

Write-Optimization for File-System Metadata

Rob Johnson
VMware Research

ExxonMobil

“A user does ls -l on a directory with a million files”

Aaron Steichen

“Want fast find on trillions of files”



Chris Beecroft



“Flattening the namespace gives great performance but renames are expensive”

Dave Bonnie

File systems need to perform well on many metadata operations

	Query	Update
Point	<ul style="list-style-type: none">• stat	<ul style="list-style-type: none">• atime update• creat• unlink
Recursive	<ul style="list-style-type: none">• find• grep	<ul style="list-style-type: none">• rm -rf• recursive chmod

Full-path indexing

Write optimization

Custom data structures

rename

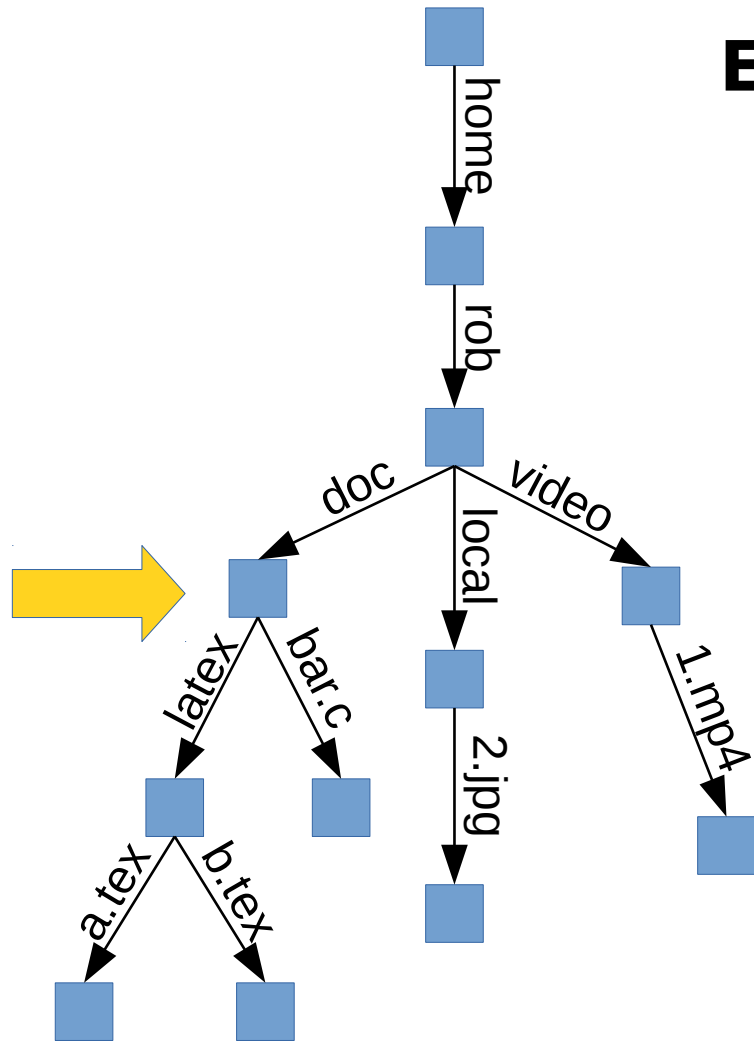
Full-path indexing

Full-path indexing in BetrFS

- BetrFS maintains two indexes
 - Metadata index: `path → struct stat`
 - Data index: `(path, blk#) → data[4096]`
- **Paths sorted in DFS order (i.e. full-path indexing)**
- Implications:
 - Data blocks laid out sequentially
 - **Directory scans → range queries**

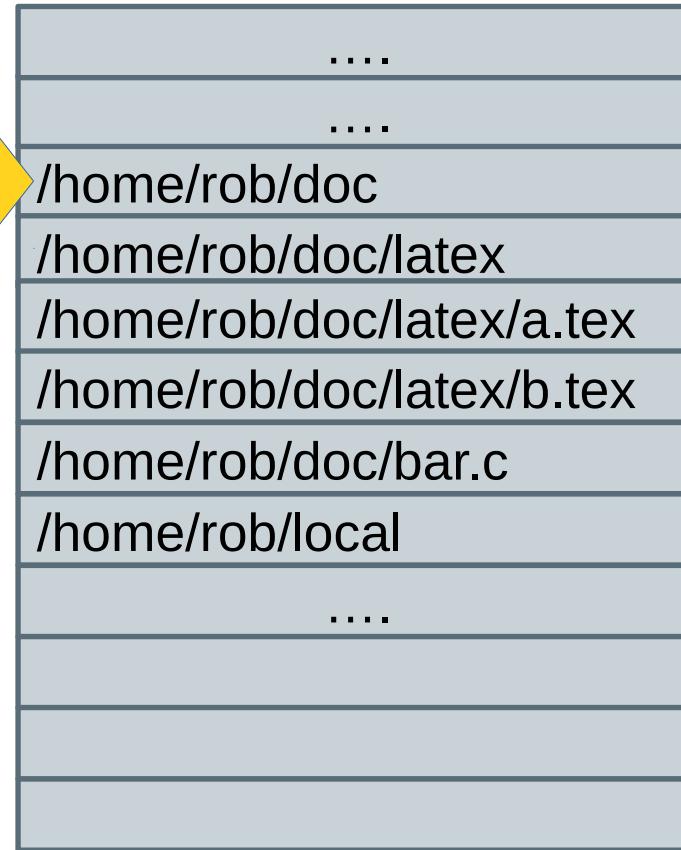
Full-path indexing yields fast directory scans

Example: `grep -r "key" /home/rob/doc/`



Directory Tree (logical)

