# Scalable Range Locks for Scalable Address Spaces And Beyond

Alex Kogan        Dave Dice        **Shady Issa**

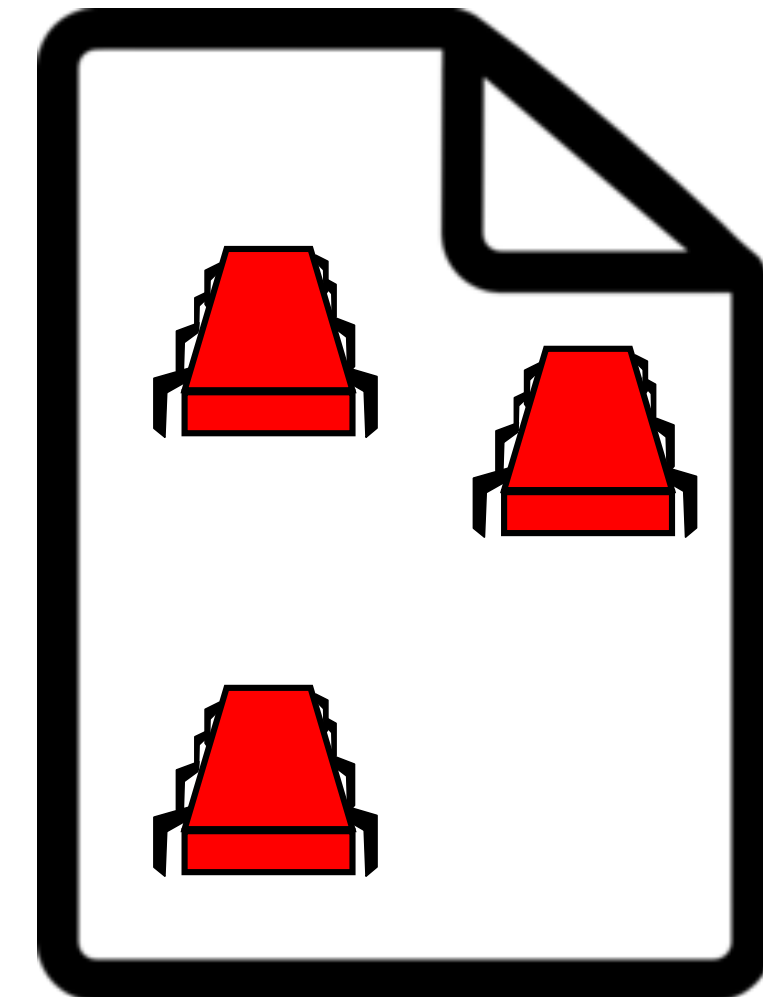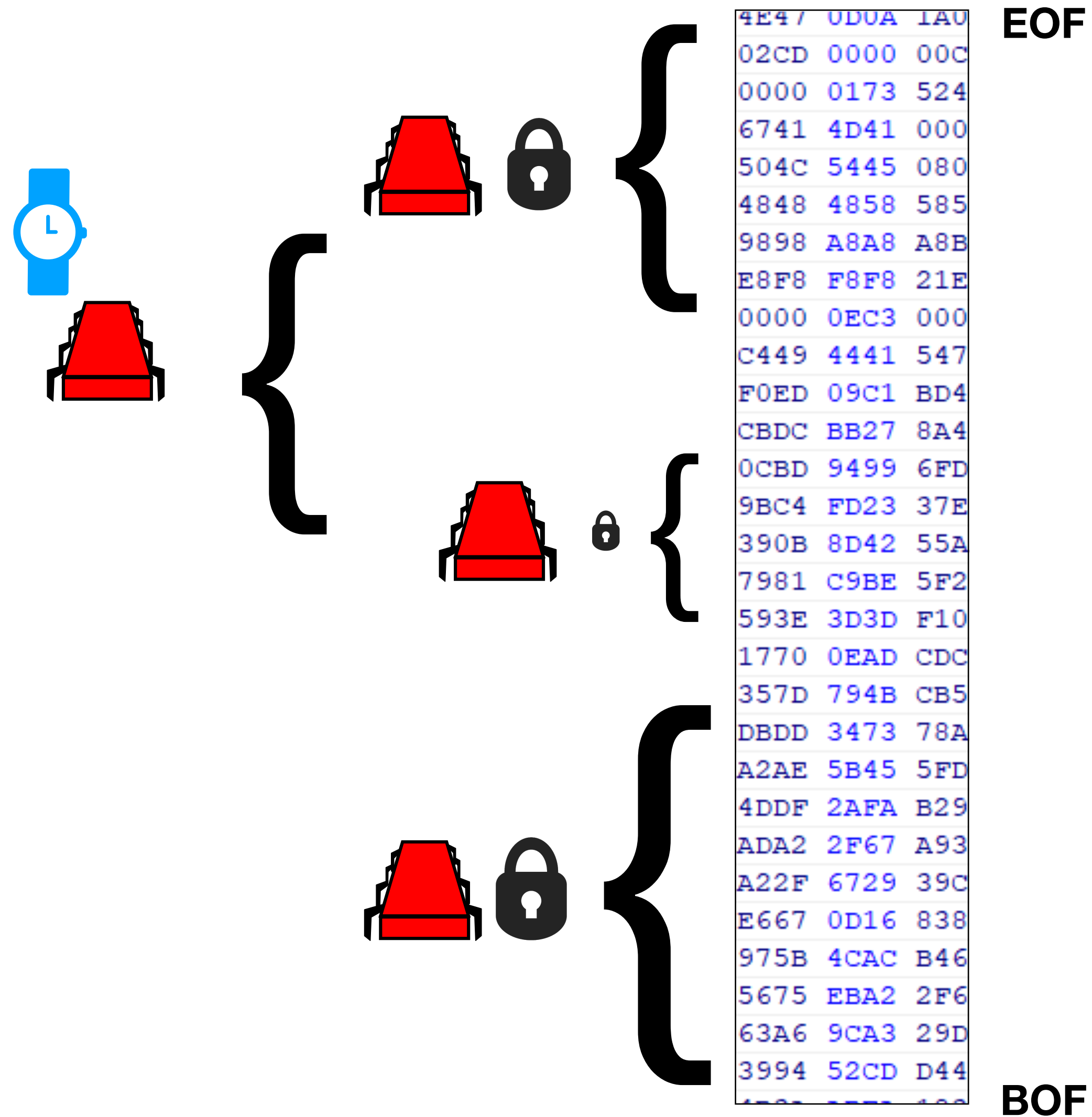Oracle Labs                    U. Lisboa & INESC-ID

