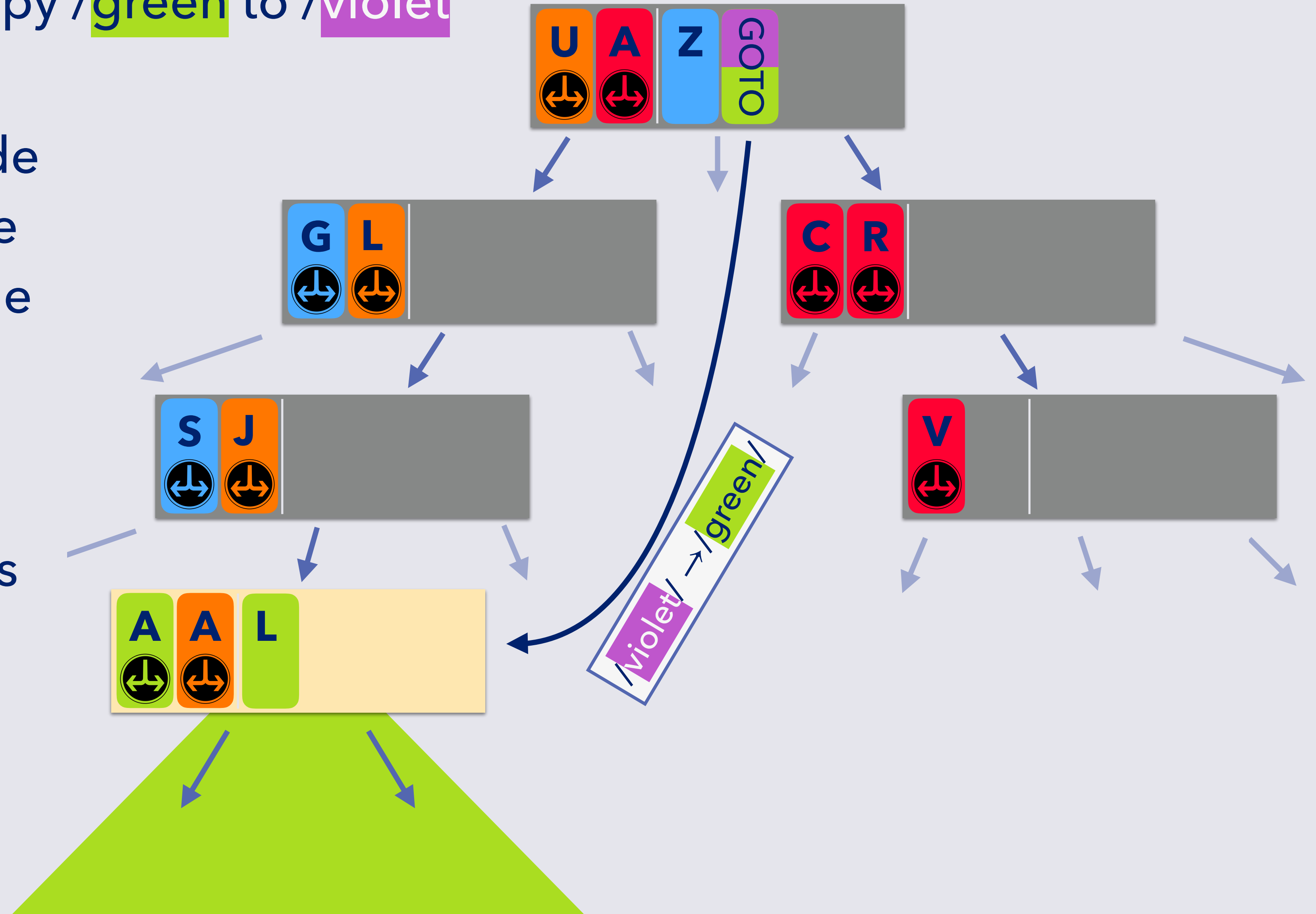
Reducing Copy Latency in B^ϵ -DAGs

Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message



A GOTO message changes the structure of the tree



Reducing Copy Latency in B^ϵ -DAGs

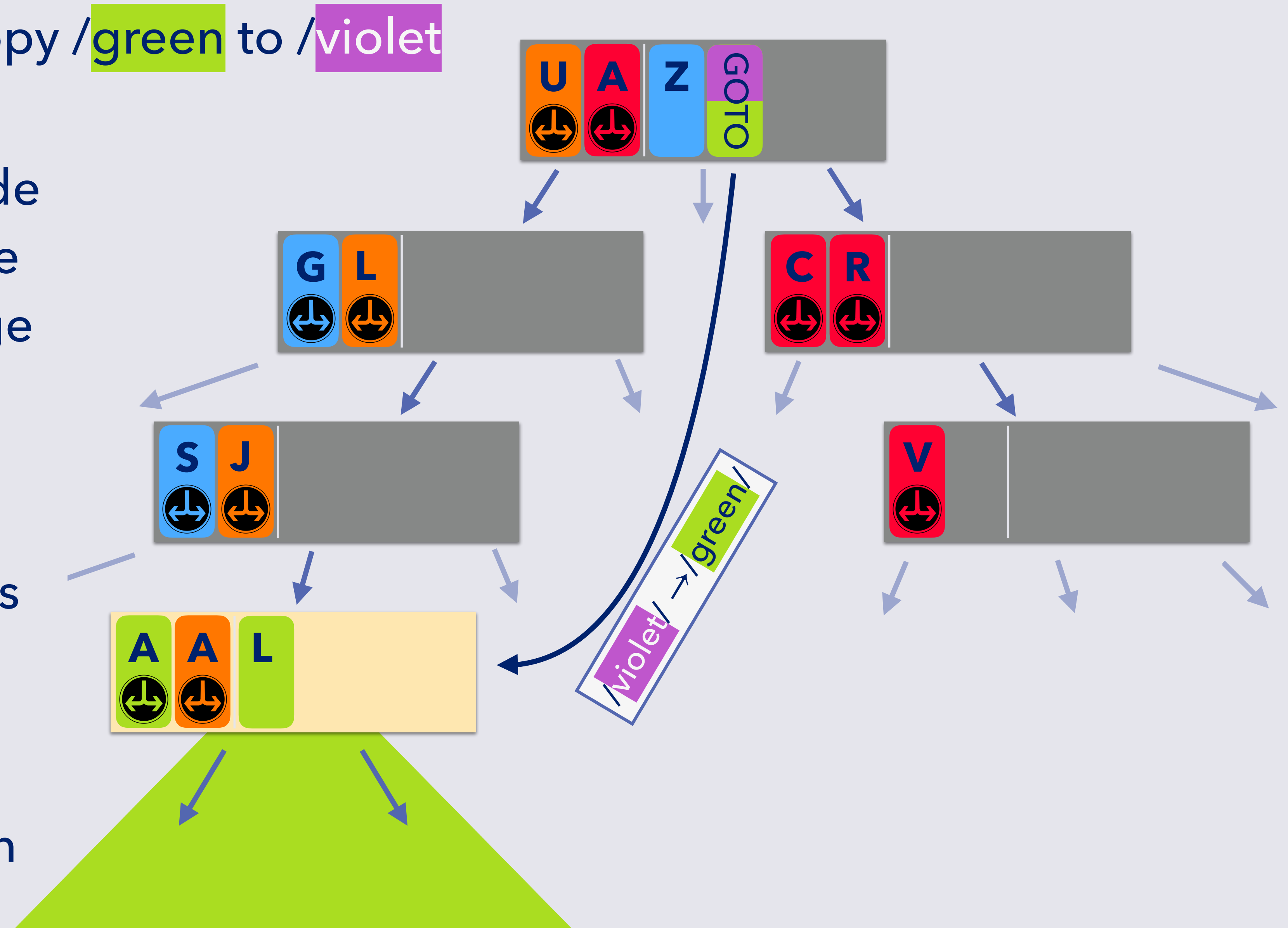
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B ϵ -DAGs

Copy /green to /violet

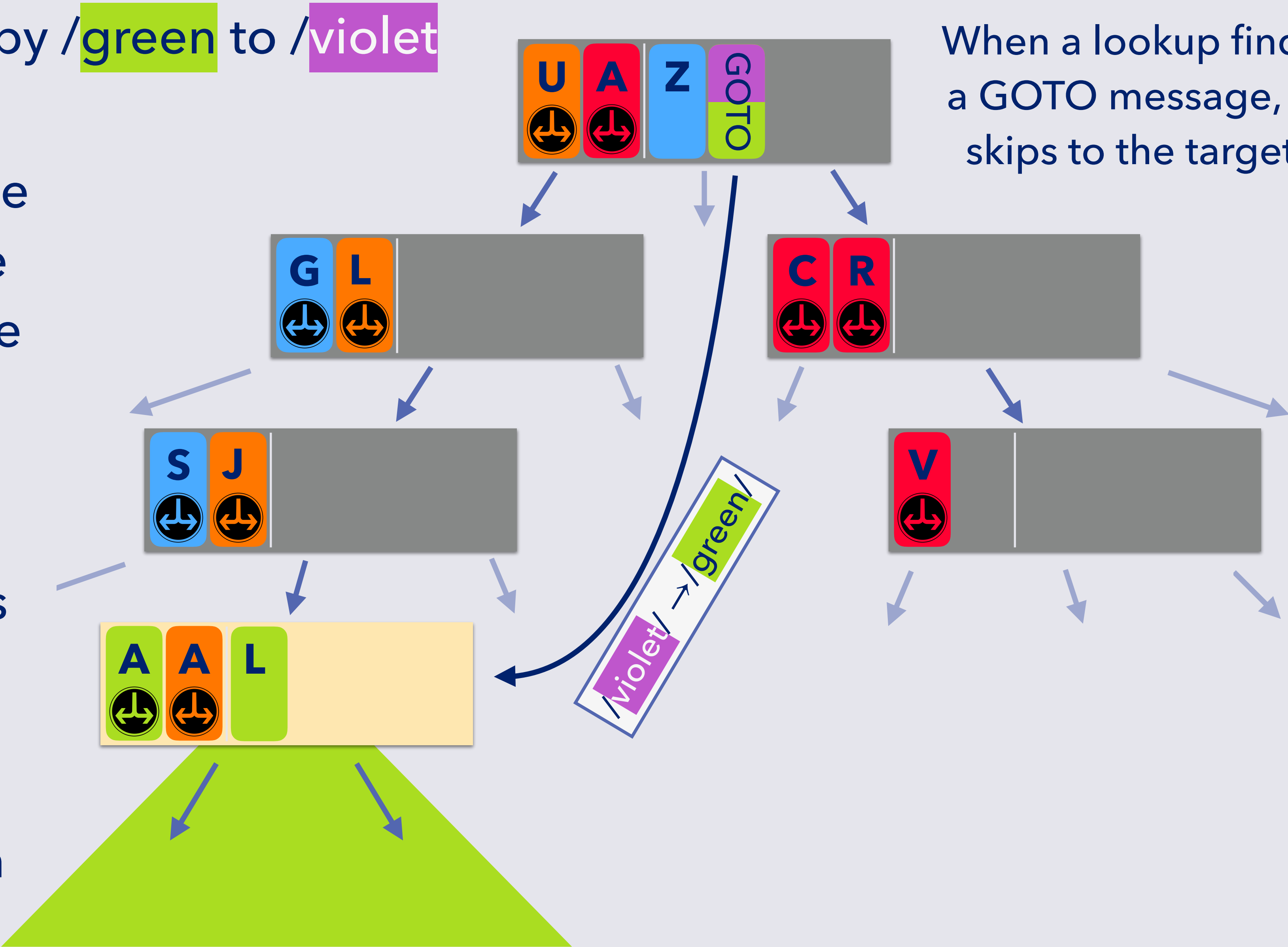
- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

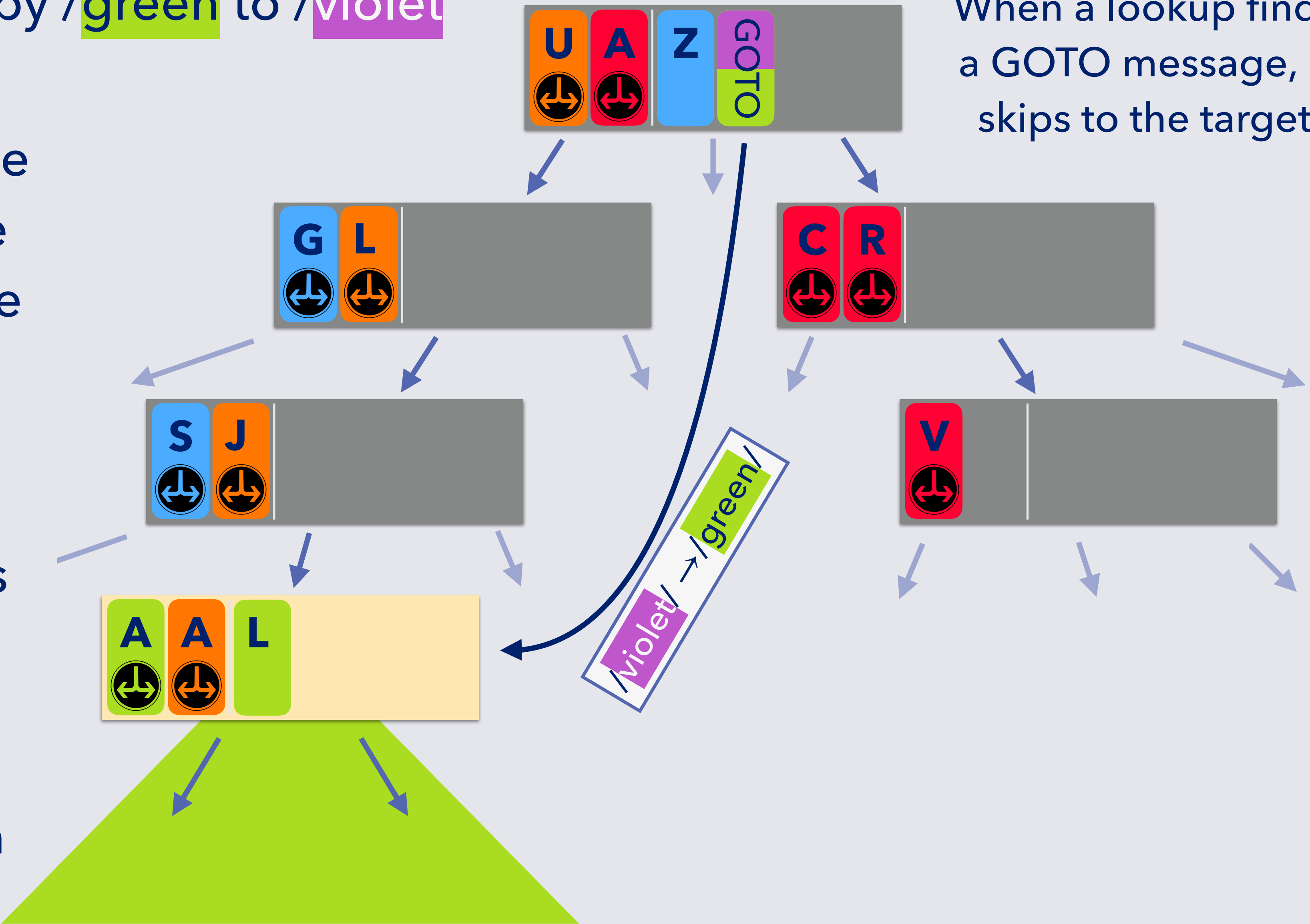
1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