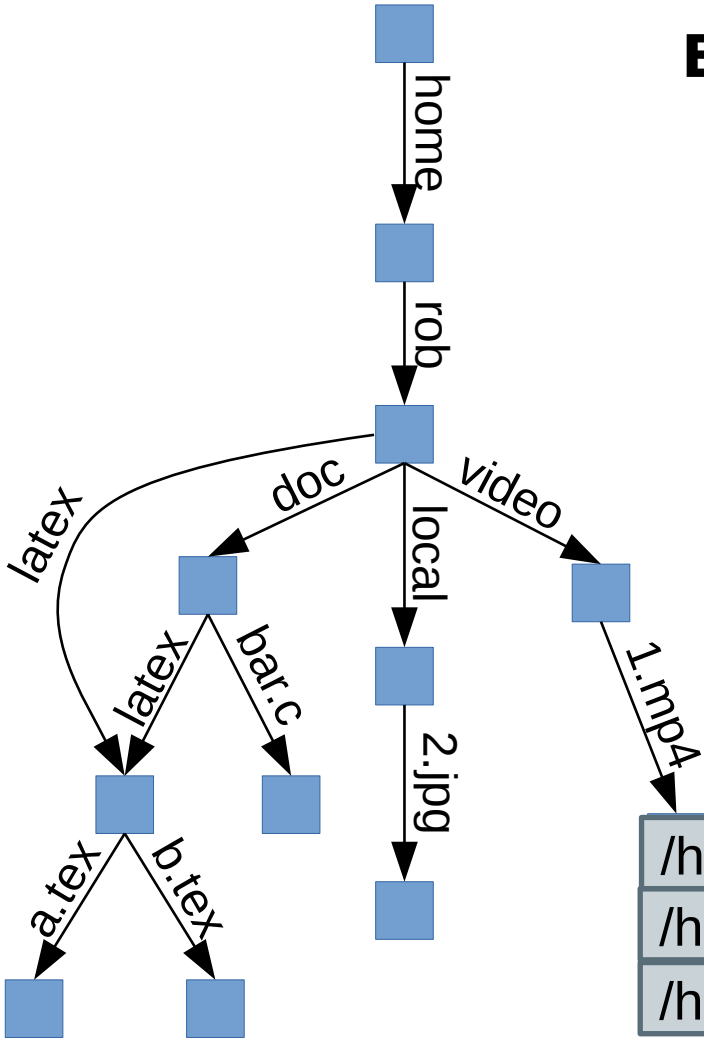
disk head



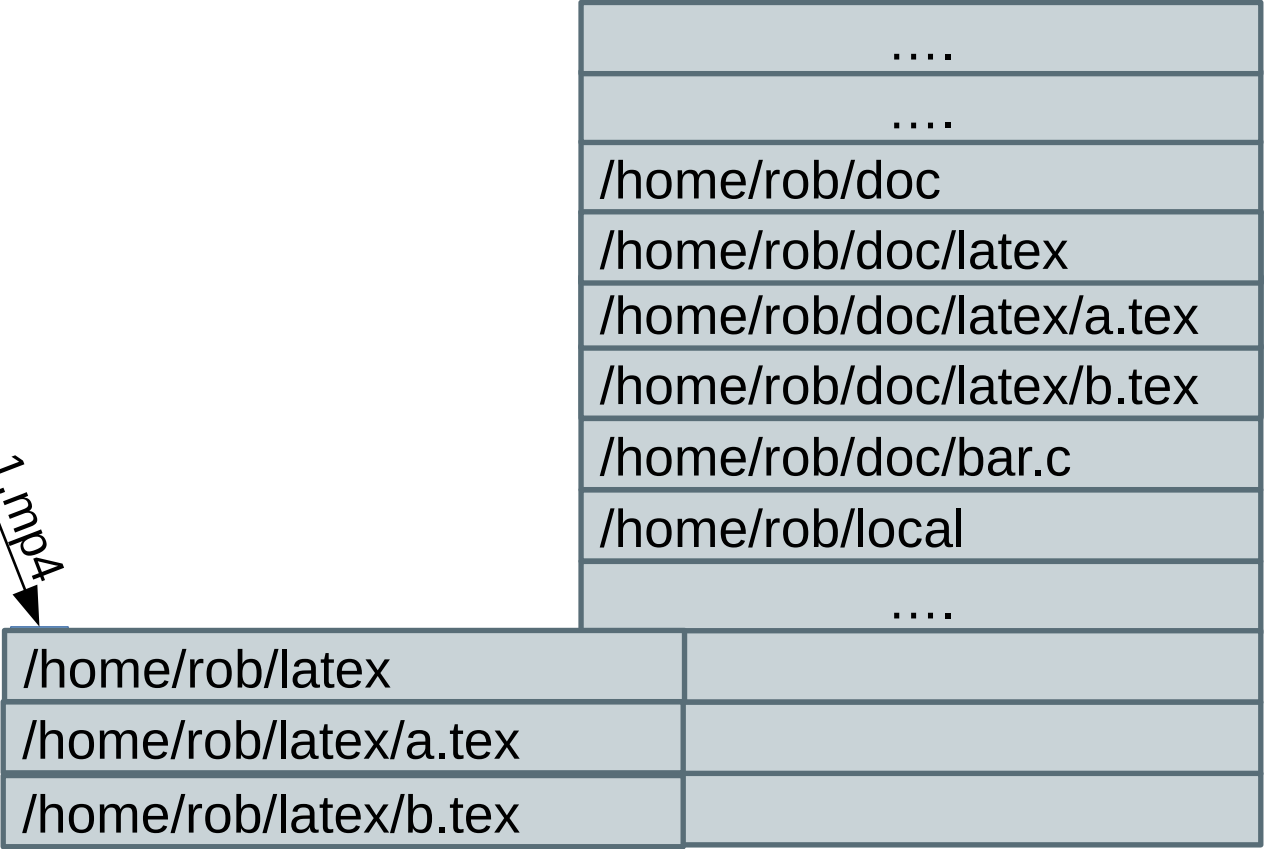
Disk (physical)

Rename is expensive when using full-path indexing

Example: `mv /home/rob/doc/latex /home/rob`



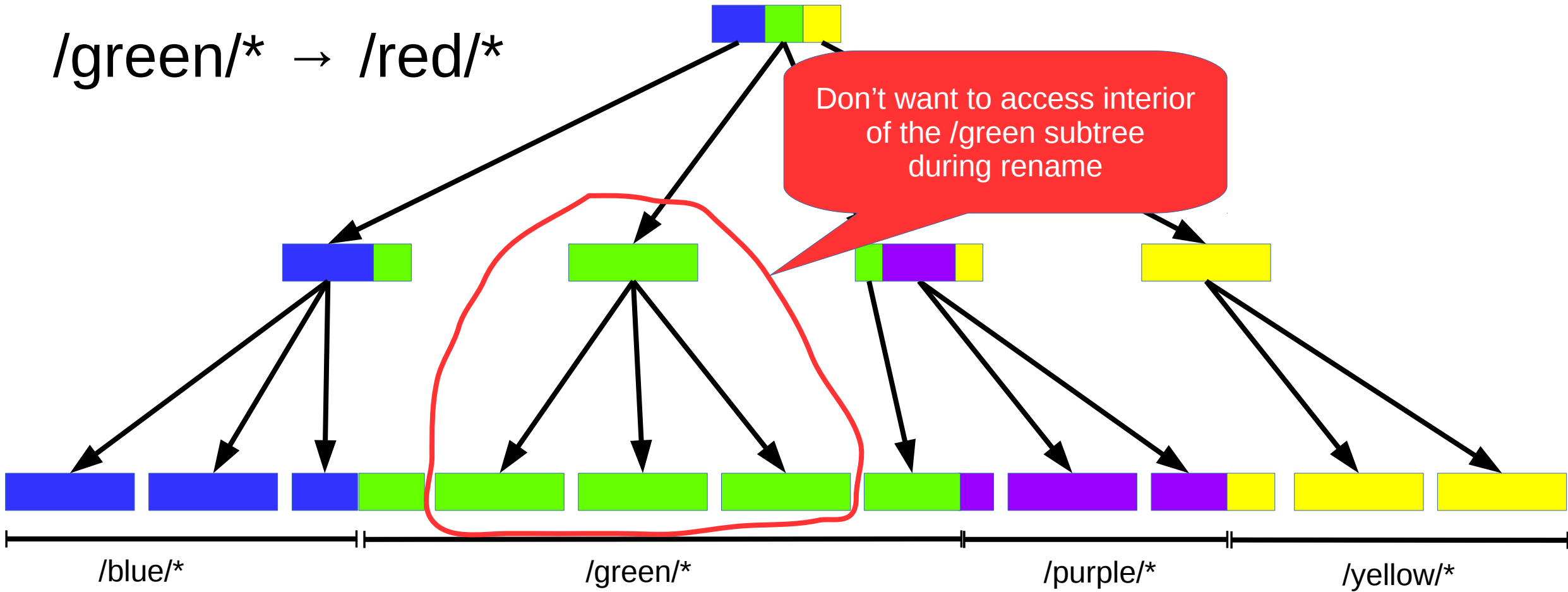
Directory Tree (logical)