# Range Locks

- Conceived in parallel filesystems

- Allow concurrent access to shared resources

  - e.g.: writing to the same file

# Range Locks

# Linux kernel Scalability Bottleneck

**How to get rid of mmap_sem**

**Please consider subscribing to LWN**

Subscriptions are the lifeblood of LWN.net. If you appreciate this conte[nt] and would like to see more of it, your subscription will help to ensure t[hat] LWN continues to thrive. Please visit this page to join up and keep LW[N on] the net.

By **Jonathan Corbet**
May 8, 2019

LSFMM

The mmap_sem lock used in the memory-management subsystem has been a known scalability problem for years, but it has proved difficult to rem[ove]. management track of the 2019 Linux [...] mmap_sem and how it might be elimina[ted...] vague at best.

## [RFC 0/4] Replace mmap_sem by a range lock

Laurent Dufour  |  Wed, 19 Apr 2017 05:19:10 -0700

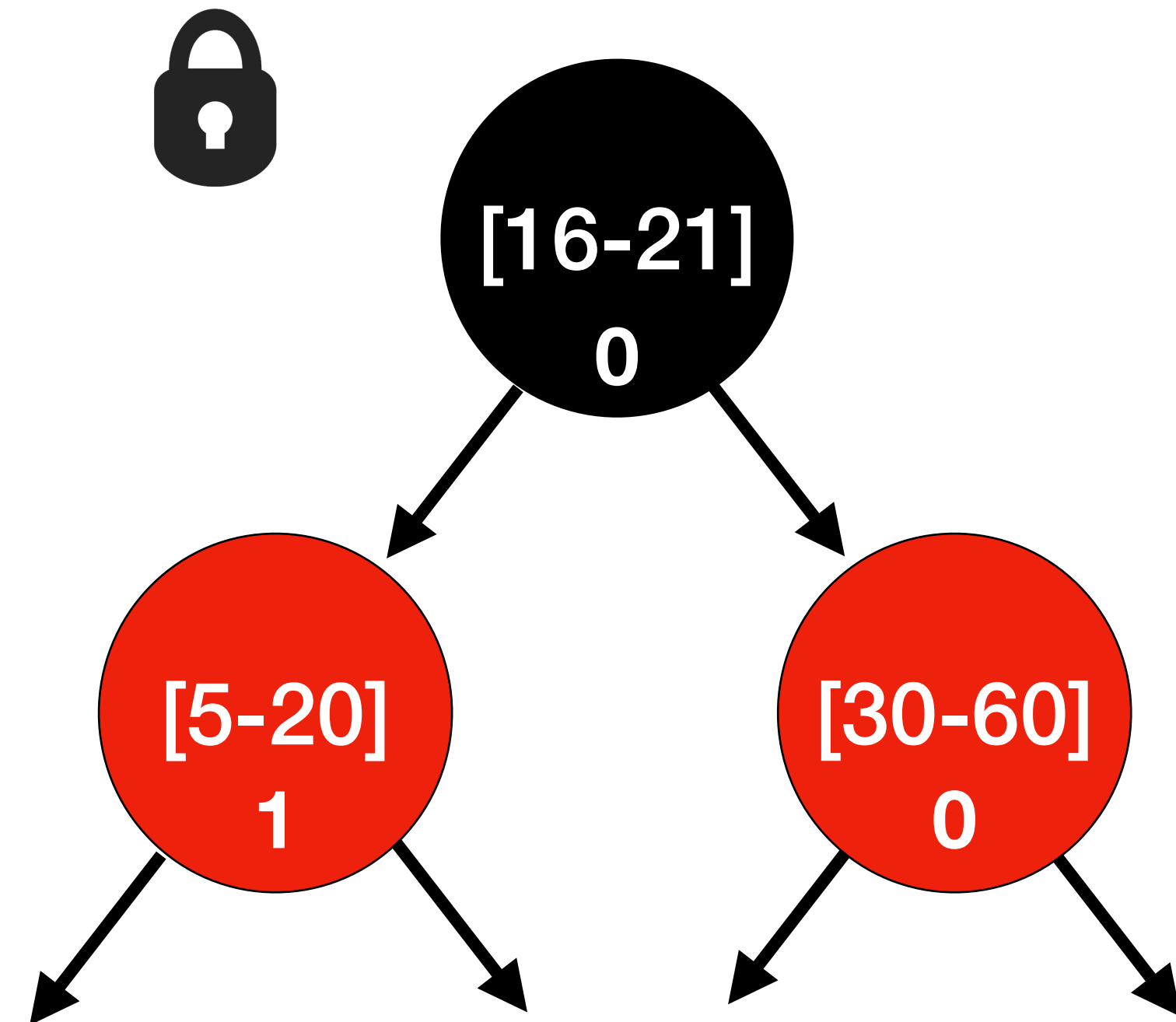## [RFC PATCH 00/64] mm: towards parallel address space operations