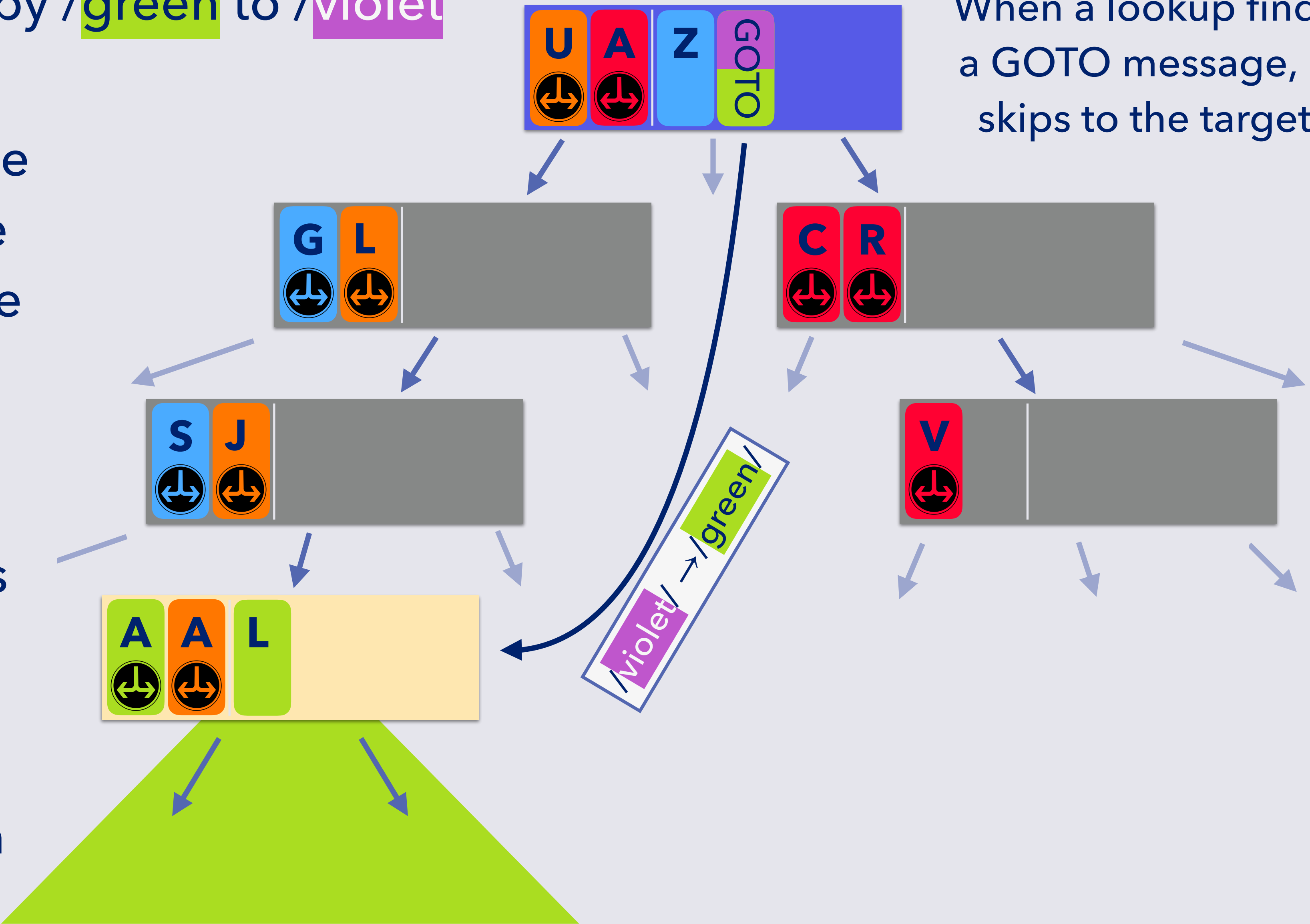
- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

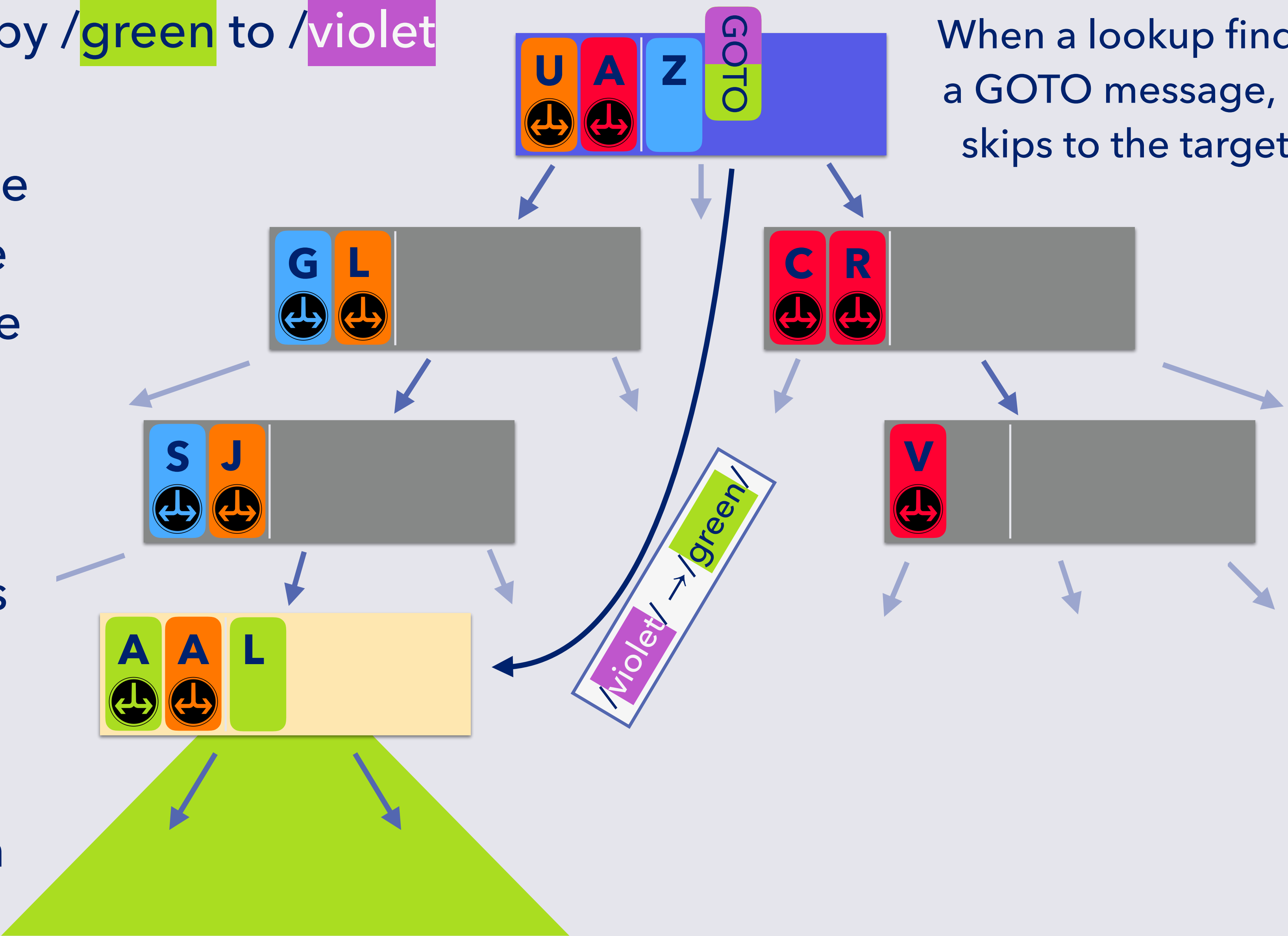
1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

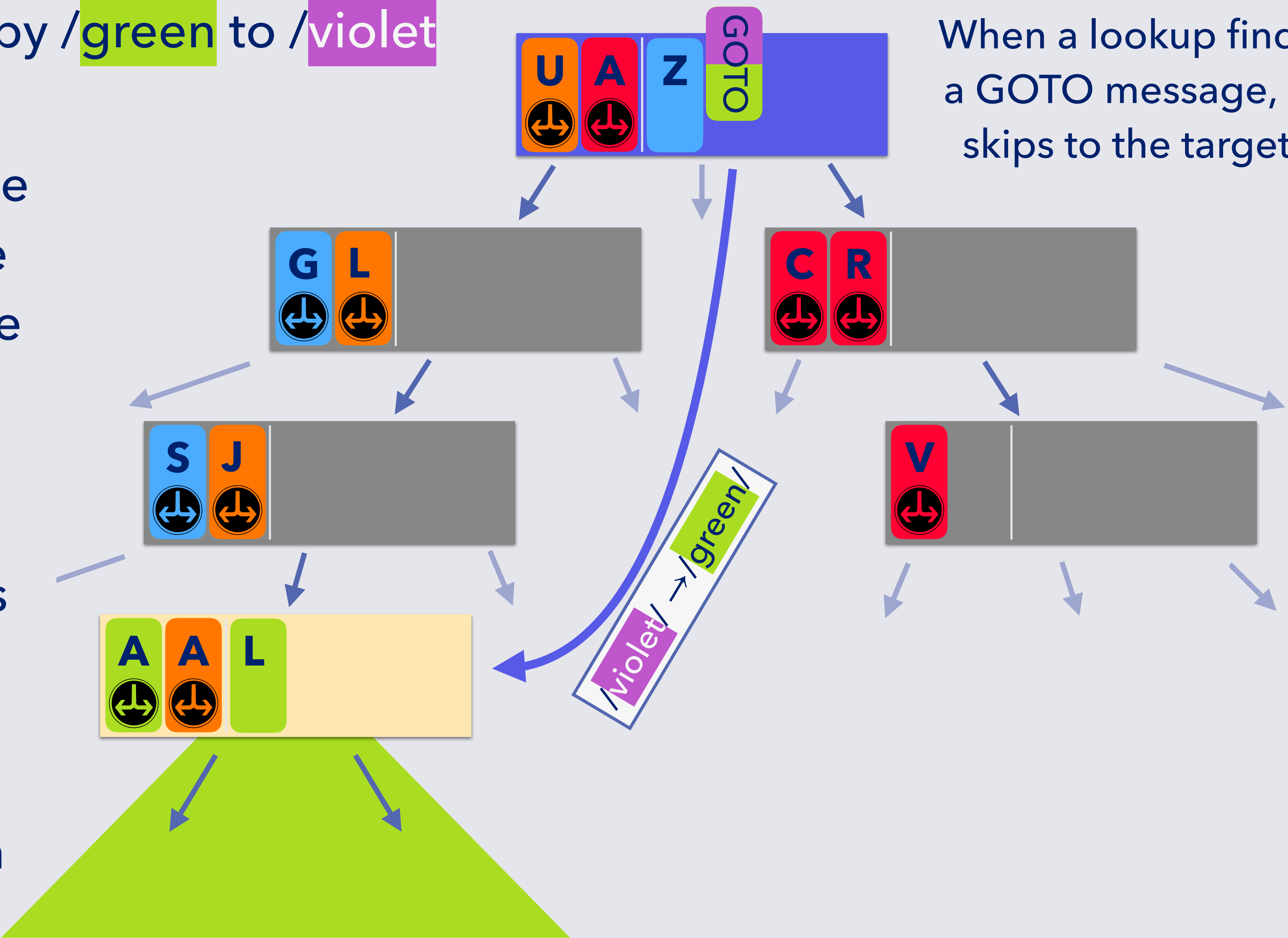
1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

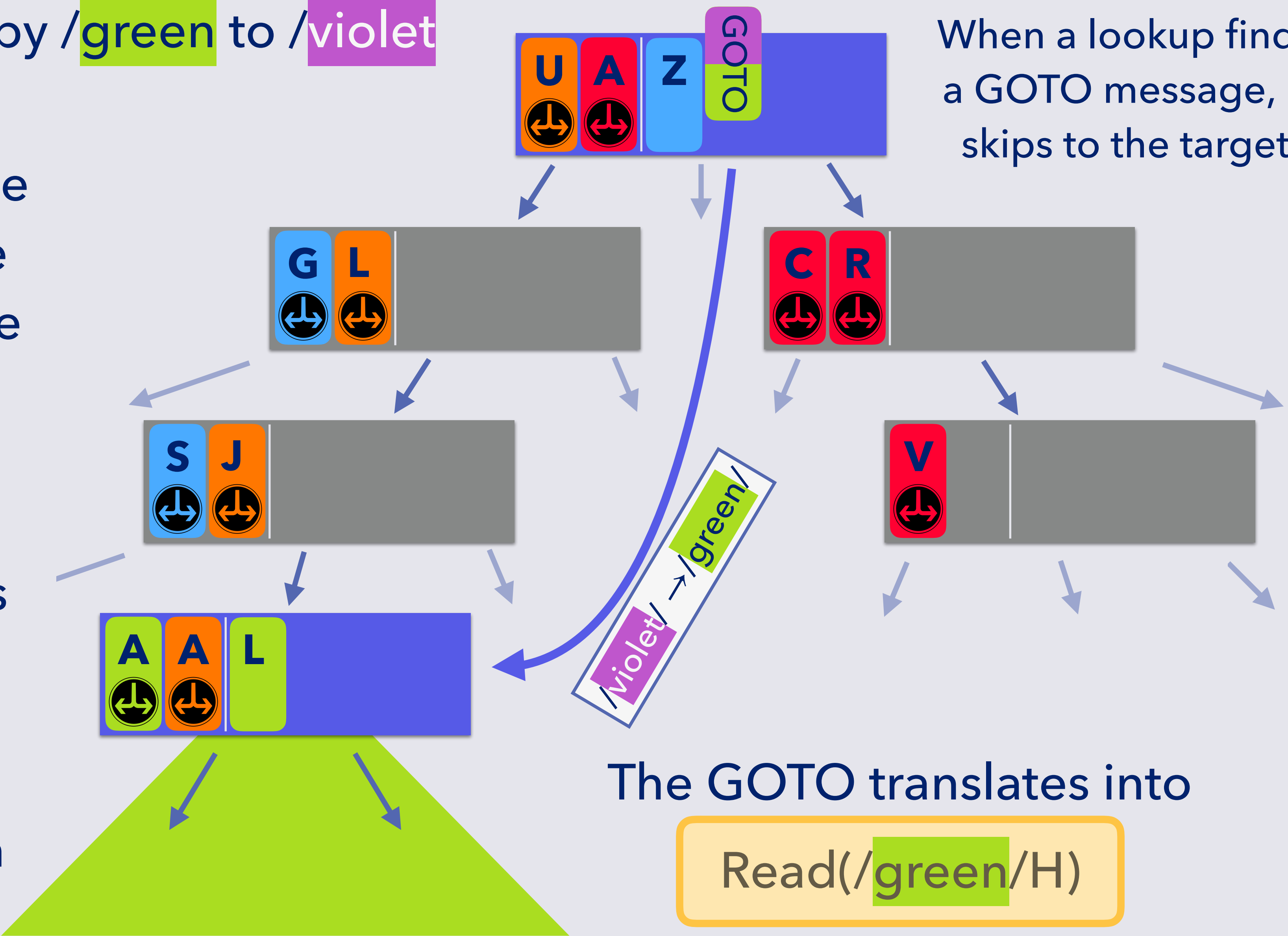
1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



The GOTO translates into

Read(/green/H)

Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

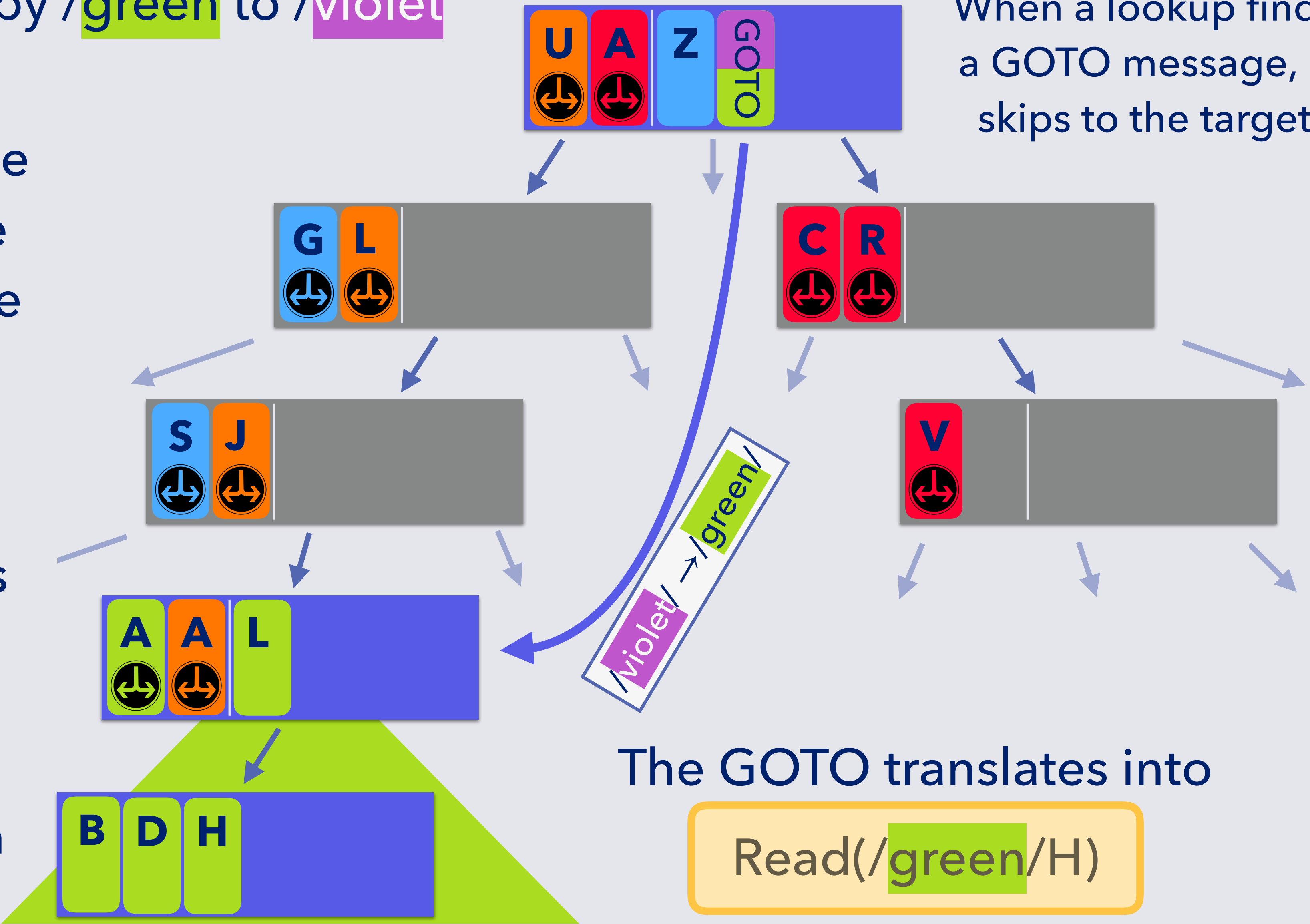
1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Read(/green/H)

Reducing Copy Latency in B ϵ -DA

Read(/violet/H)

Copy /green to /violet

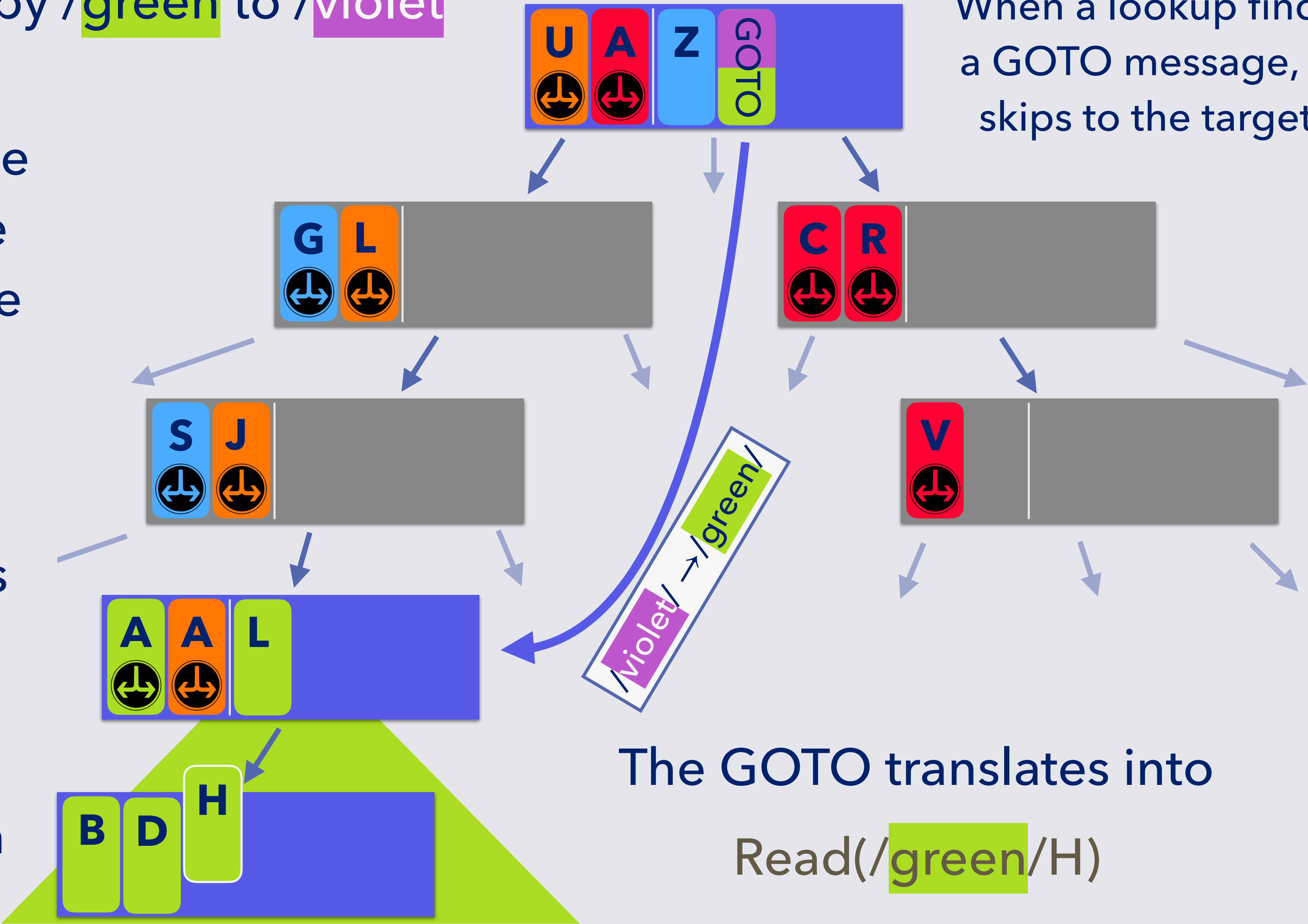
- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation

When a lookup finds a GOTO message, it skips to the target



Reducing Copy Latency in B^ϵ -DAGs

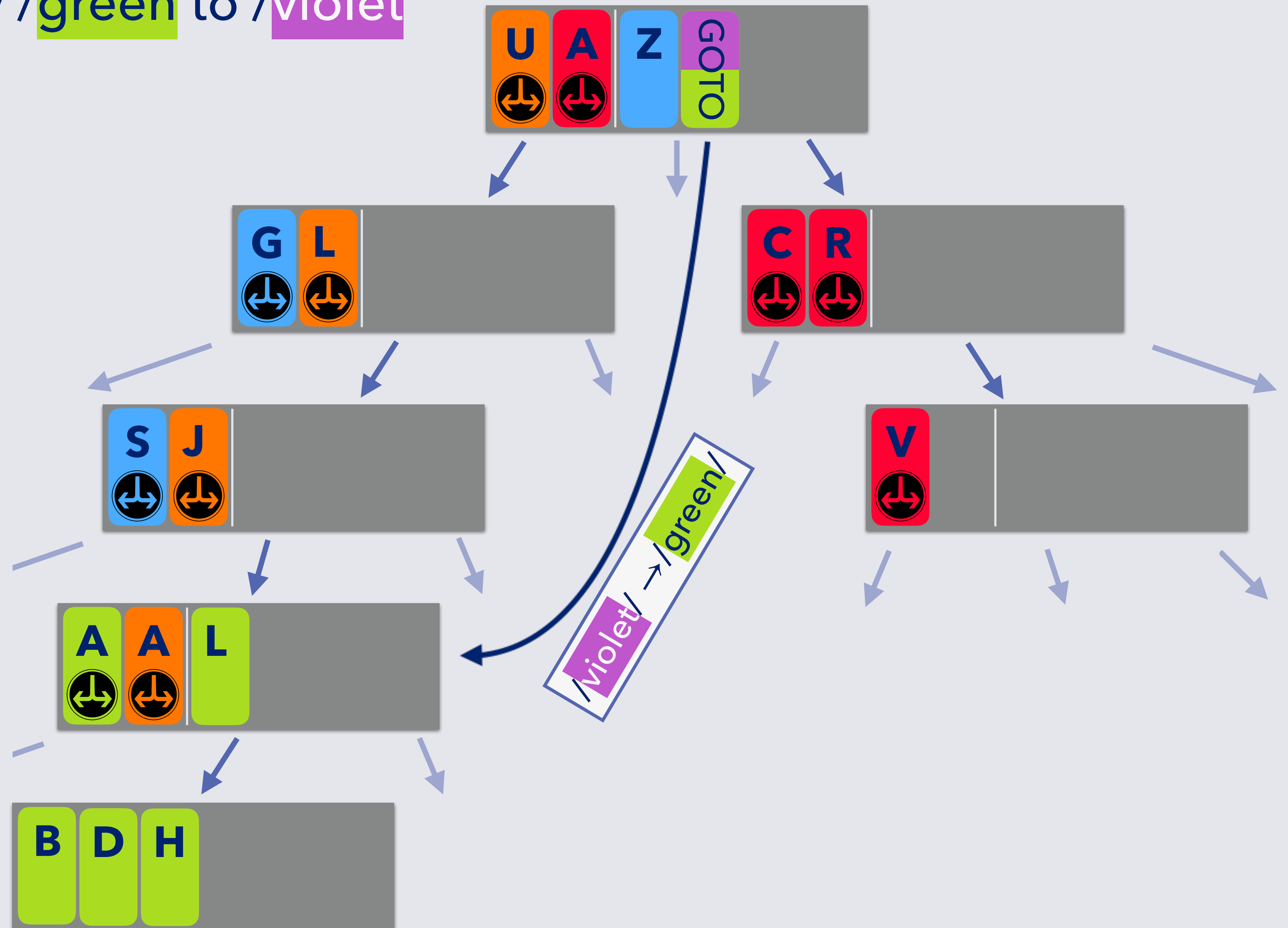
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

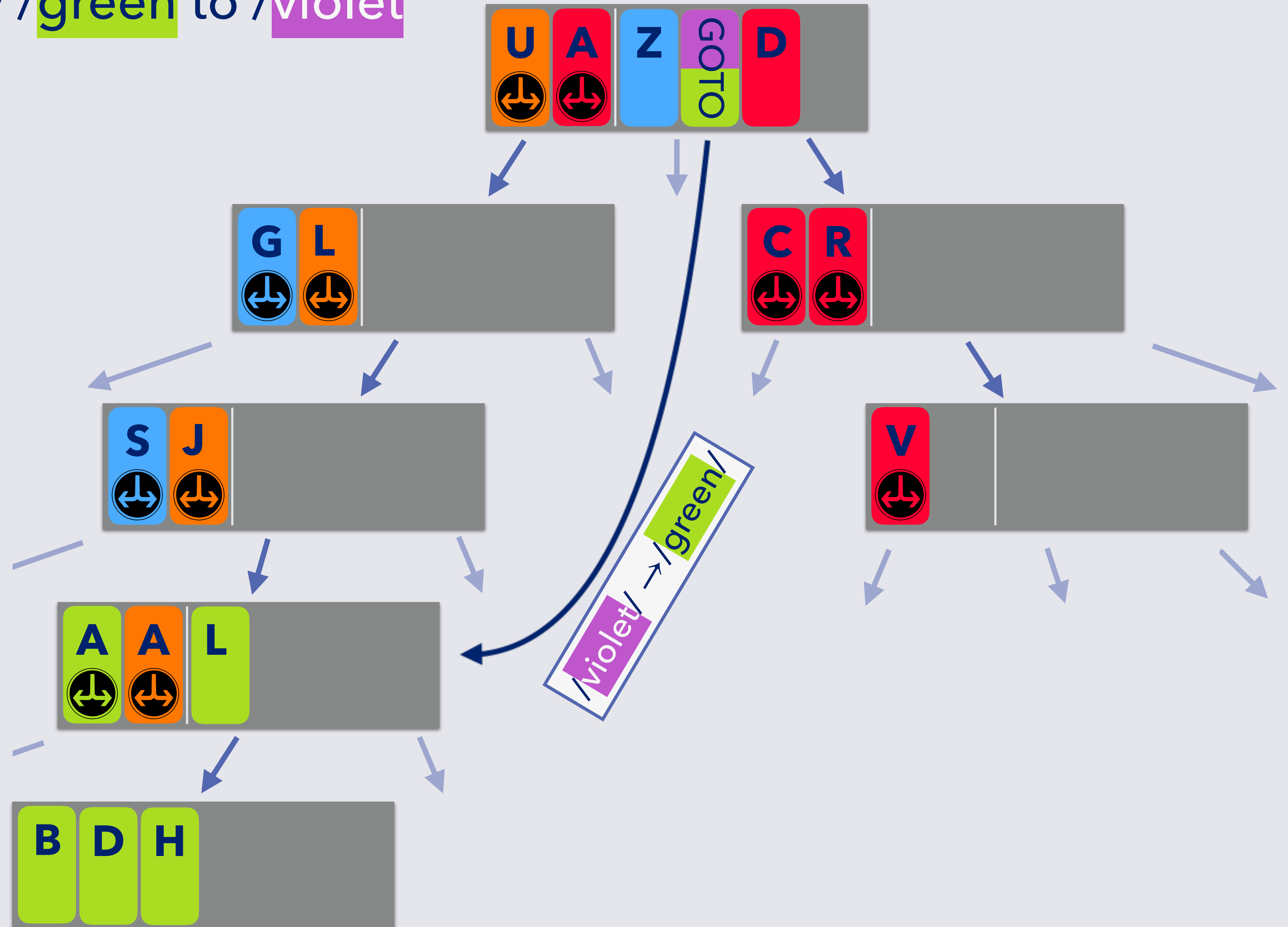
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