Disk (physical)

Rename in BetrFS

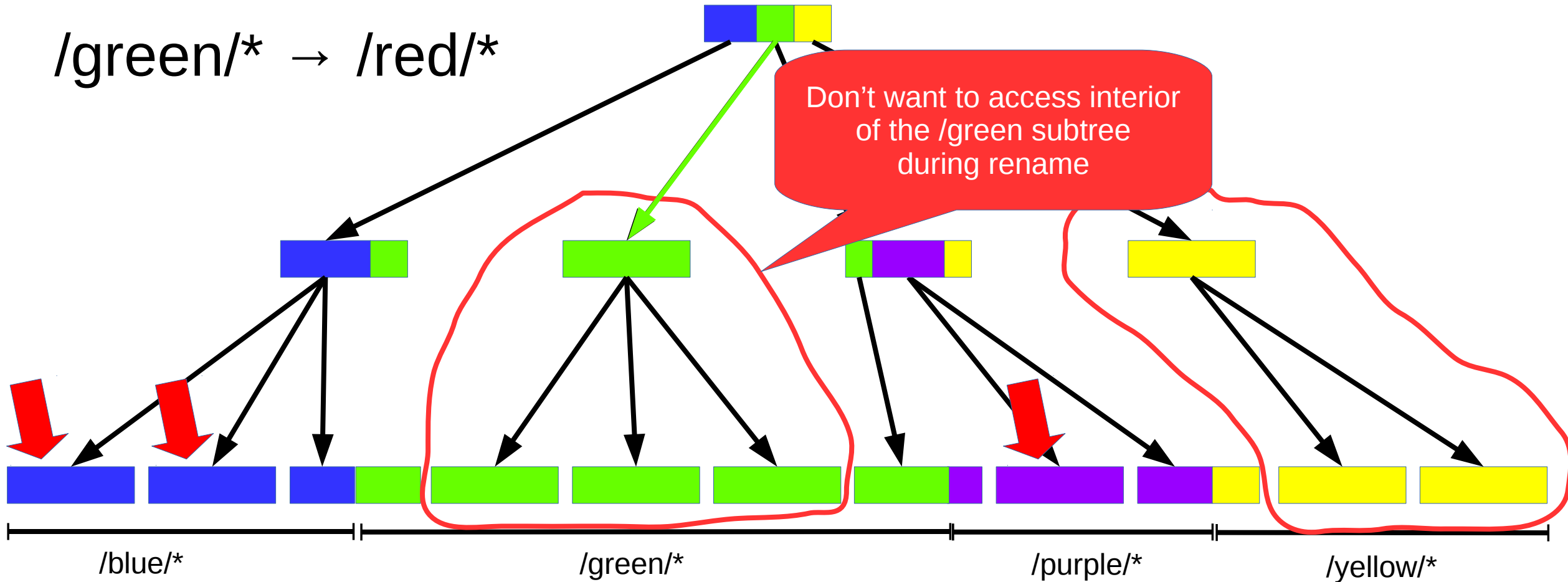
`/green/*` → `/red/*`



Lifted B^ϵ -trees

Idea: omit common prefixes from nodes and sub-trees

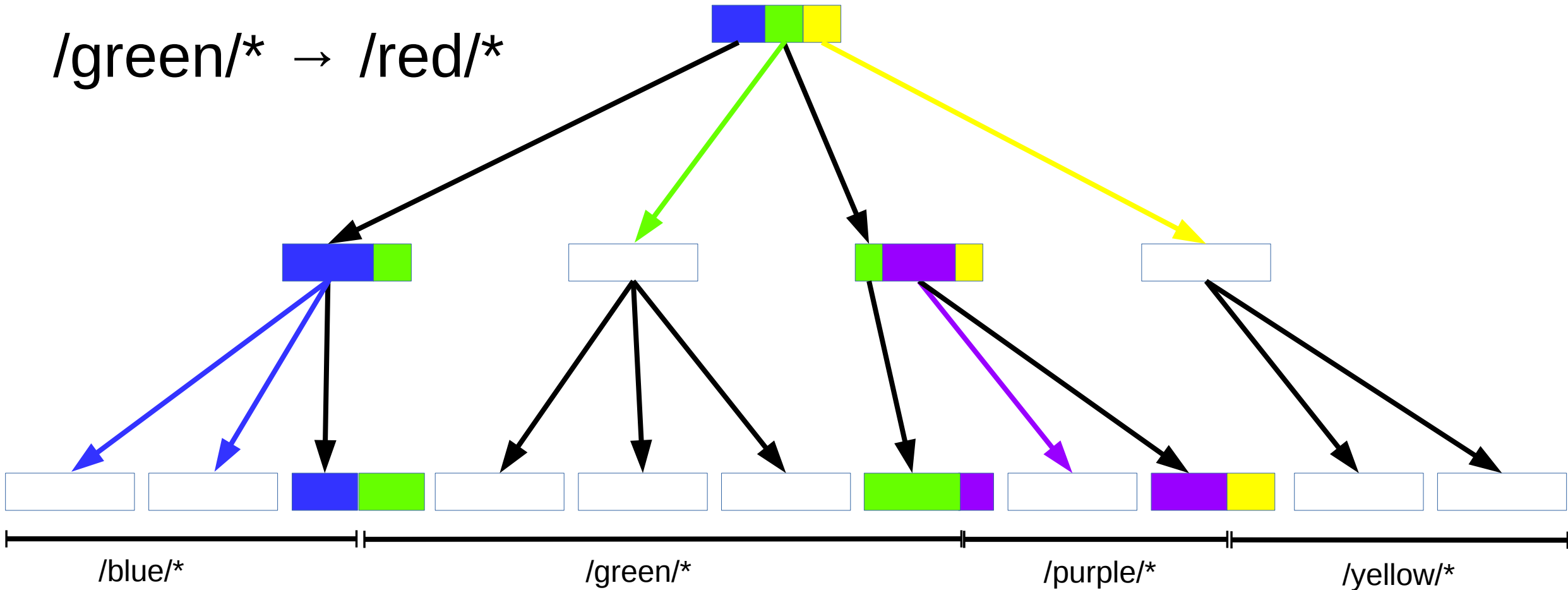
$/\text{green}/^* \rightarrow / \text{red}/^*$