Davidlohr Bueso  |  Sun, 04 Feb 2018 17:30:07 -0800

## [PATCH v3 -tip 0/6] locking: Introduce range reader/writer lock

Davidlohr Bueso  |  Mon, 15 May 2017 02:08:54 -0700

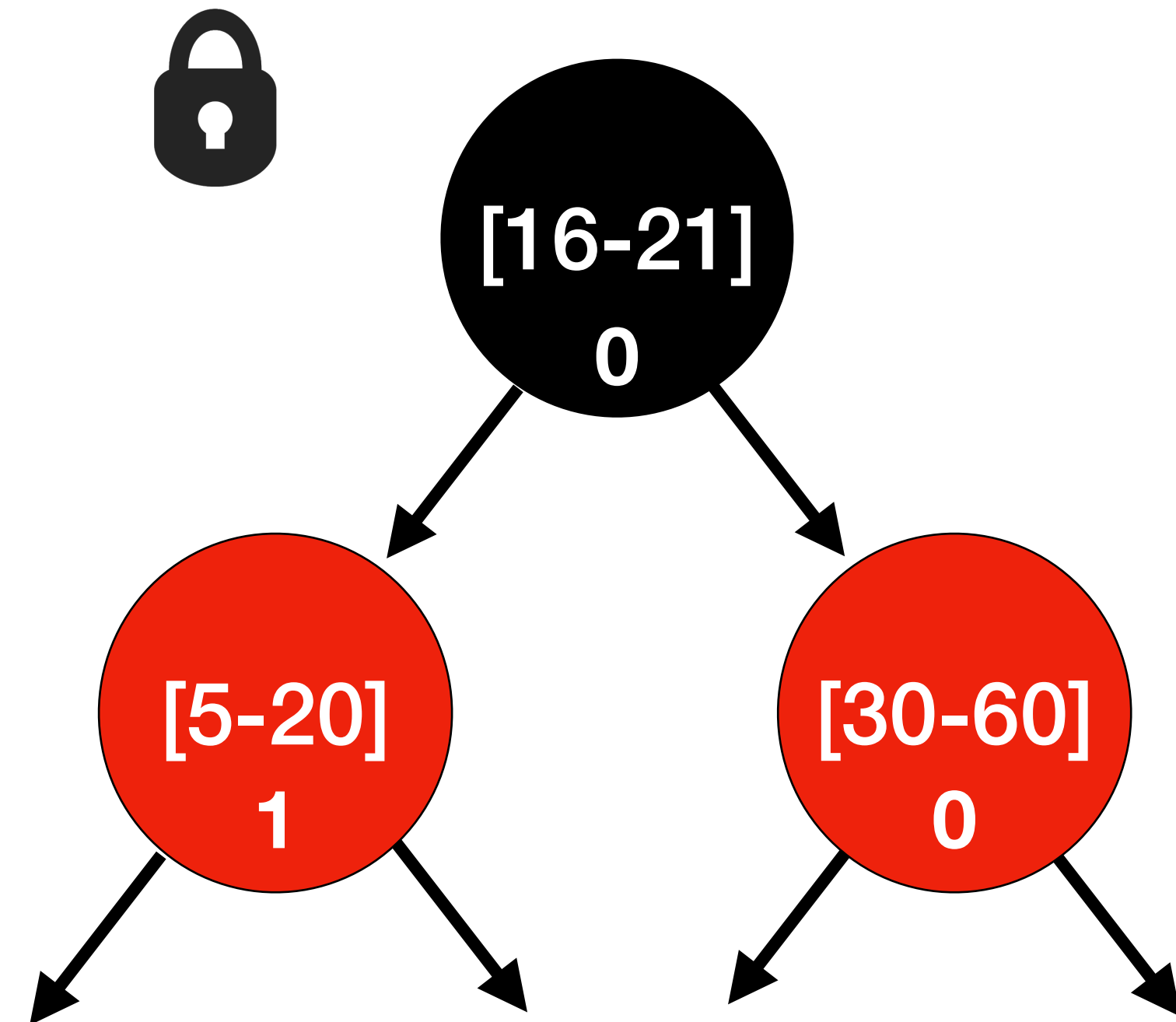4

# Existing Range Locks

- Auxiliary red-black tree

  - Ranges sorted by starting address

- Protected by spin-lock

  - contention even for shared access

# Existing Range Locks

Current RL are not scalable