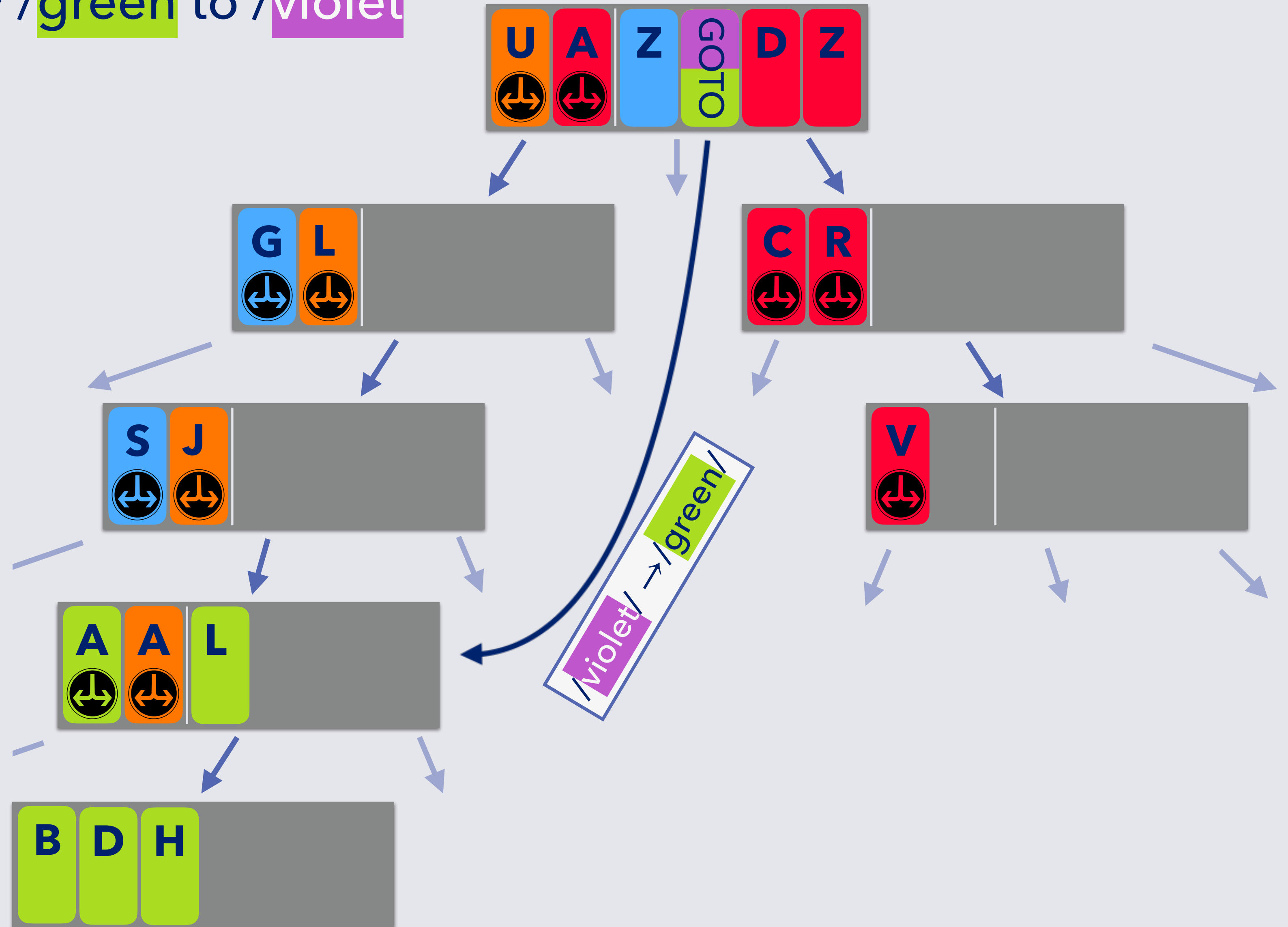
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

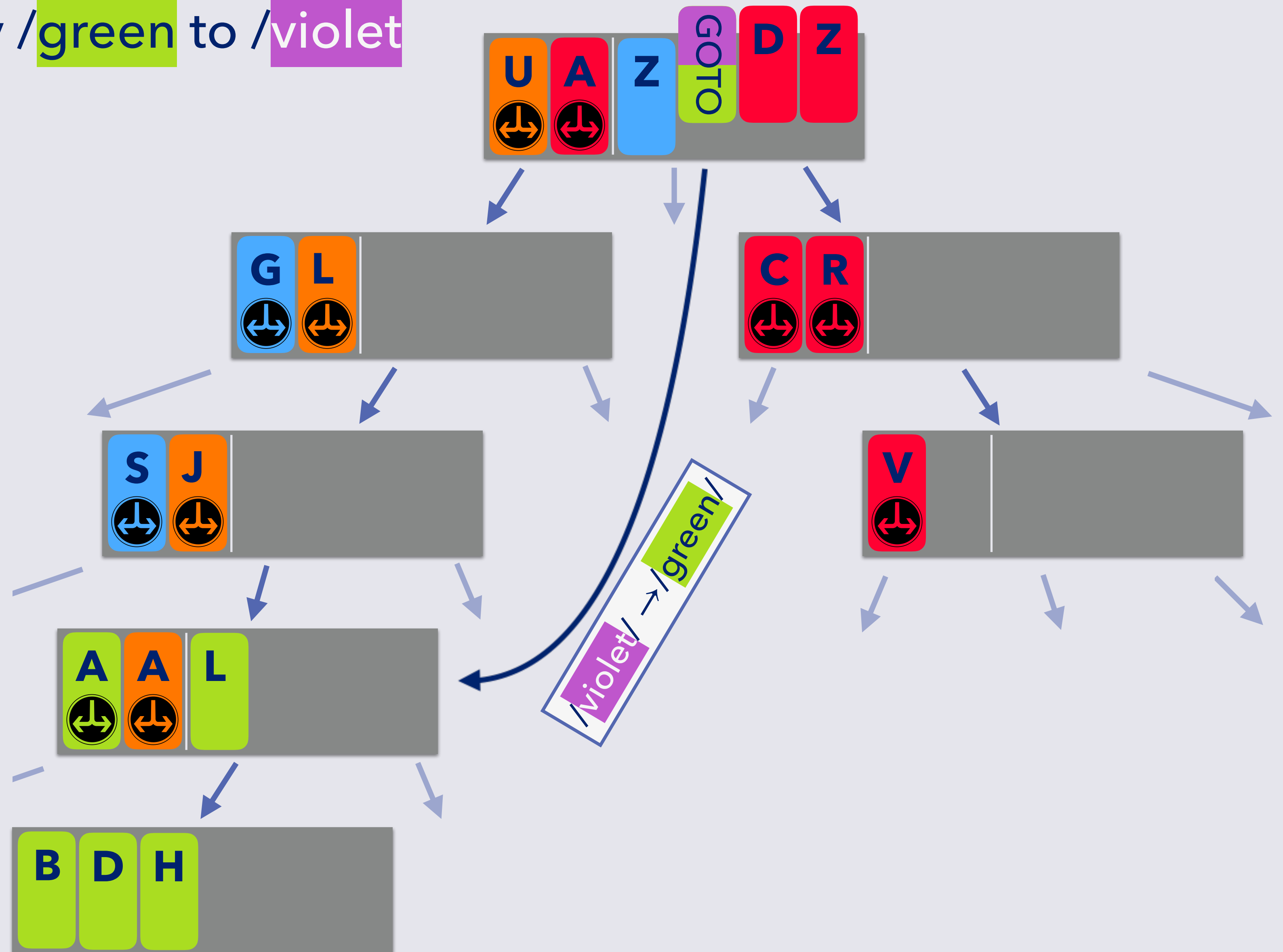
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

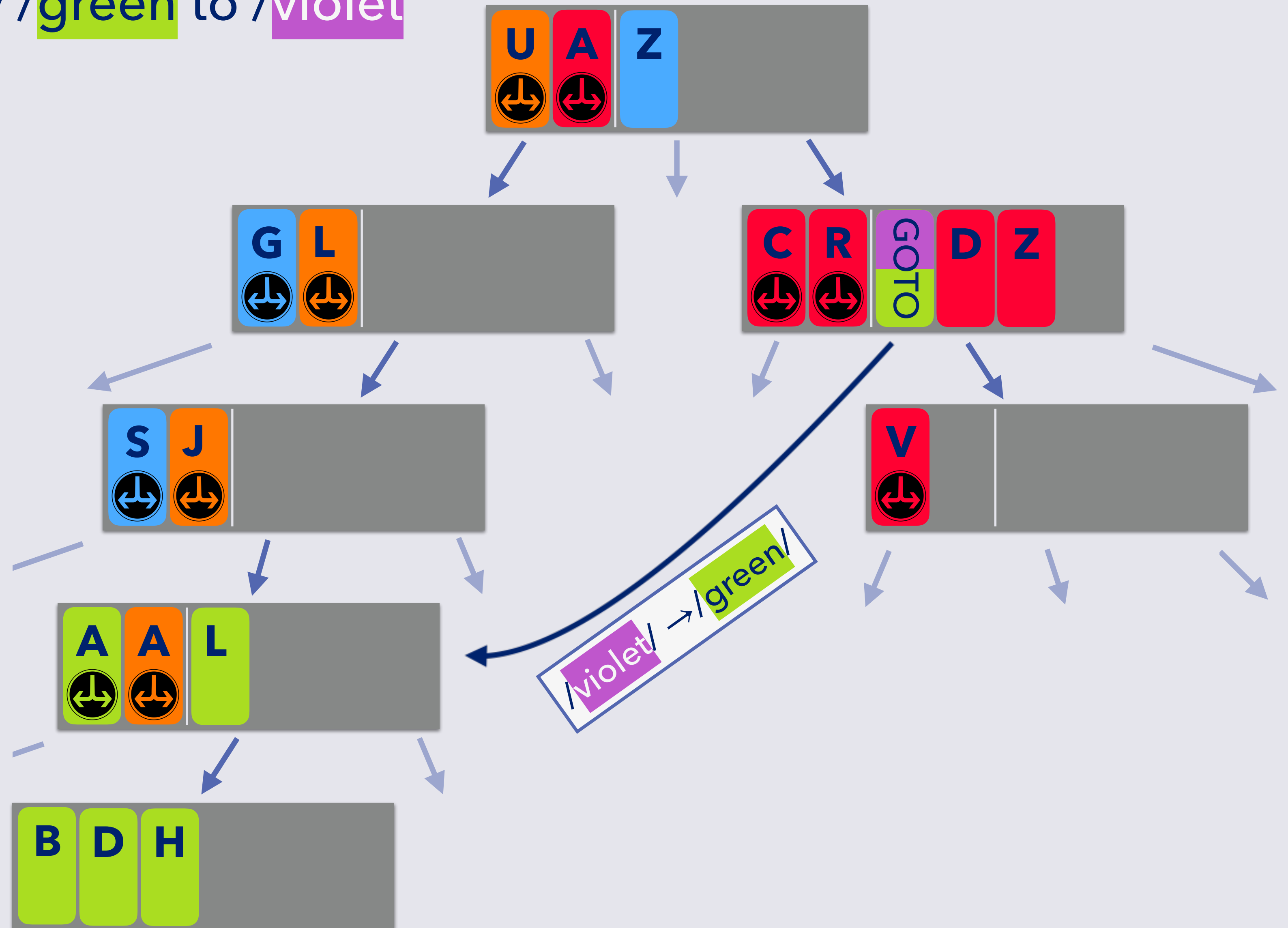
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B ϵ -DAGs

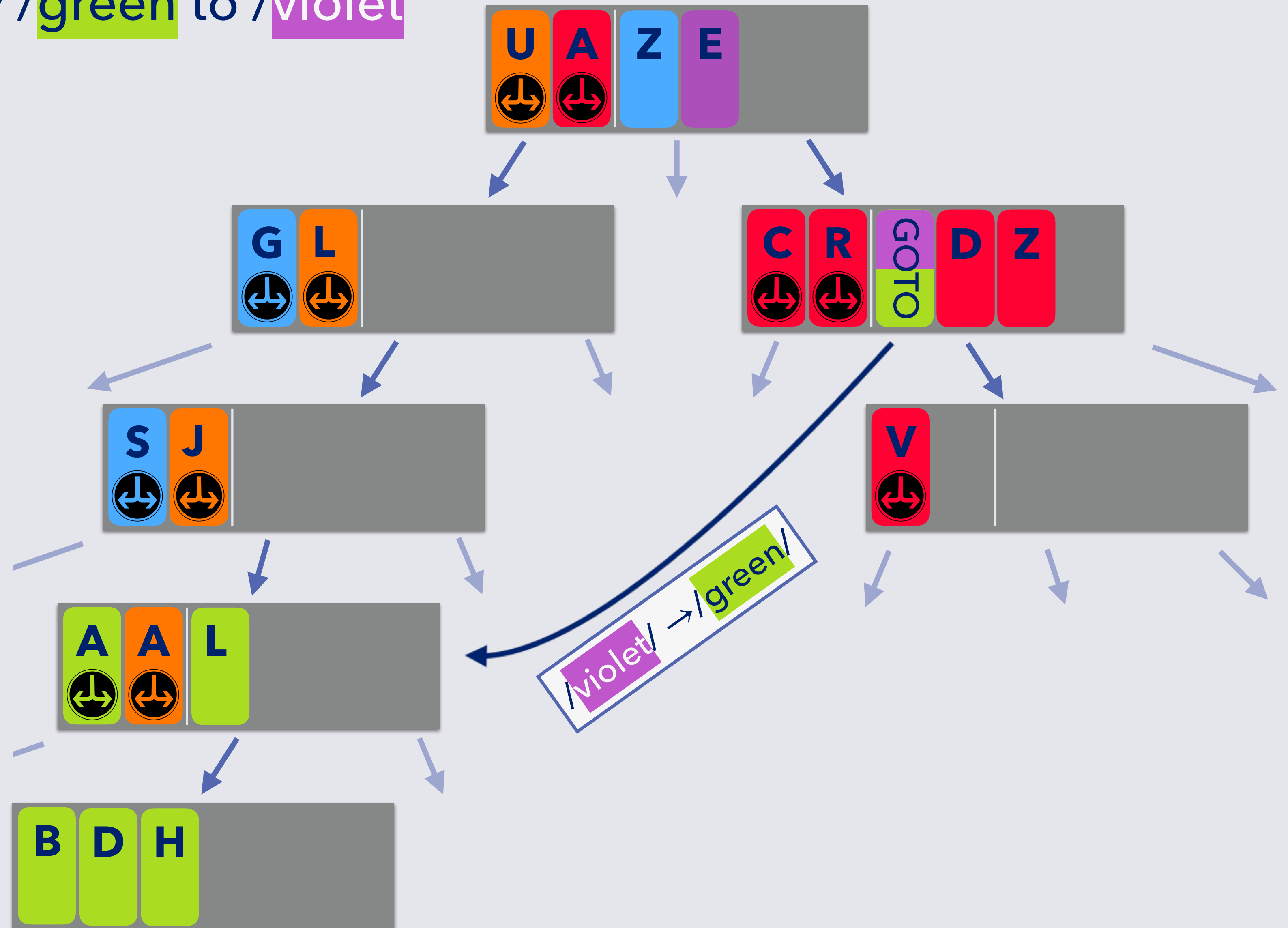
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

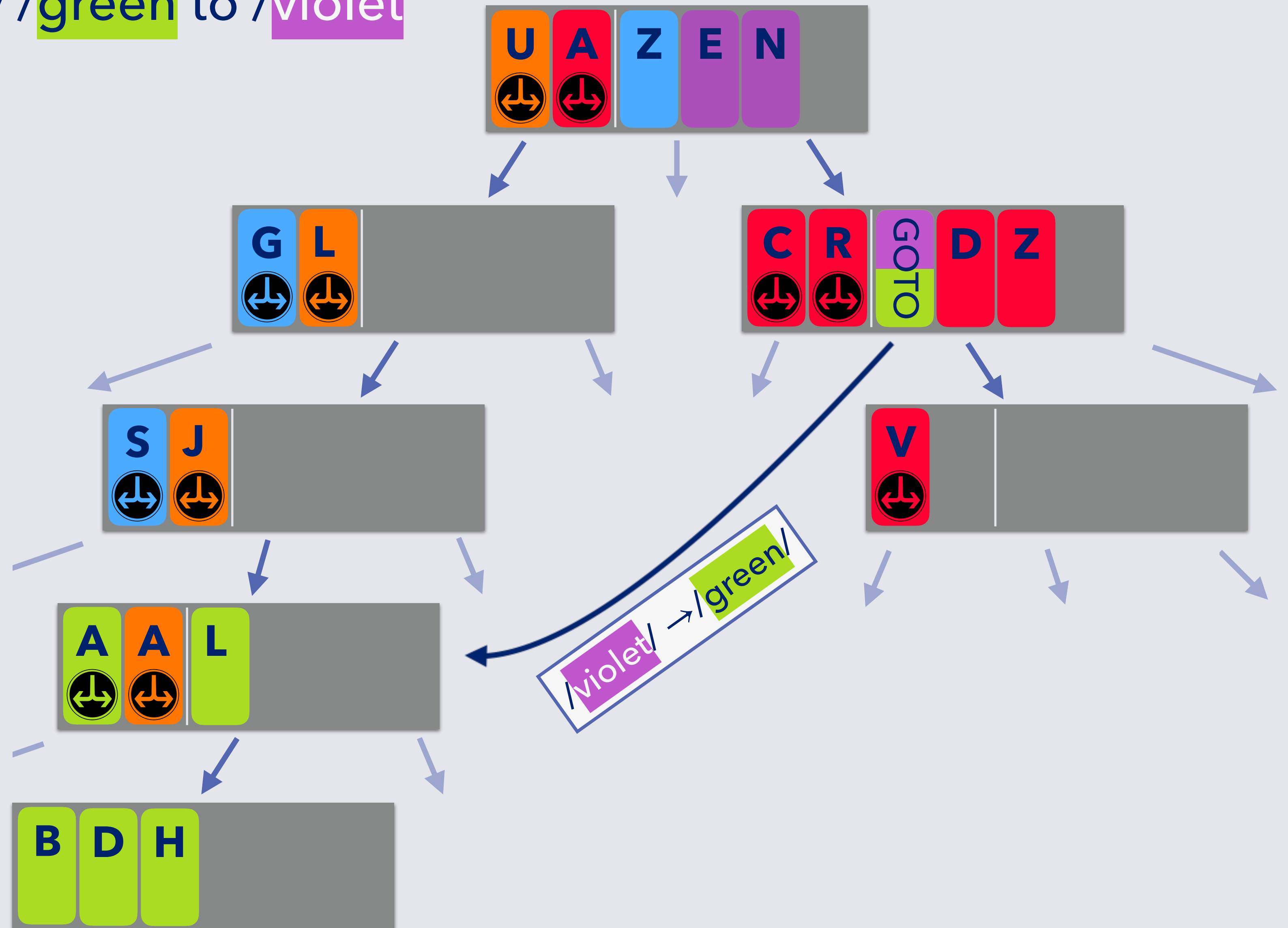
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