Lifted B^ϵ -trees

Idea: omit common prefixes from nodes and sub-trees

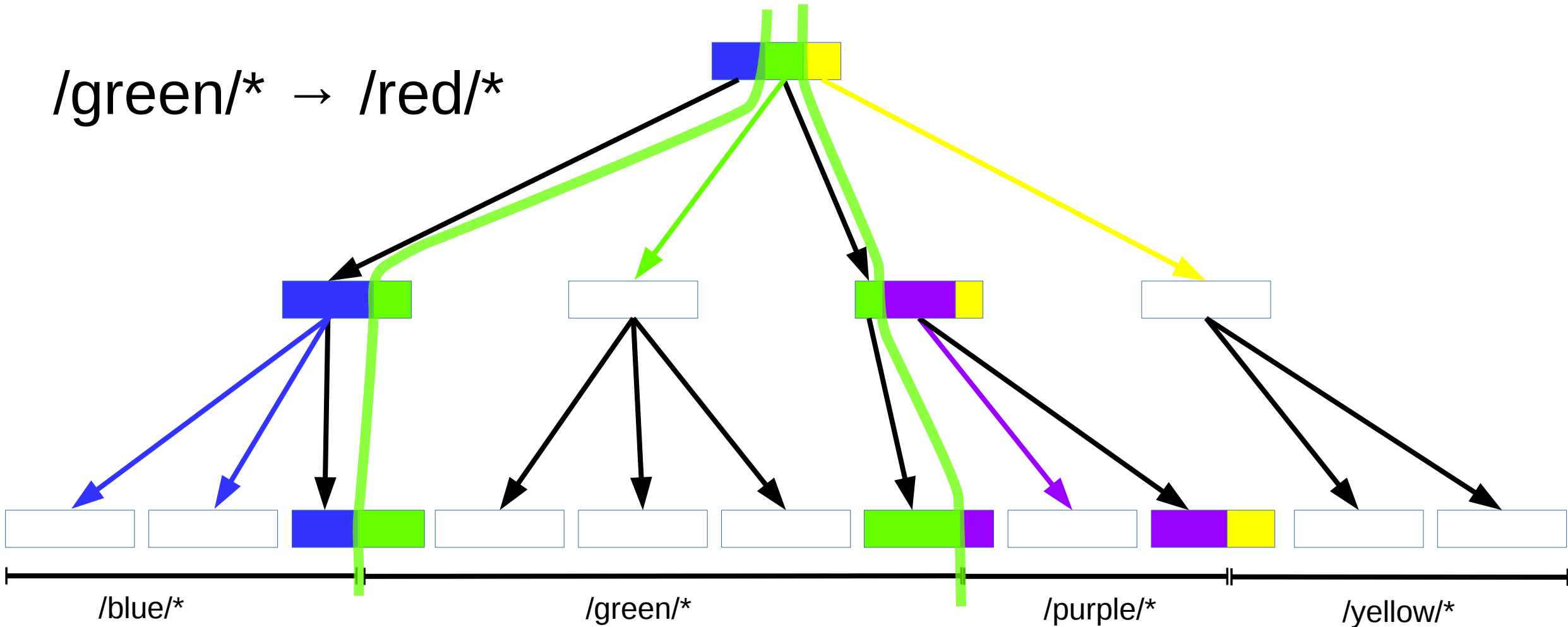
$/\text{green}/^* \rightarrow / \text{red}/^*$