- Aux... ...k tree

  - Ranges sorted by starting address

- Protected by spin-lock

  - contention even for shared access

[16-21]
0

[5-20]
1

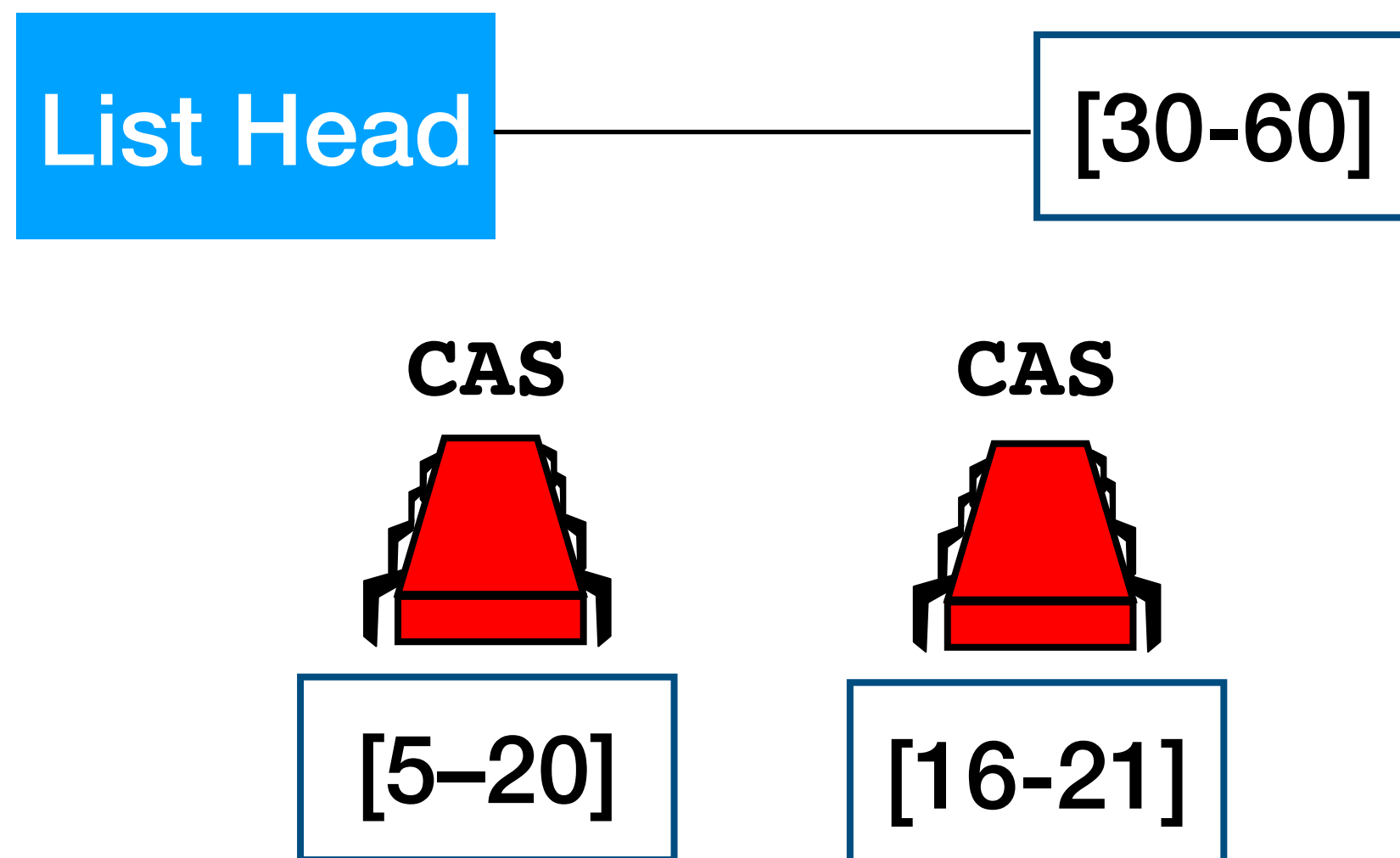[30-60]
0

# Our Contributions

- New design for Range locks

  - Lock-free in the common case

  - Scales up to **144 threads**

- Speculative approach for VM operations in the Linux kernel