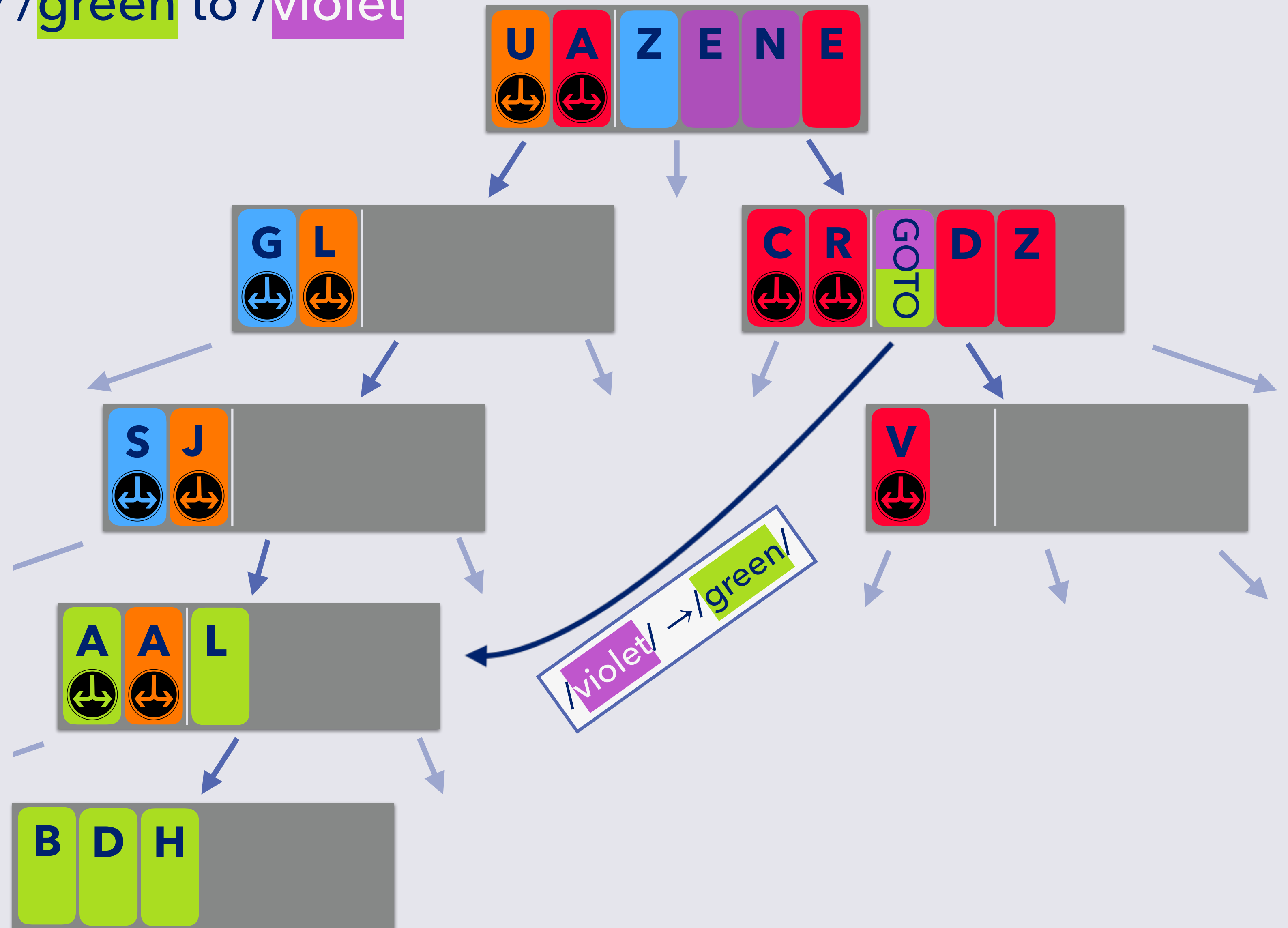
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

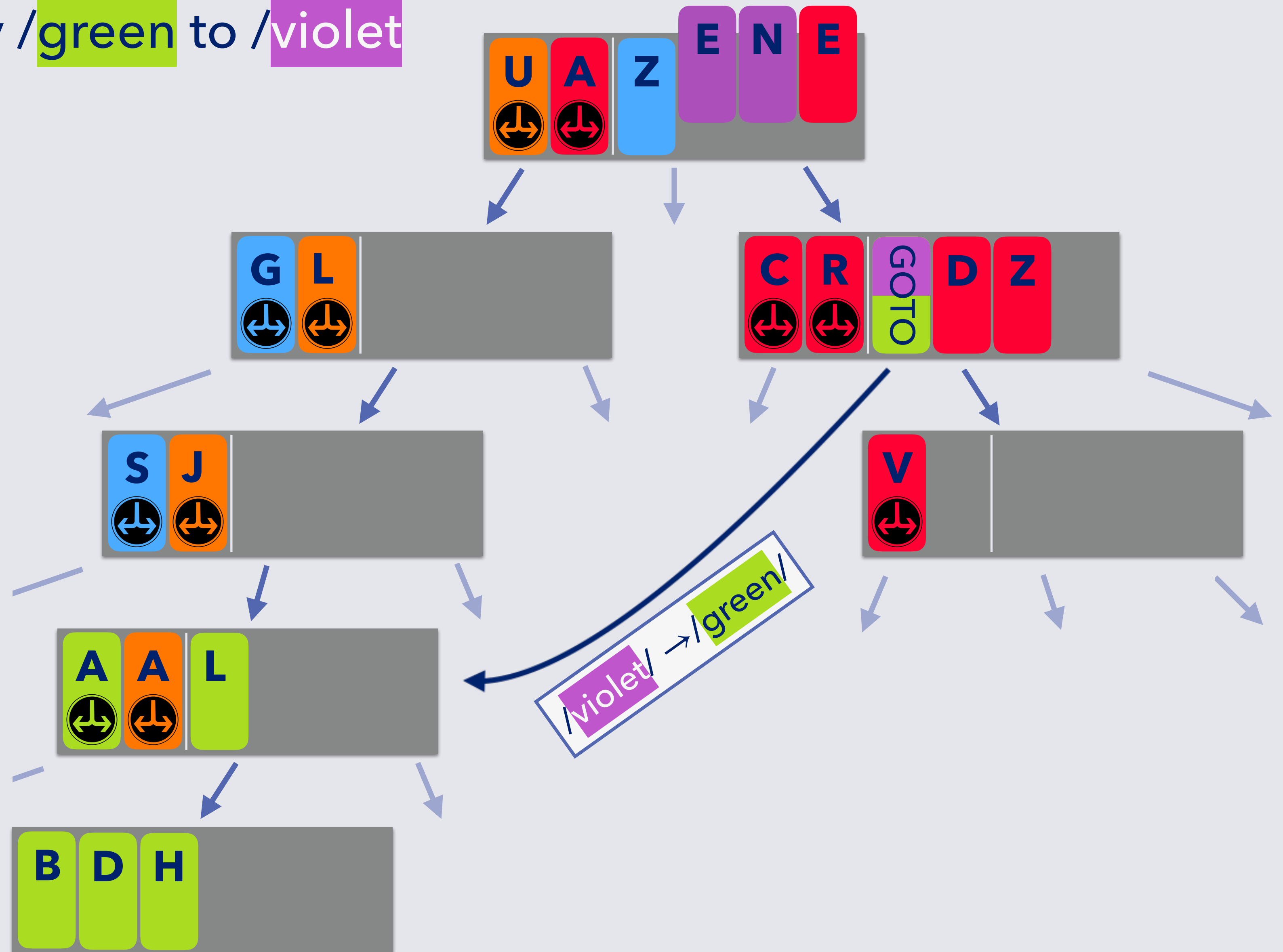
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

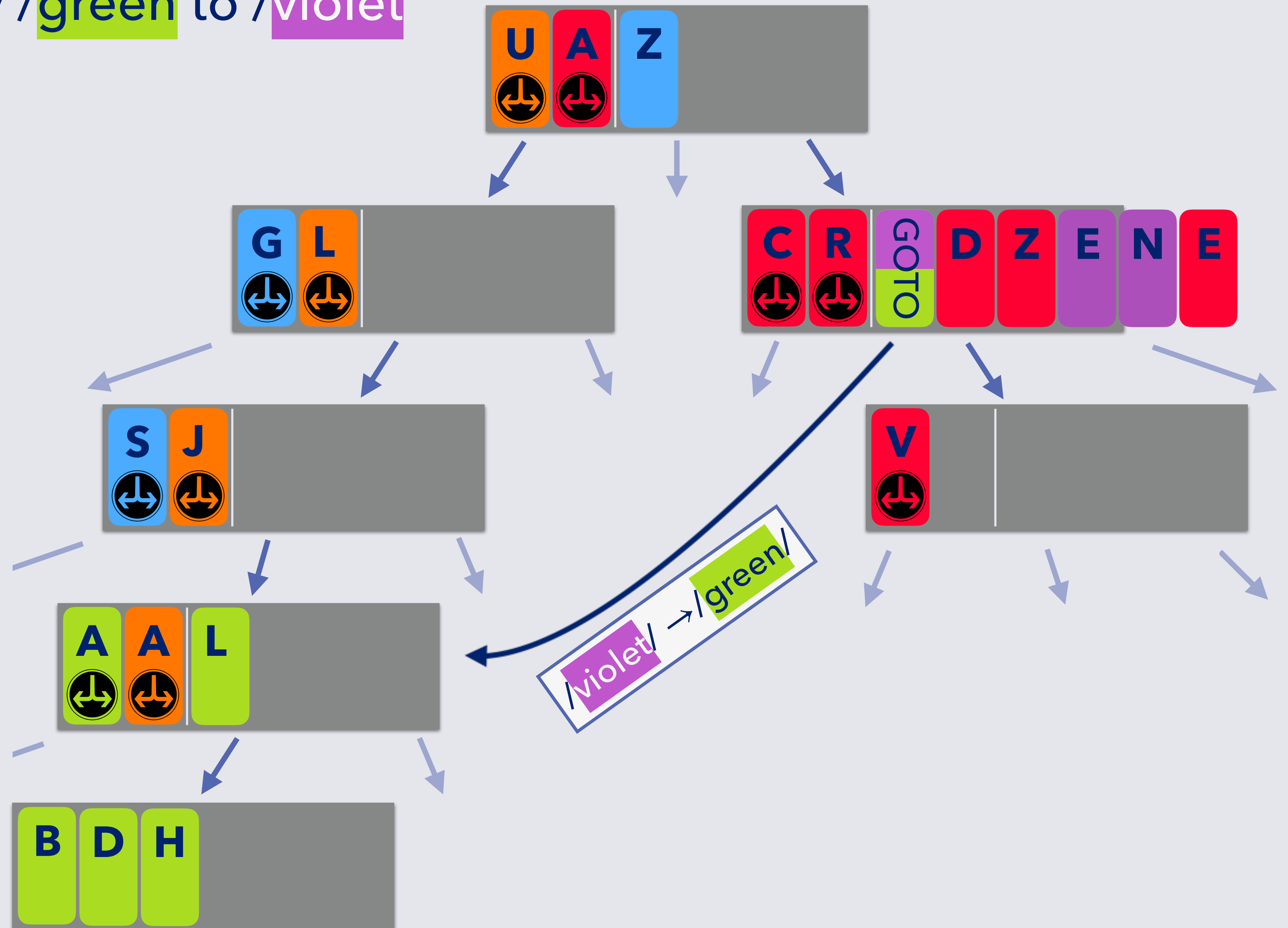
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B ϵ -DAGs