Range rename in lifted B^ϵ -trees

1. Slice out the /green subtree

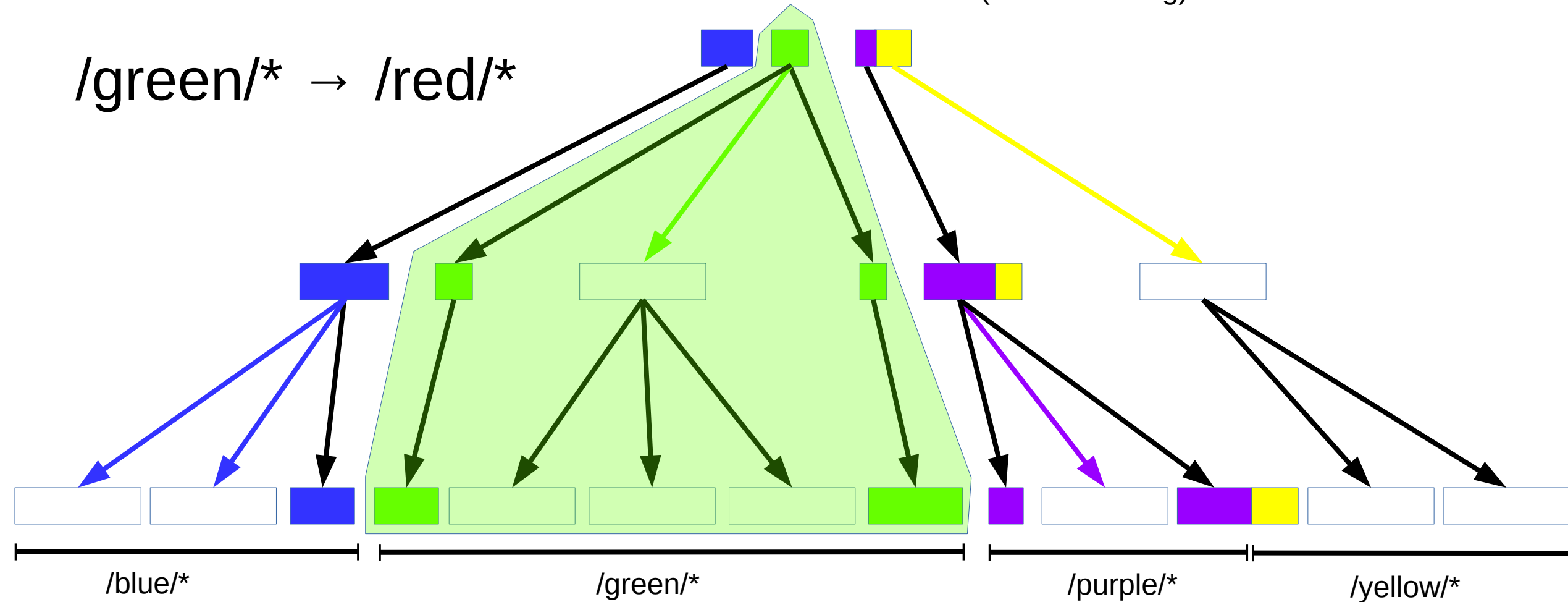
/green/* \rightarrow /red/*



Range rename in lifted B^ϵ -trees

1. Slice out the /green subtree
(maintain lifting)

/green/* \rightarrow /red/*

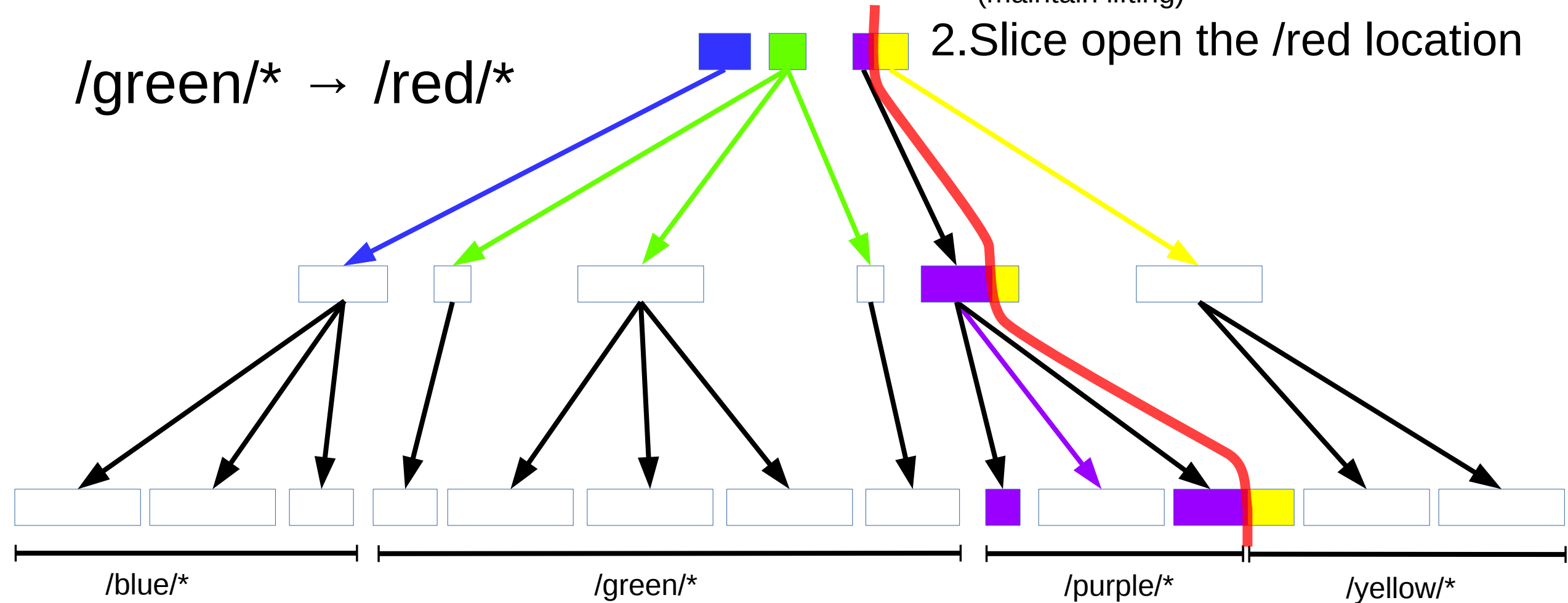


Range rename in lifted B^ϵ -trees

$/\text{green}/^* \rightarrow / \text{red}/^*$

1. Slice out the $/\text{green}/^*$ subtree
(maintain lifting)

2. Slice open the $/\text{red}/^*$ location

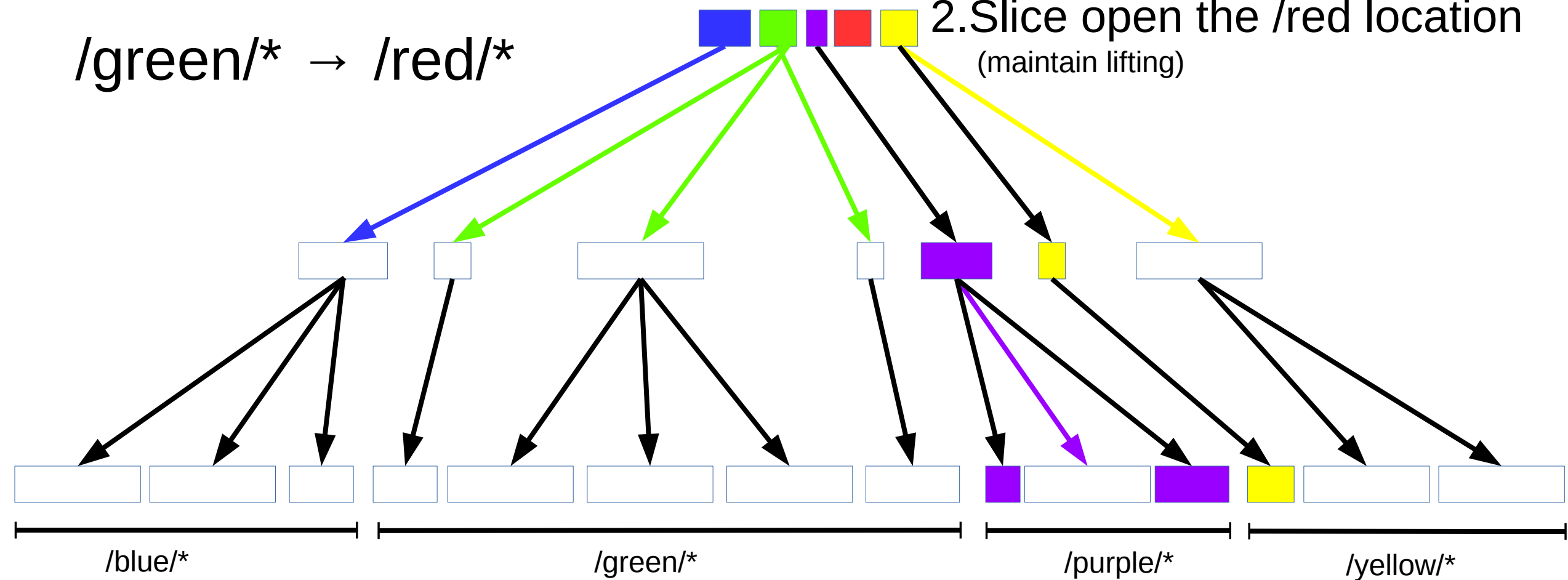


Range rename in lifted B^ϵ -trees

$/\text{green}/^* \rightarrow / \text{red}/^*$

1. Slice out the $/\text{green}/$ subtree
(maintain lifting)

2. Slice open the $/\text{red}/$ location
(maintain lifting)