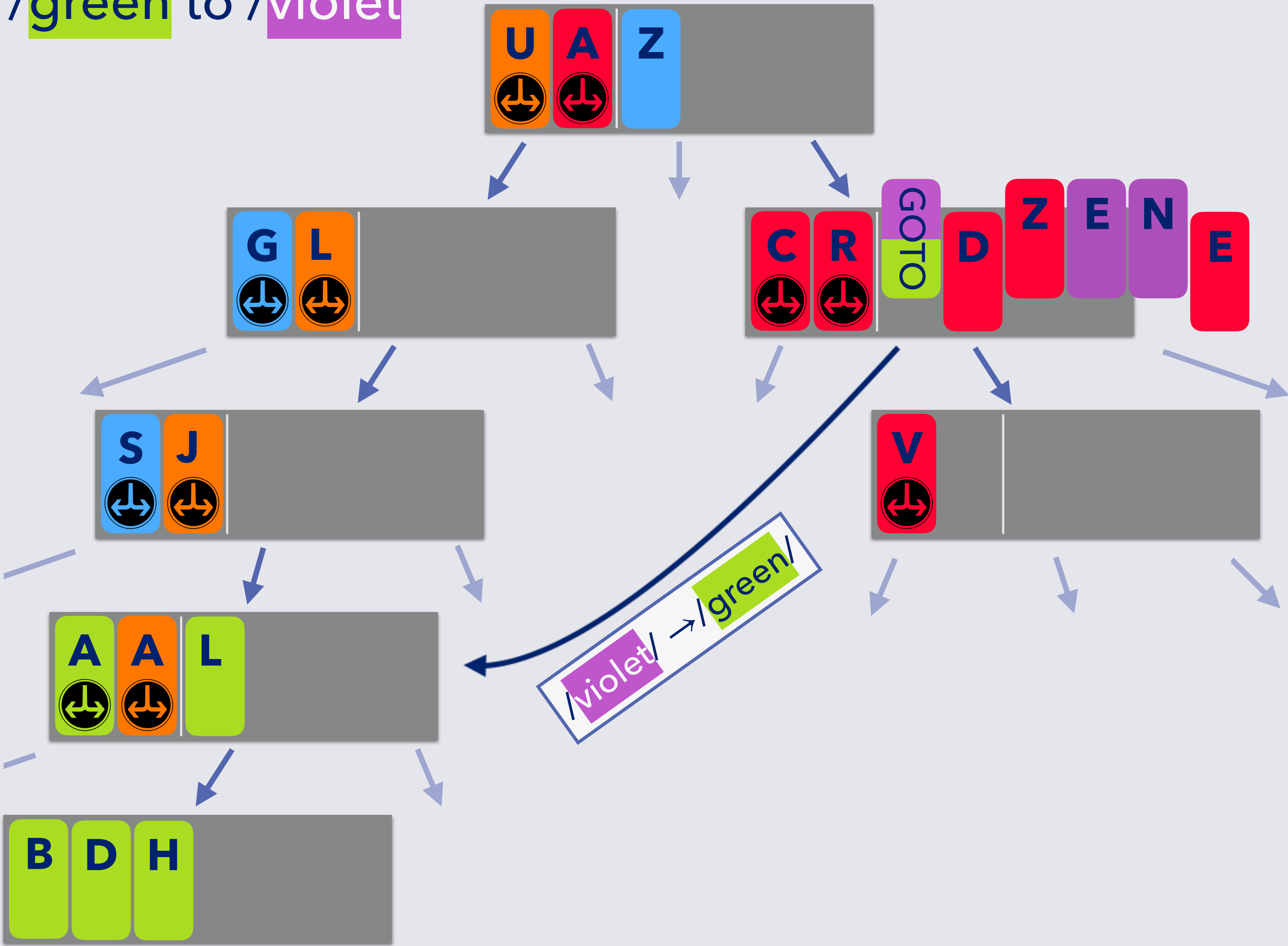
Copy /green to /violet

- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message



A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B^ϵ -DAGs

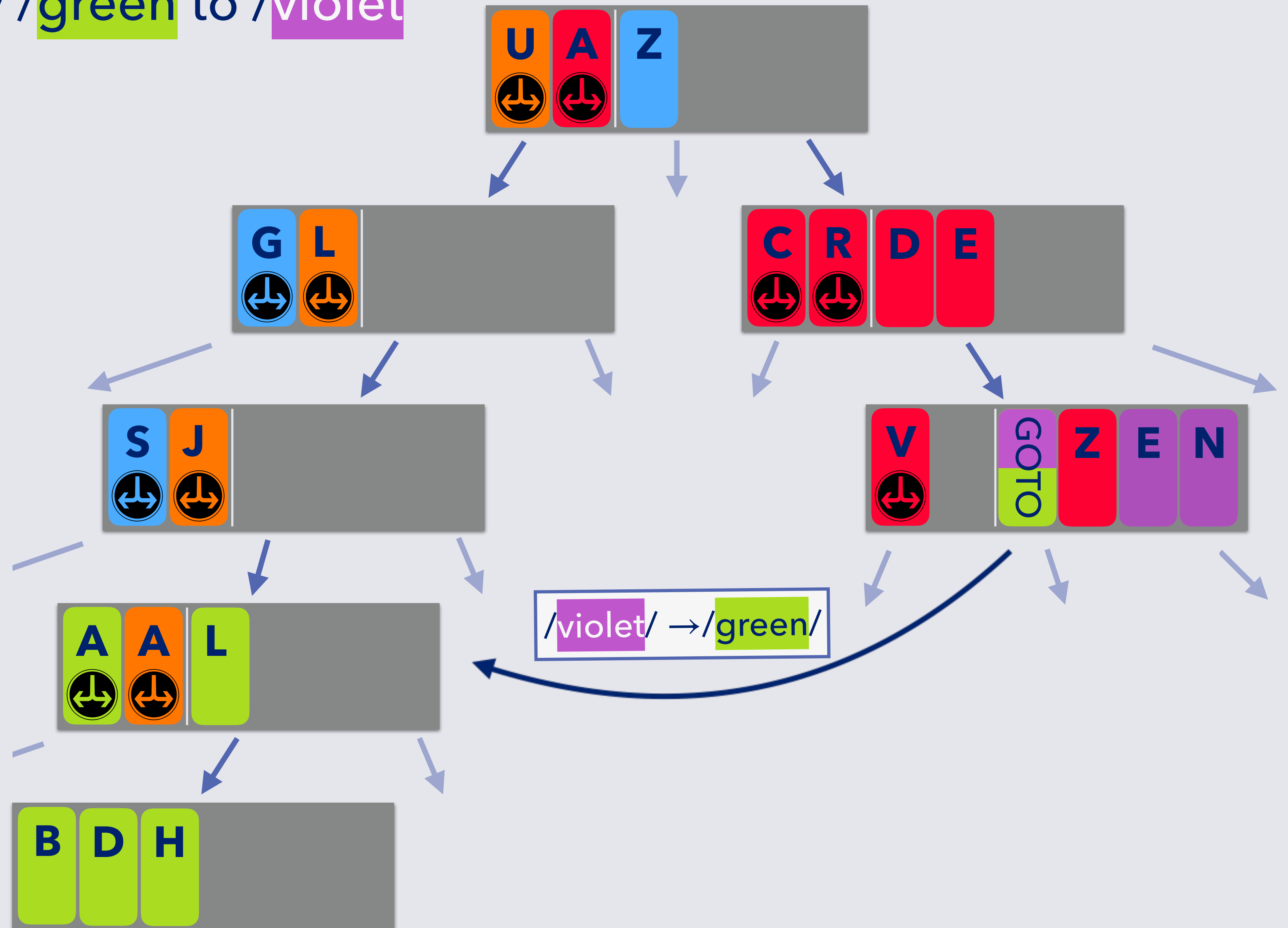
Copy /green to /violet

1. Flush messages to node covering /green subtree
2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B ϵ -DAGs

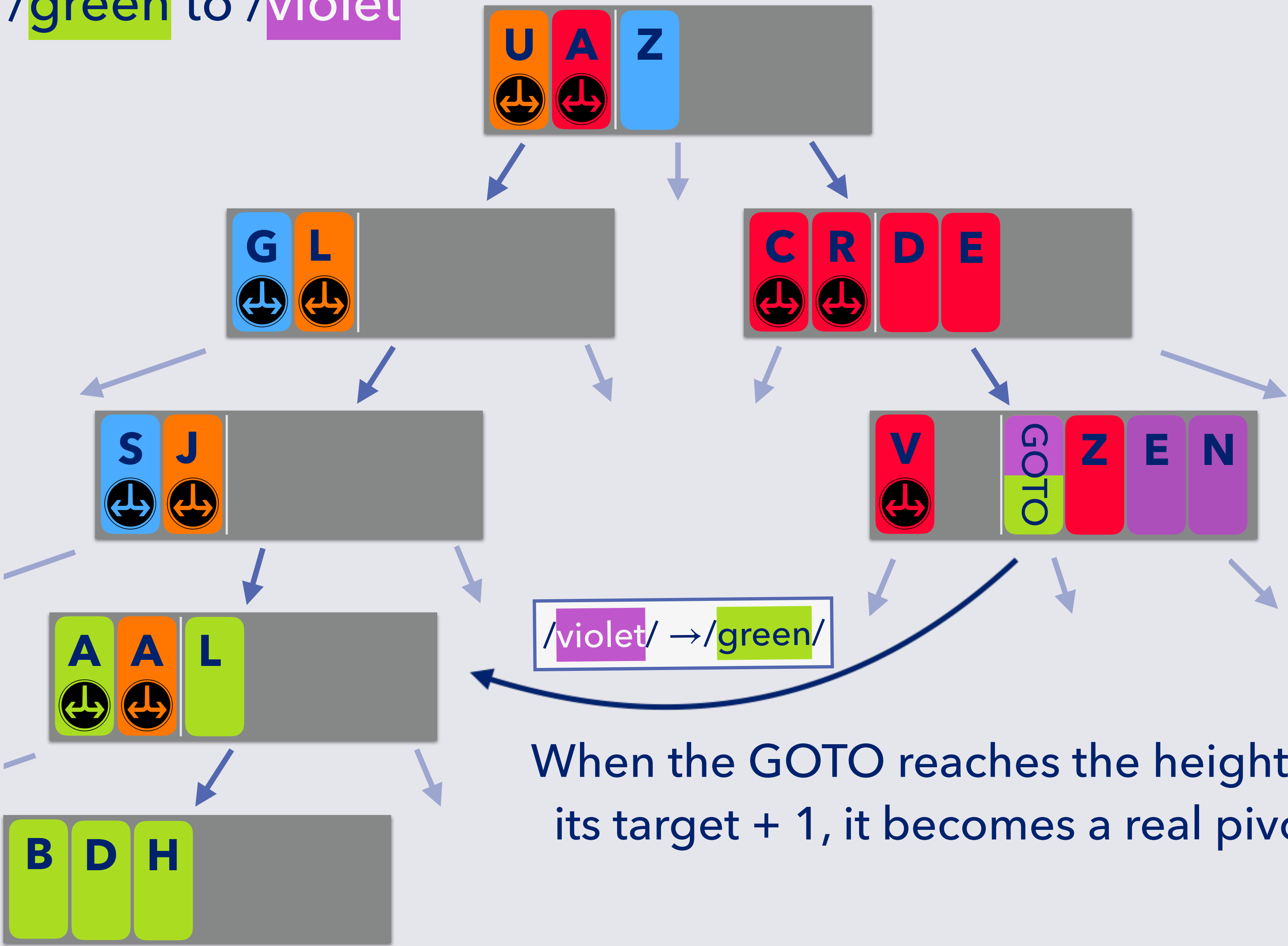
Copy /green to /violet

- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



When the GOTO reaches the height of its target + 1, it becomes a real pivot

Reducing Copy Latency in B ϵ -DAGs

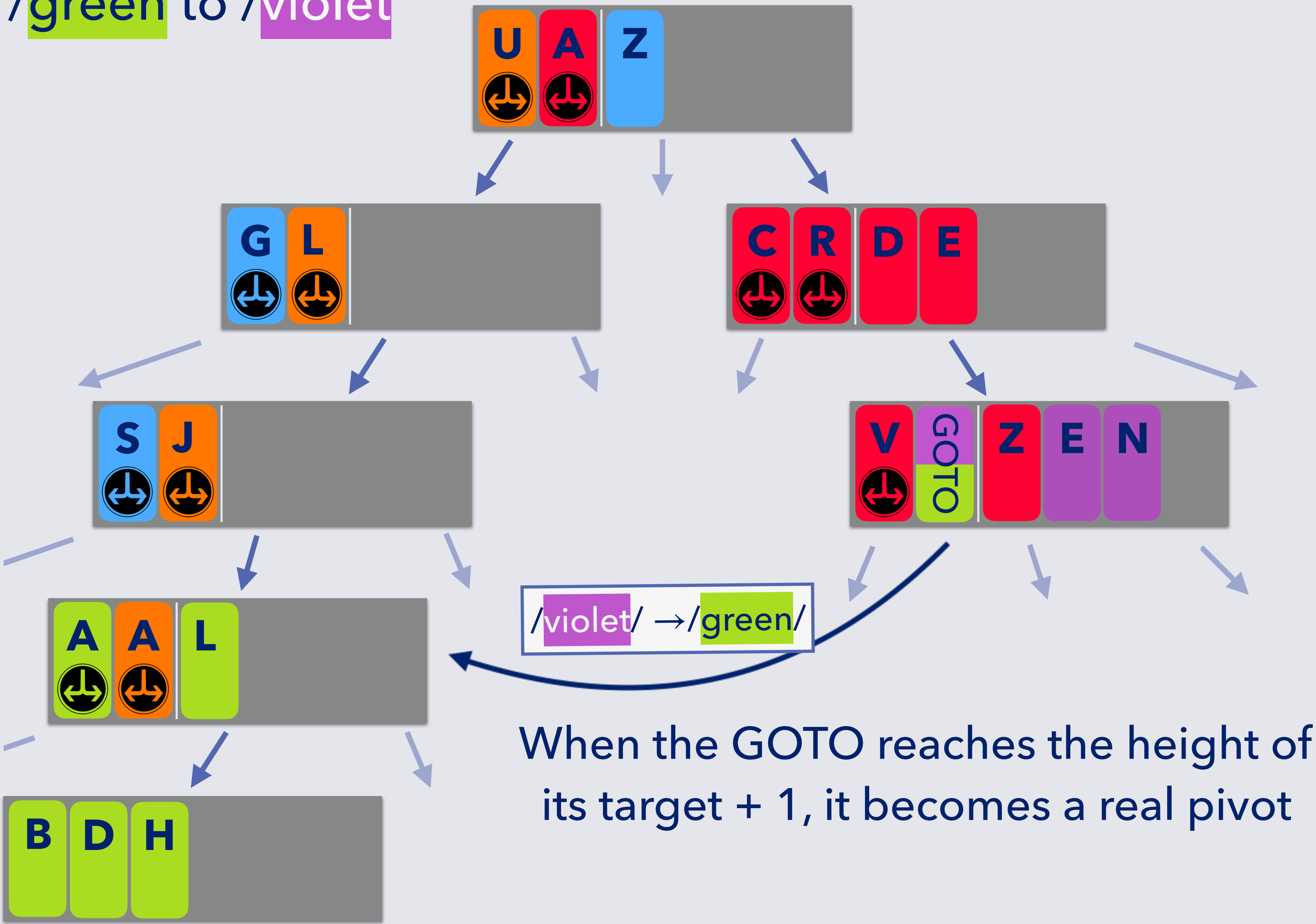
Copy /green to /violet

- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



Reducing Copy Latency in B ϵ -DAGs

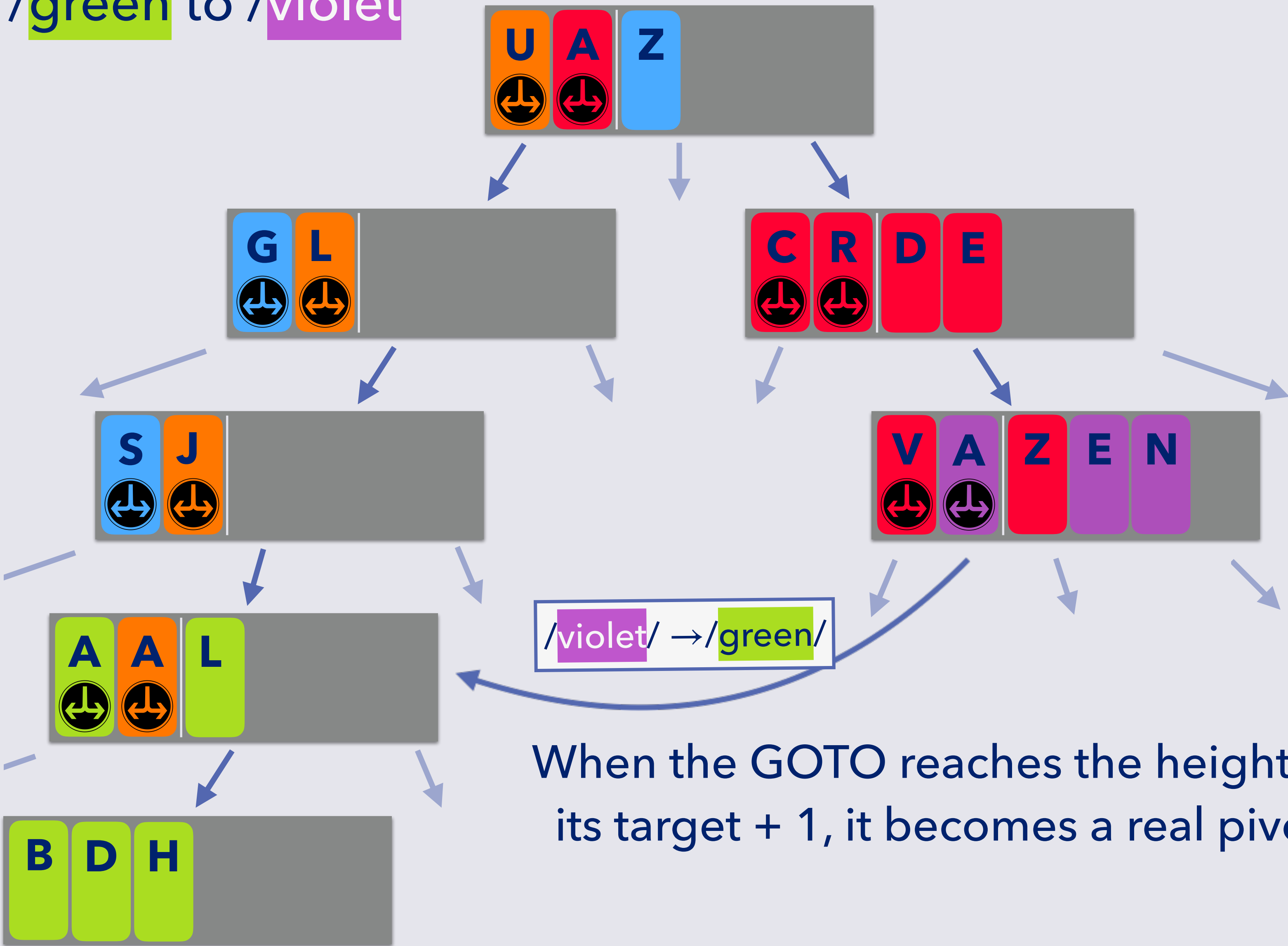
Copy /green to /violet

- 1. Flush messages to node covering /green subtree
- 2. Insert a GOTO message

GOTO

A GOTO message changes the structure of the tree

It functions as a pivot including prefix translation



When the GOTO reaches the height of its target + 1, it becomes a real pivot

Logical Copy with B^ϵ -DAGs

Performance Goals

Copy-on-Abundant-Write

Space efficient



Low latency



Copy-specific

B^ϵ -trees have large nodes