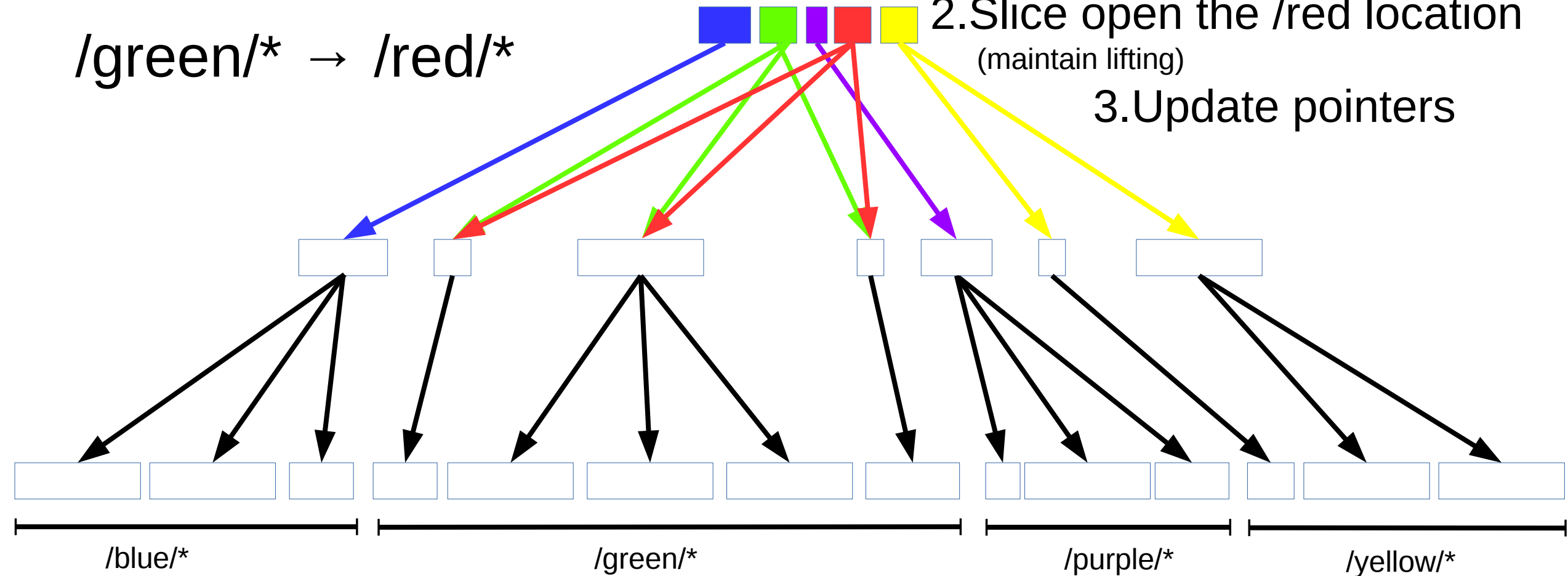
Range rename in lifted B^ϵ -trees

$/\text{green}/^* \rightarrow / \text{red}/^*$

1. Slice out the $/\text{green}/$ subtree
(maintain lifting)

2. Slice open the $/\text{red}/$ location
(maintain lifting)

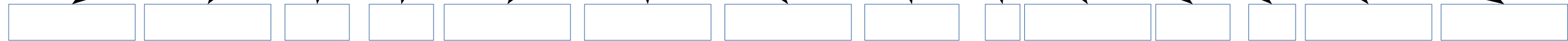
3. Update pointers



Range rename in lifted B^ϵ -trees

$/green/* \rightarrow /red/*$

1. Slice out the $/green$ subtree
(maintain lifting)
2. Slice open the $/red$ location
(maintain lifting)
3. Update pointers



$/blue/*$

$/red/*$

$/purple/*$

$/yellow/*$

IO complexity of range rename

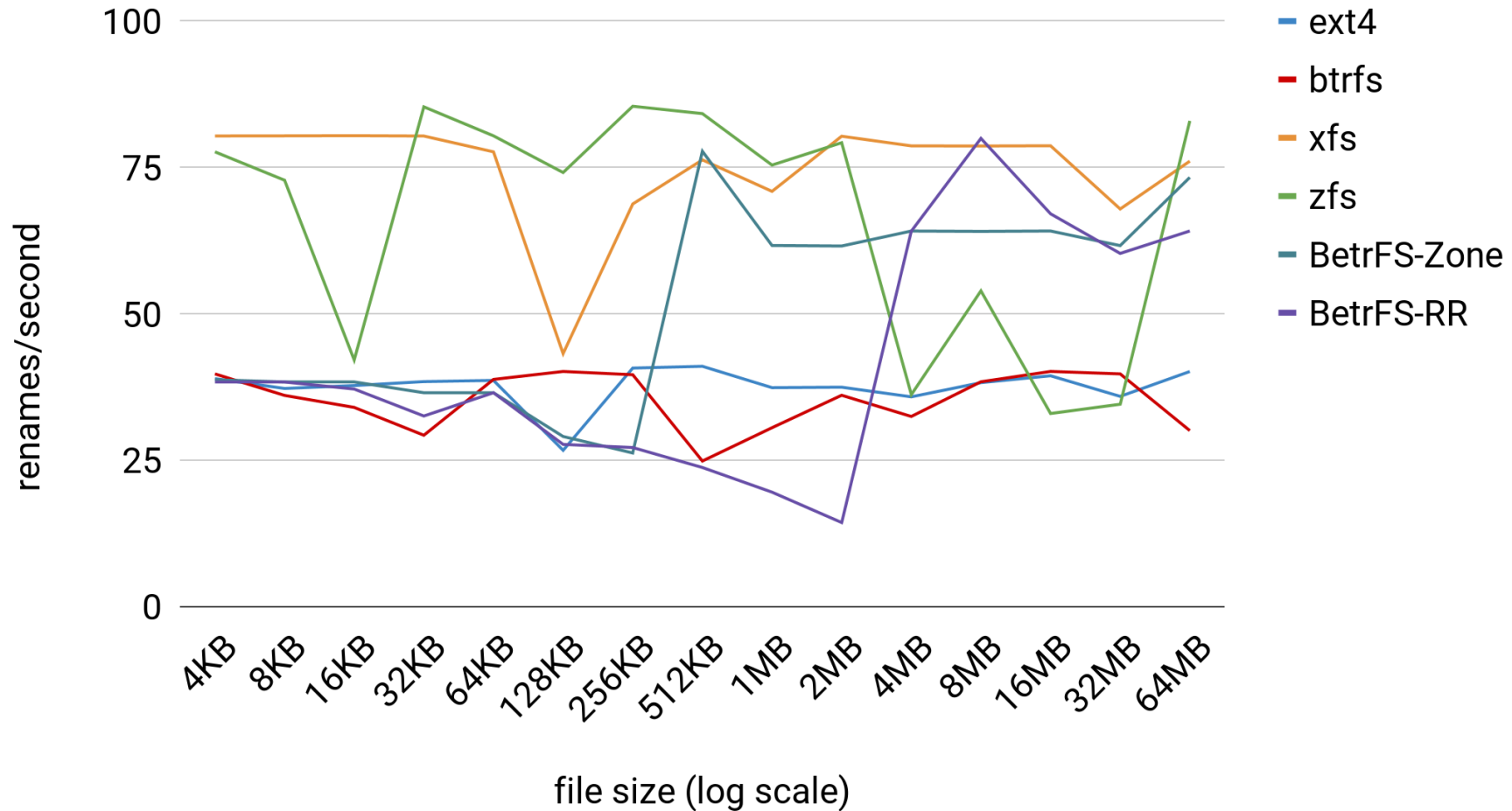
- Each slice touches a root-to-leaf path
- Maintaining lifting doesn't need to access any additional nodes