Fast writes



\Rightarrow Locality*

Fast reads



General file system

**File Systems Fated for Senescence? Nonsense, Says Science!,
Conway et al, FAST 2017*

Logical Copy with B^ϵ -DAGs

Performance Goals

Copy-on-Abundant-Write

Space efficient



GOTO messages

Low latency



Copy-specific

B^ϵ -trees have large nodes

Fast writes



\Rightarrow Locality*

Fast reads



General file system

**File Systems Fated for Senescence? Nonsense, Says Science!,
Conway et al, FAST 2017*

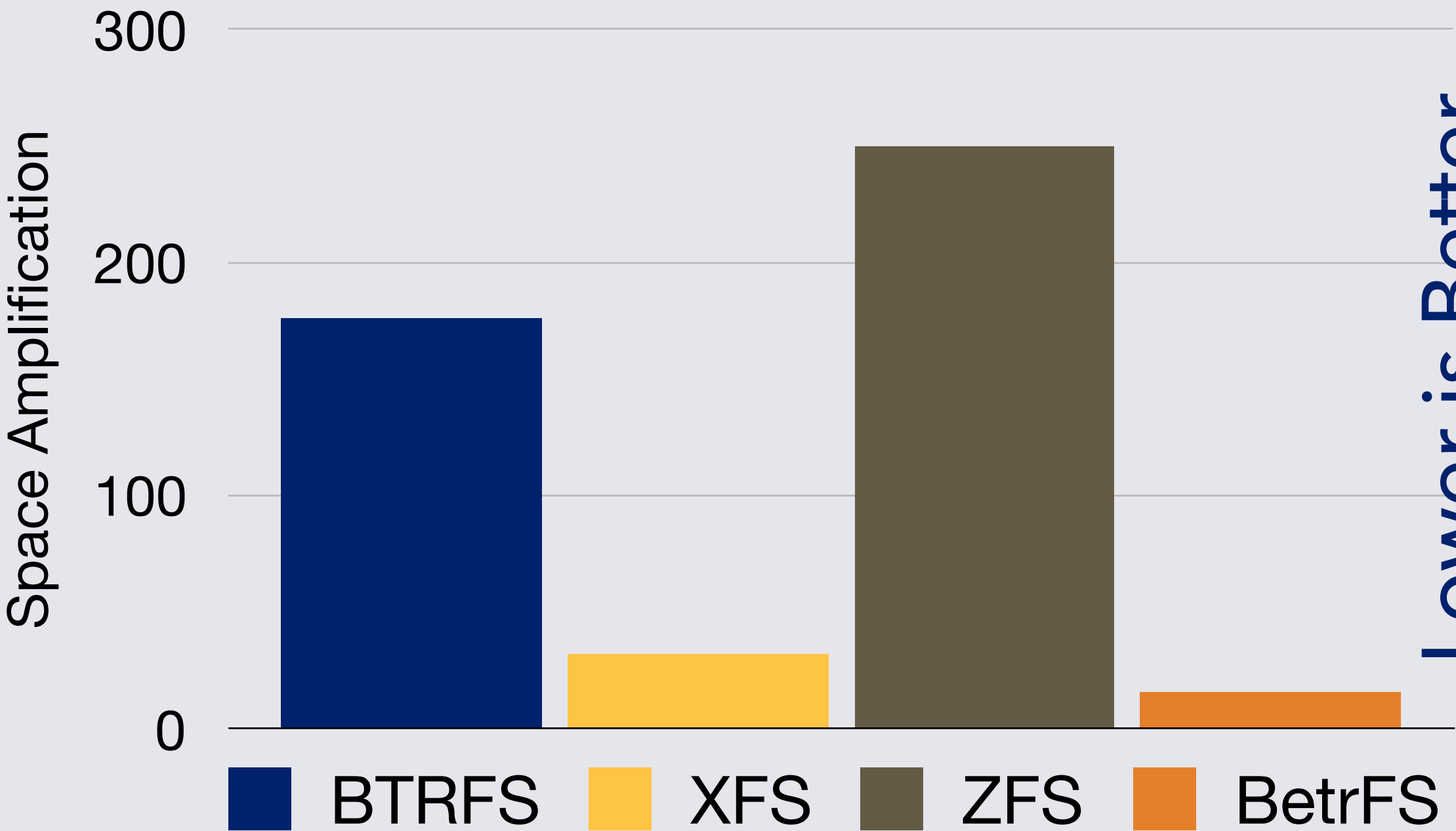
Evaluation

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

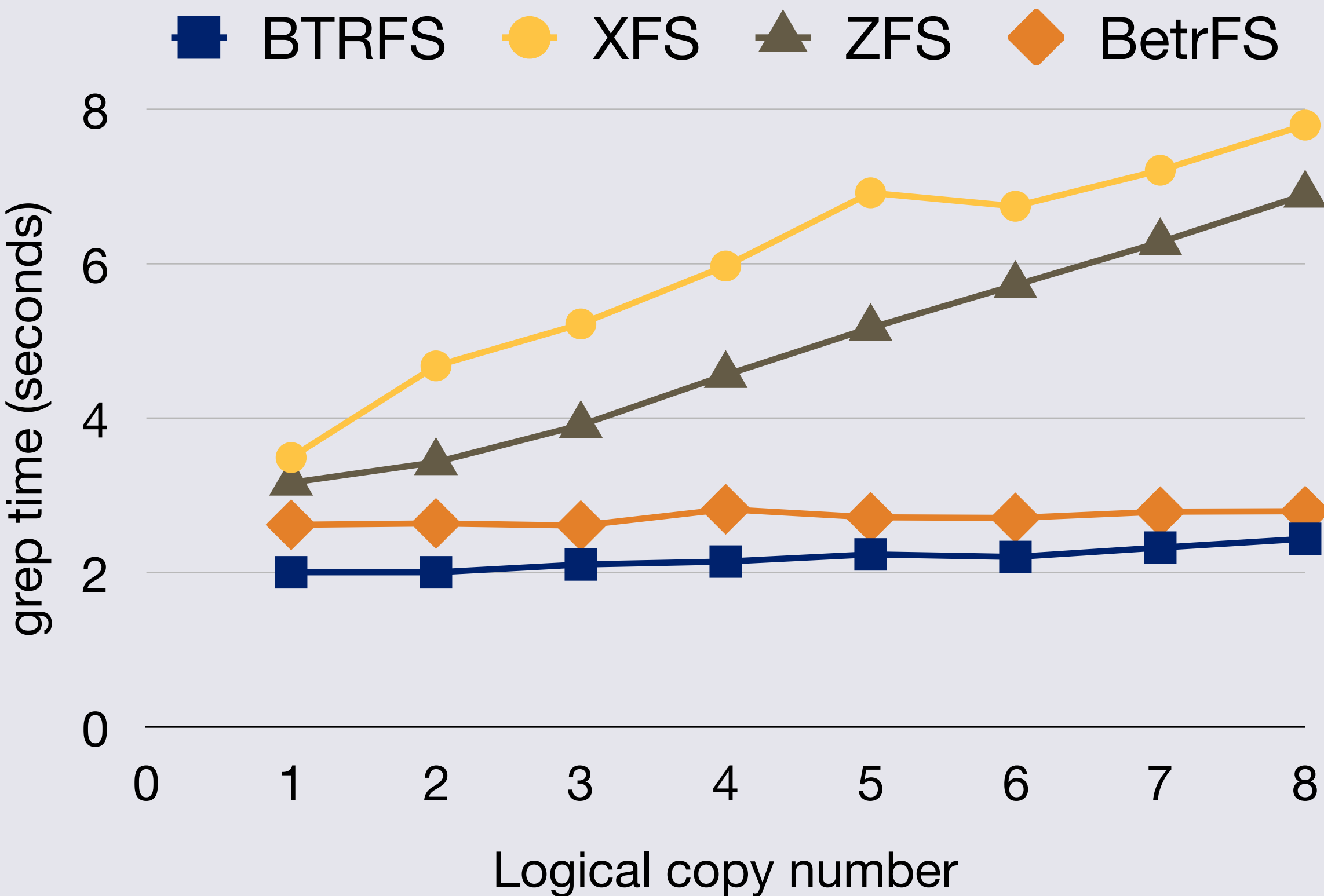
Init: 64 4MiB files with random data.

Each round: logically copy all files, then change 16B in each file (1KiB total)



Space Amplification

additional file system size / added data (1KiB)



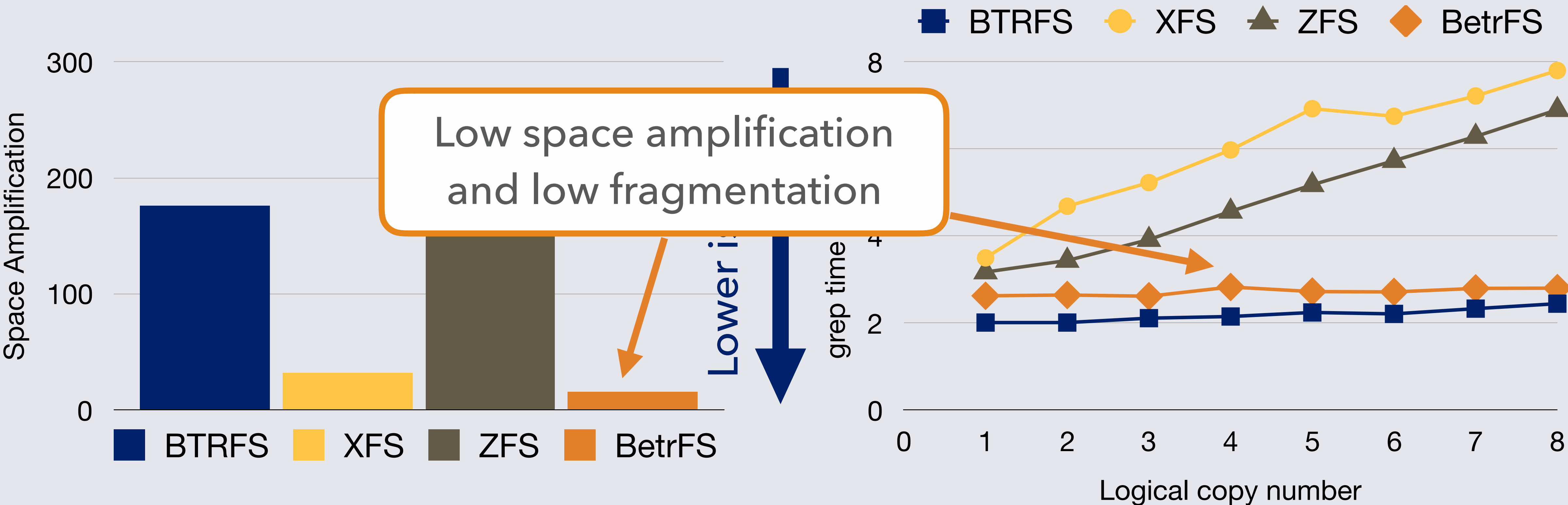
Fragmentation

measured by timing a grep over the file system

Space Amplification and Fragmentation

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

Init: 64 4MiB files with random data.
Each round: logically copy all files, then change 16B in each file (1KiB total)



Low space amplification
and low fragmentation

Space Amplification

additional file system size / added data (1KiB)

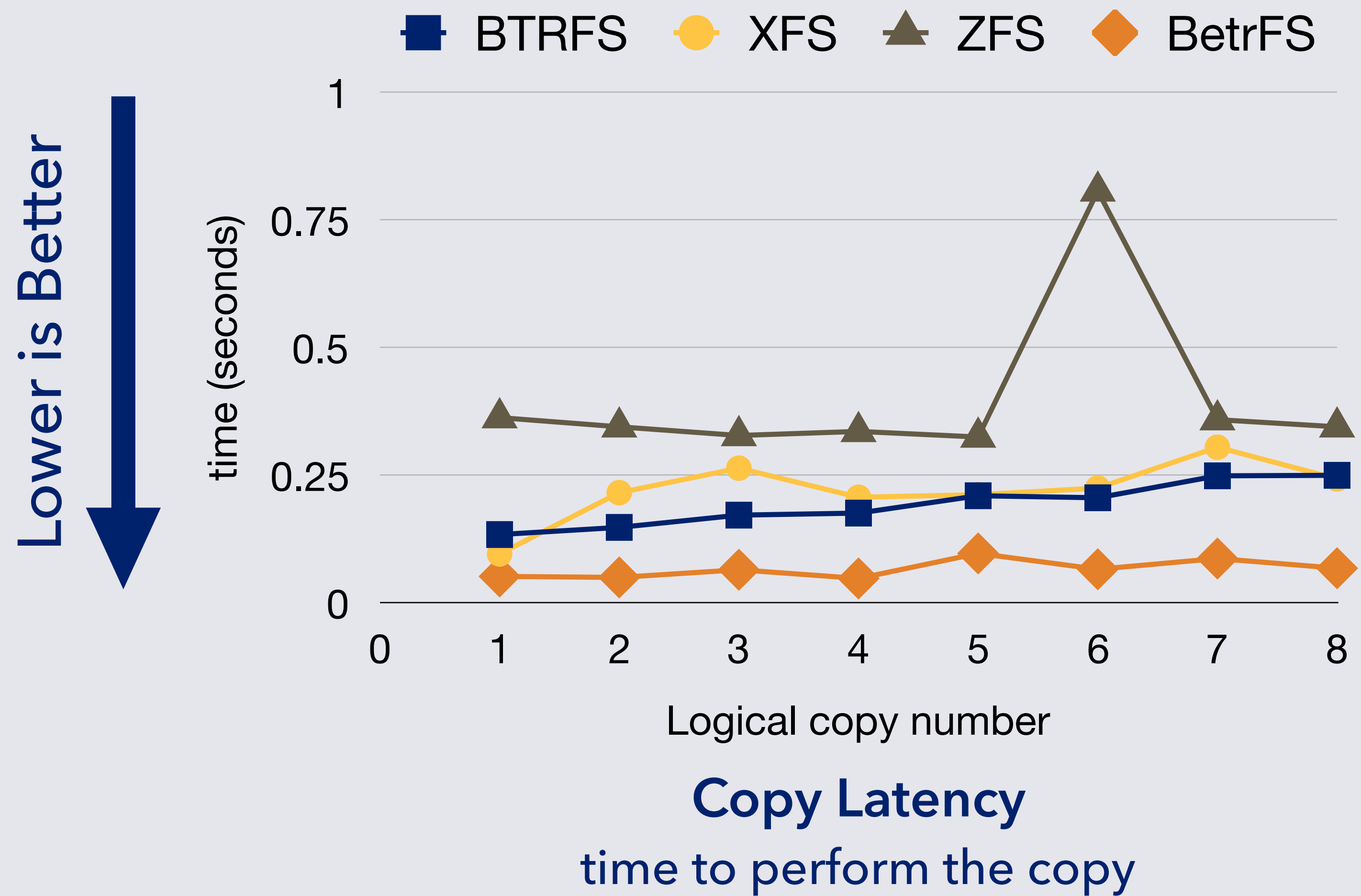
Fragmentation

measured by timing a grep over the file system

Copy Latency

Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

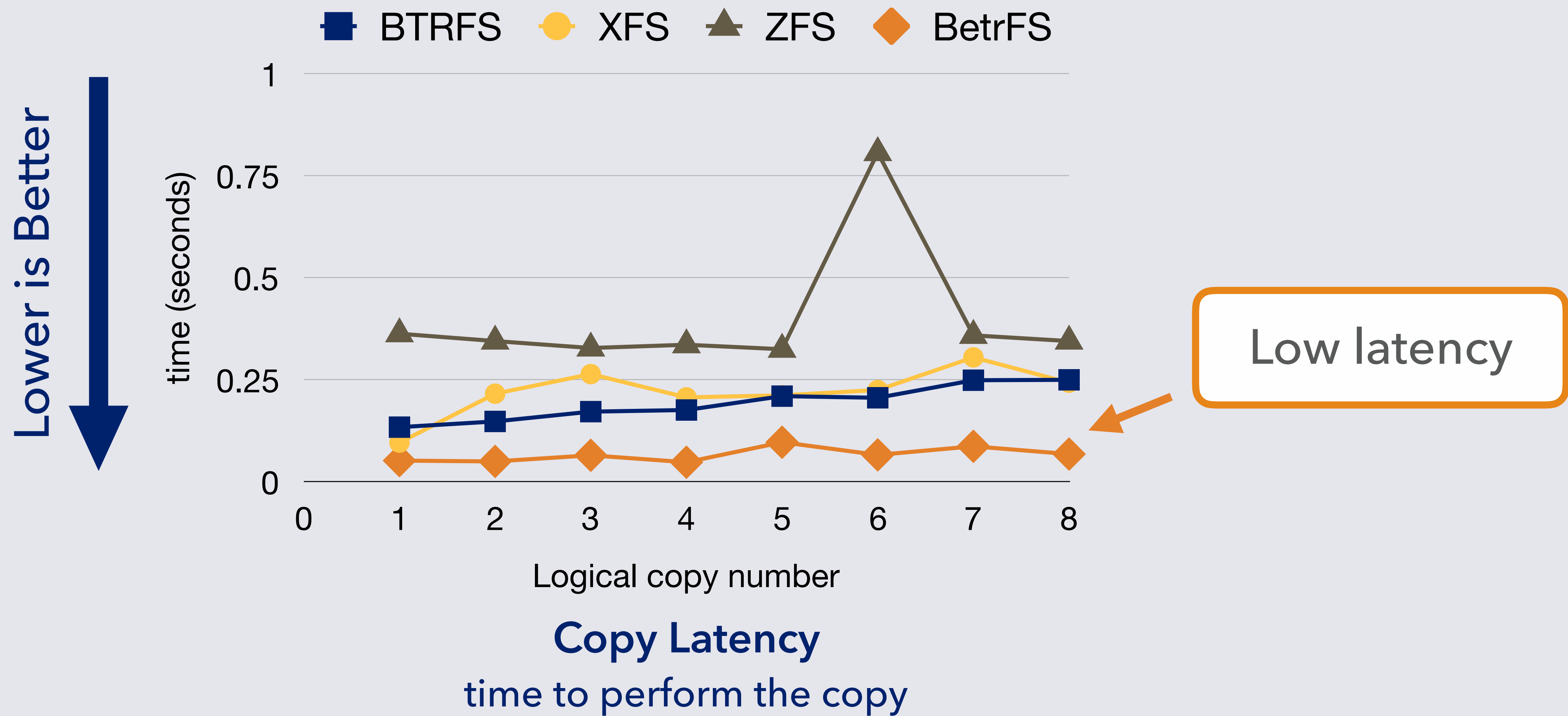
Init: 64 4MiB files with random data.
Each round: logically copy all files, then change 16B in each file (1KiB total)



Copy Latency

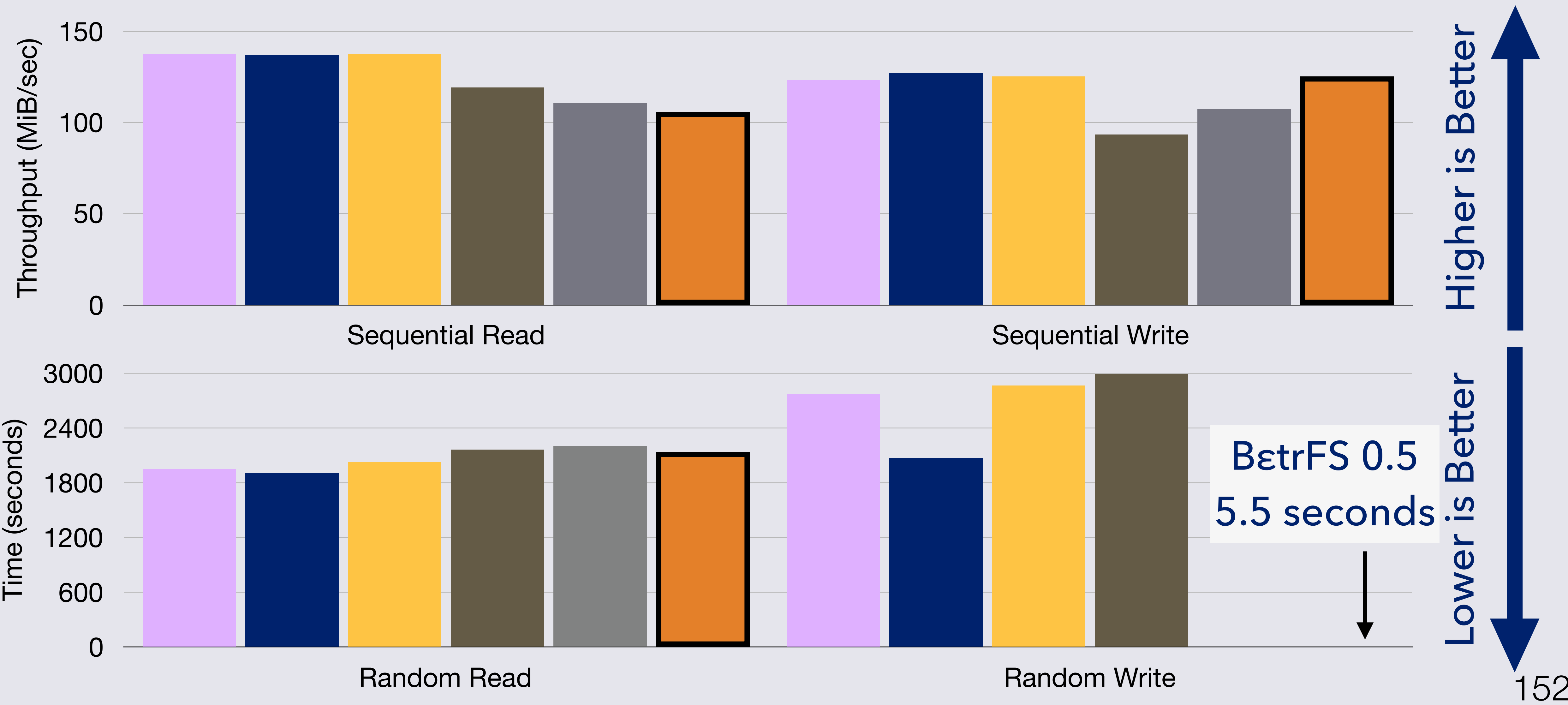
Dell Optiplex 790
4-core 3.40 GHz Intel Core i7 CPU
4GiB RAM
500GB 7200 RPM SATA disk

Init: 64 4MiB files with random data.
Each round: logically copy all files, then change 16B in each file (1KiB total)

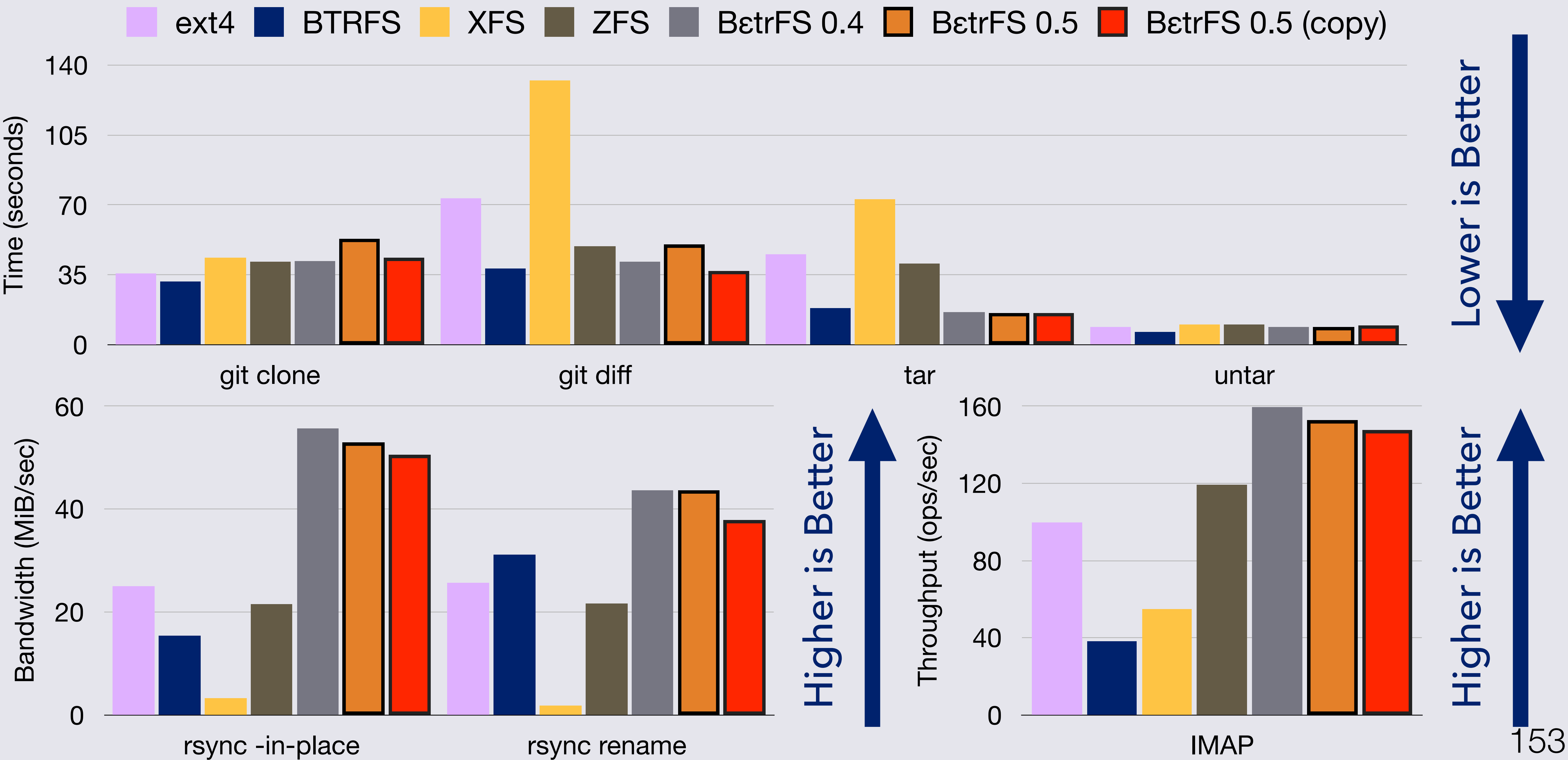


General File System Microbenchmarks

ext4 BTRFS XFS ZFS BεtrFS 0.4 BεtrFS 0.5



Application Benchmarks



Technical Conclusions

B^ϵ -DAGs transform copy-on-write into copy-on-abundant-write

- Gives strong bounds on space amplification
- Preserves locality, even with small writes
- Exploits B^ϵ -tree's batching and flushing

B^ϵ -DAGs preserve the fast reads and writes of B^ϵ -trees

- Preserve logarithmic tree height and query cost
- Preserve asymptotic costs of inserts and updates

Copies are fast and cheap

- GOTO messages enable low-latency DAG mutations
- Total work of copies is $O(\text{tree height})$

Evaluation Conclusions

BetrFS with B^ϵ -DAGs has strong copy performance in practice

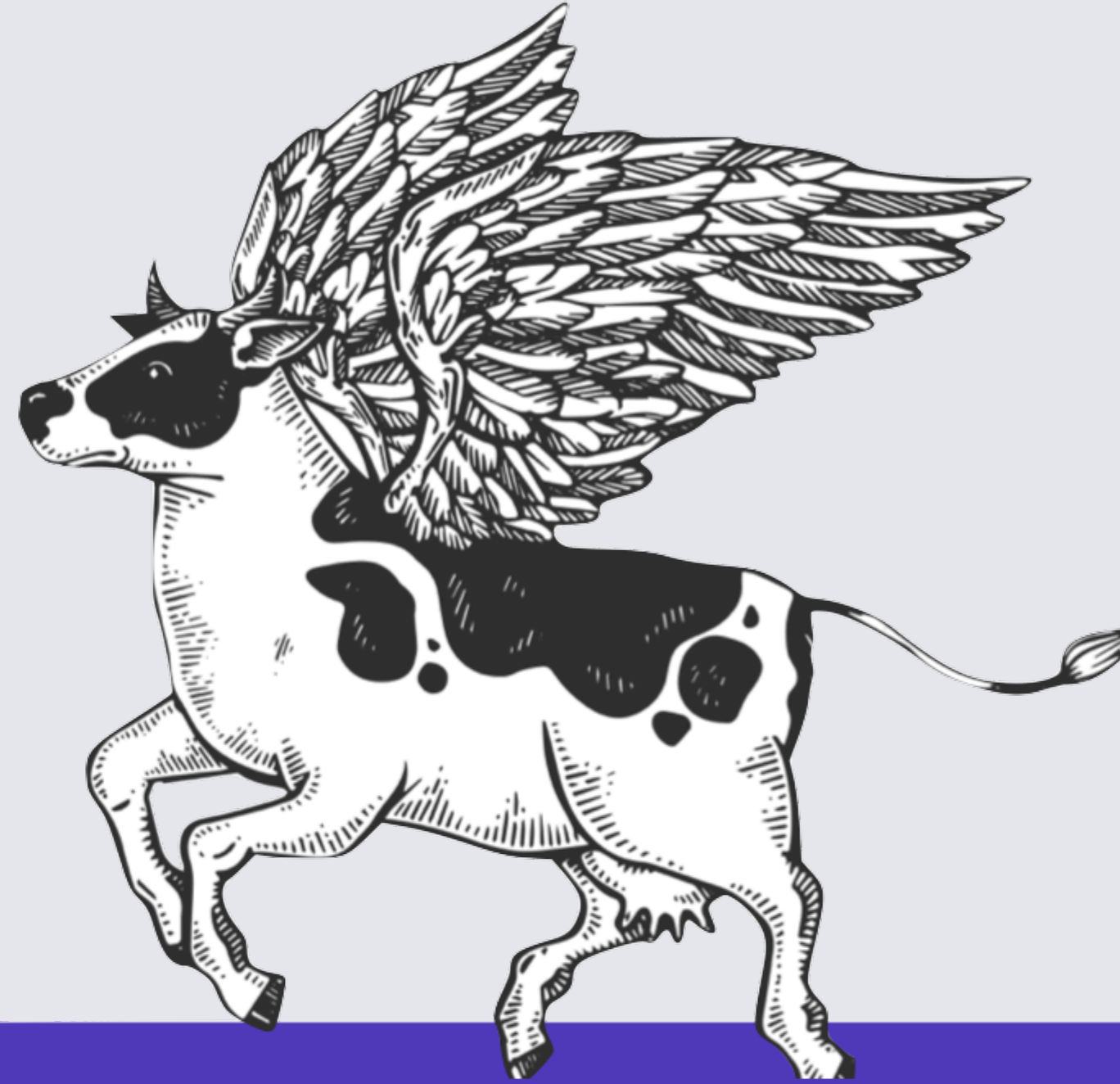
- Low space amplification
- Low latency copying
- Good locality

B^ϵ -DAGs preserve BetrFS's performance gains on other operations

- Fast random writes
- Good sequential I/O throughput
- No aging
- Strong across-the-board application benchmark performance

Thank you!

The CAW Awakens



BetrFS Episode V: Attack of the Clones

O'Really?

betrfs.org