$$O(\text{tree height})$$

$$= O(\log_B N)$$

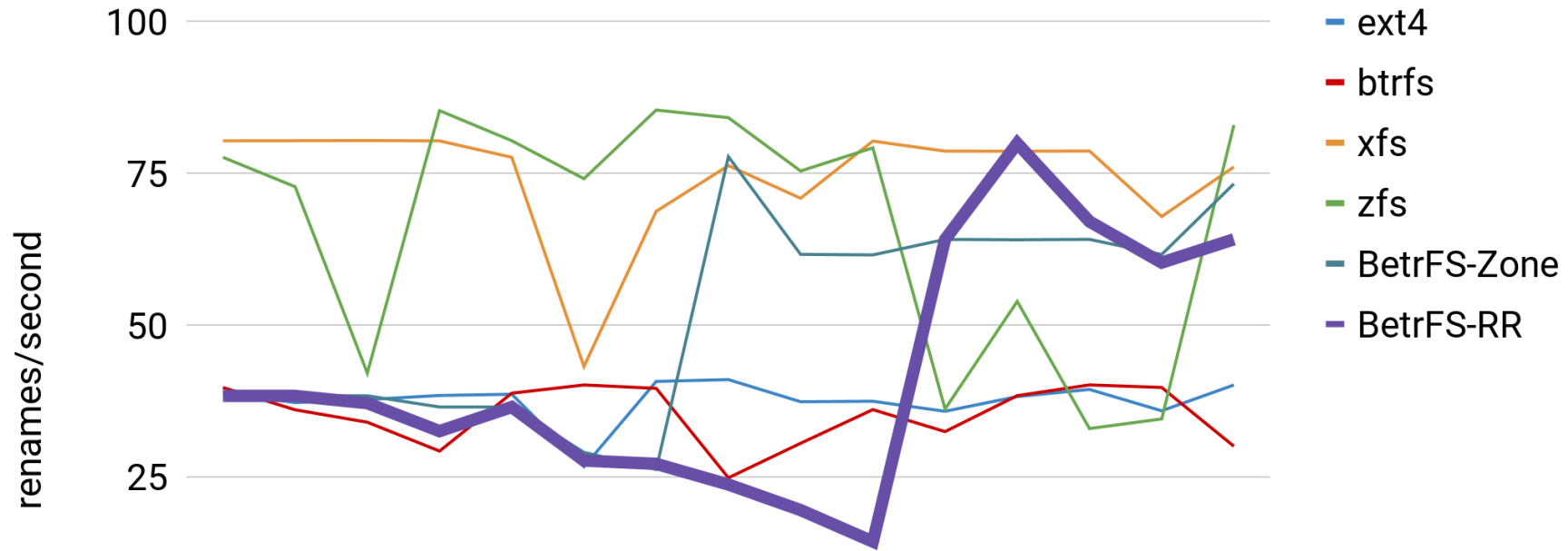
Rename Throughput

The average throughput of renaming one file 100 times (higher is better)



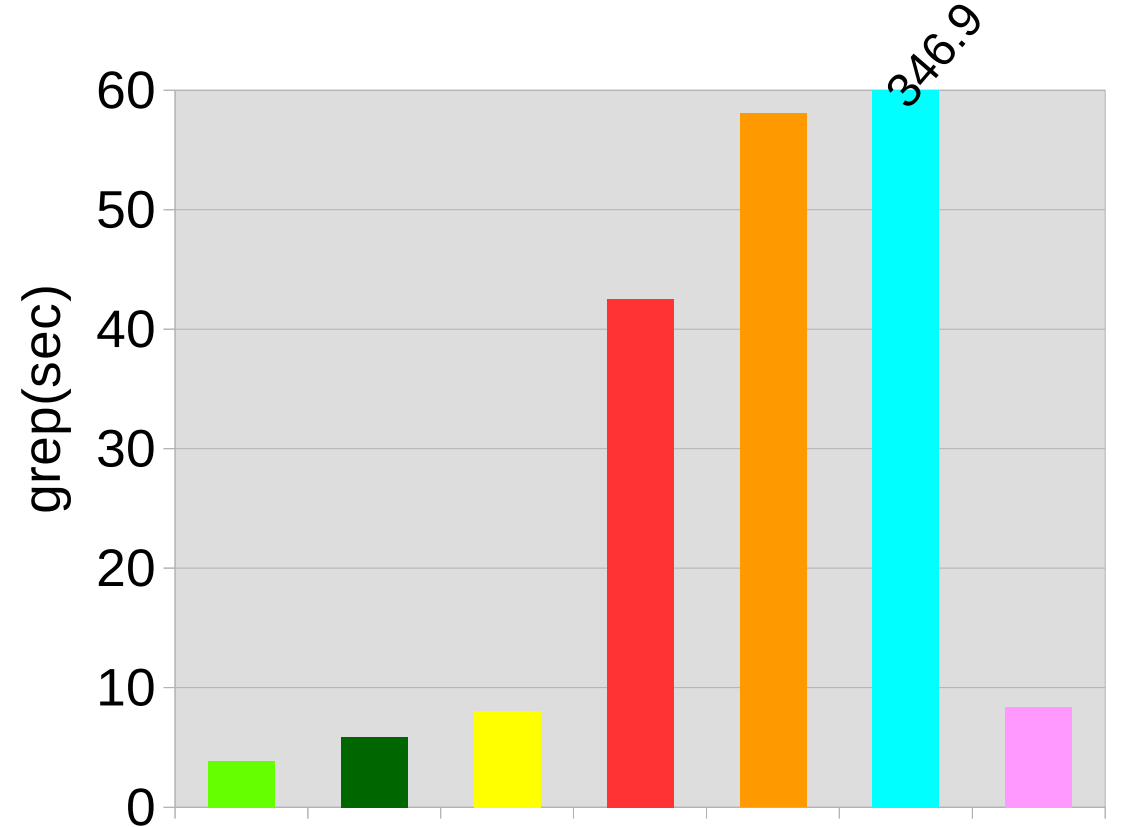
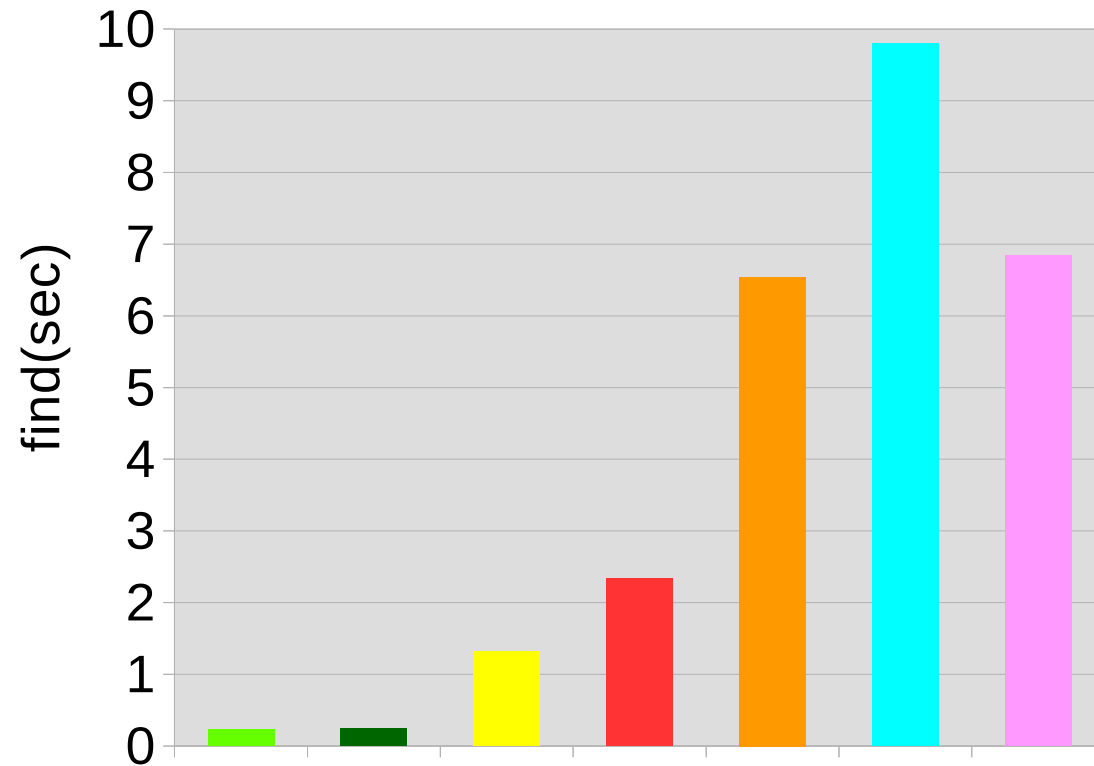
Rename Throughput

The average throughput of renaming one file 100 times (higher is better)



Rename can be as fast as inode-based file systems

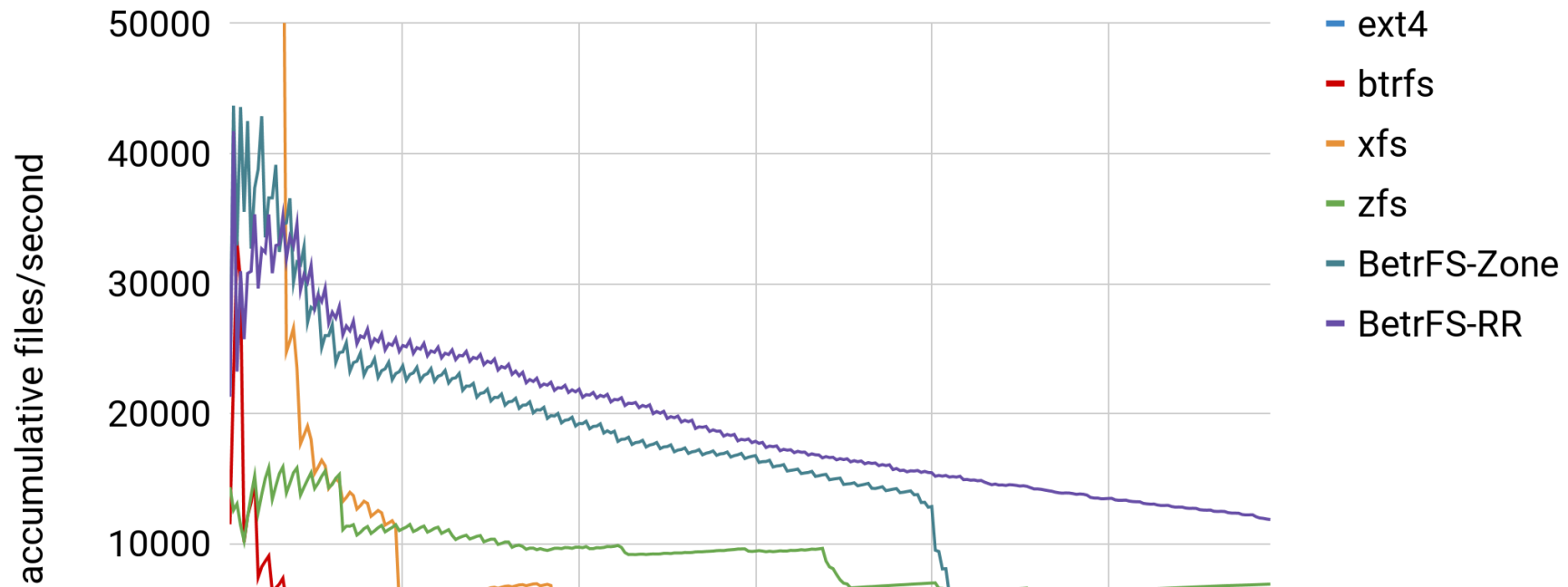
Recursive directory traversals



Full-path indexing can dramatically improve metadata scans

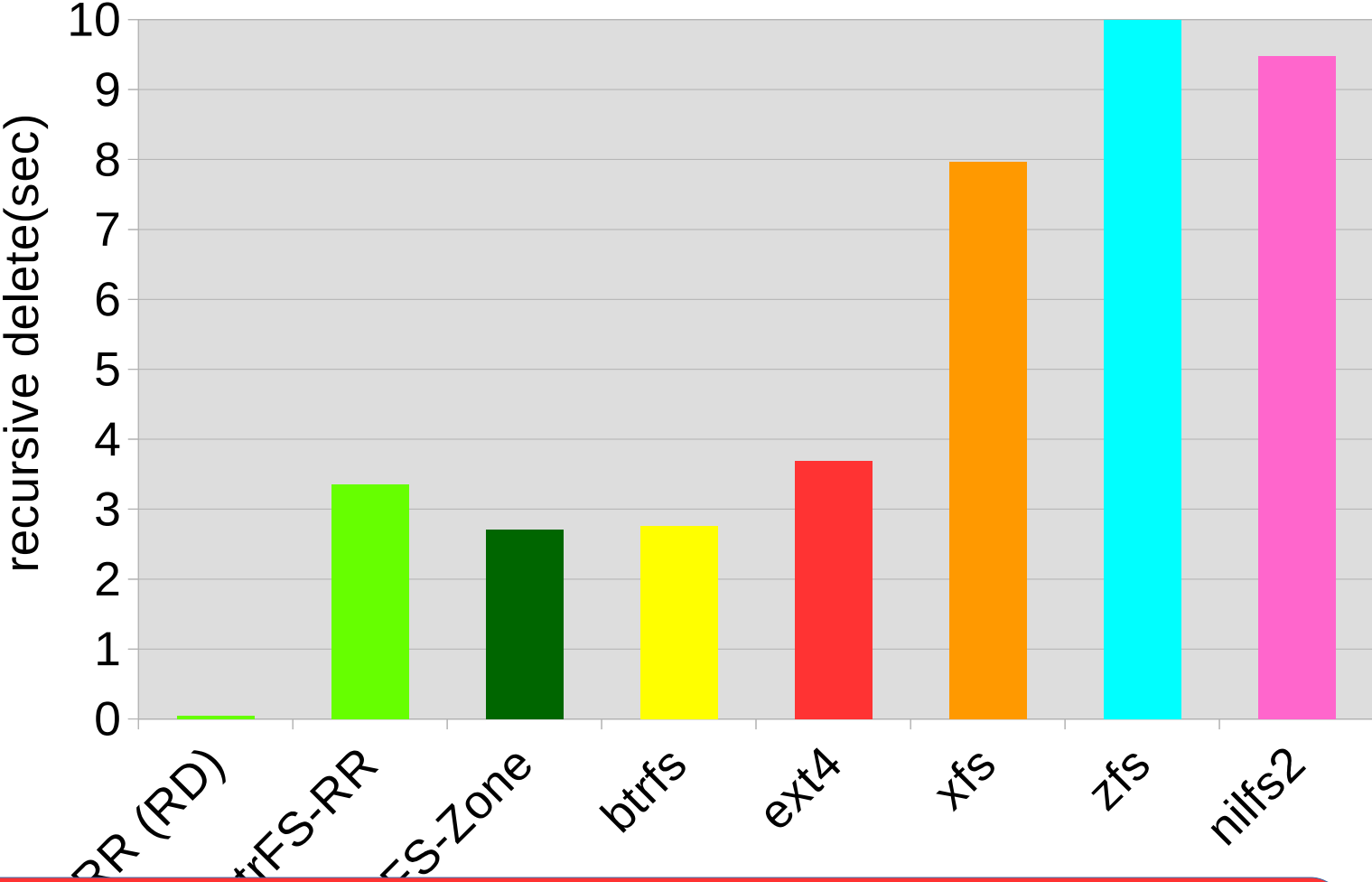
Tokubench

Tokubench: create 3 million 200-byte files in a balanced directory tree (higher is better)



Write-optimization can accelerate metadata updates

Recursive delete



Full-path indexing can unleash new optimizations

Conclusion

- Write optimization can solve a lot of metadata problems
- But write optimization can require rethinking how we organize our metadata
- And we sometimes need to extend the underlying data structures to support our metadata needs

**Code available at
betrfs.org**