- Range locks for skip lists

# List-based Range Locks

- A range lock is acquired once a range is inserted into a list

  - Sorted by their starting addresses

# List-based Range Locks

- A range lock is acquired once a range is inserted into a list

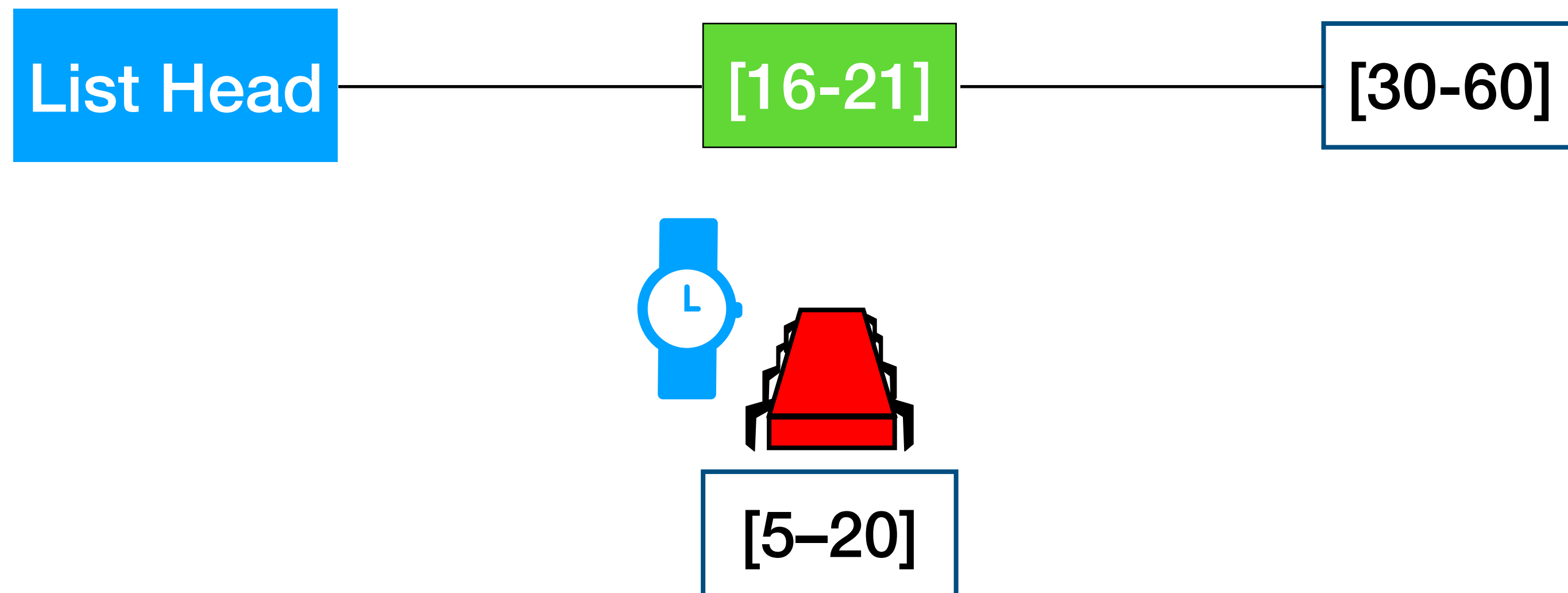  - Sorted by their starting addresses

# List-based Range Locks

- A range lock is acquired once a range is inserted into a list

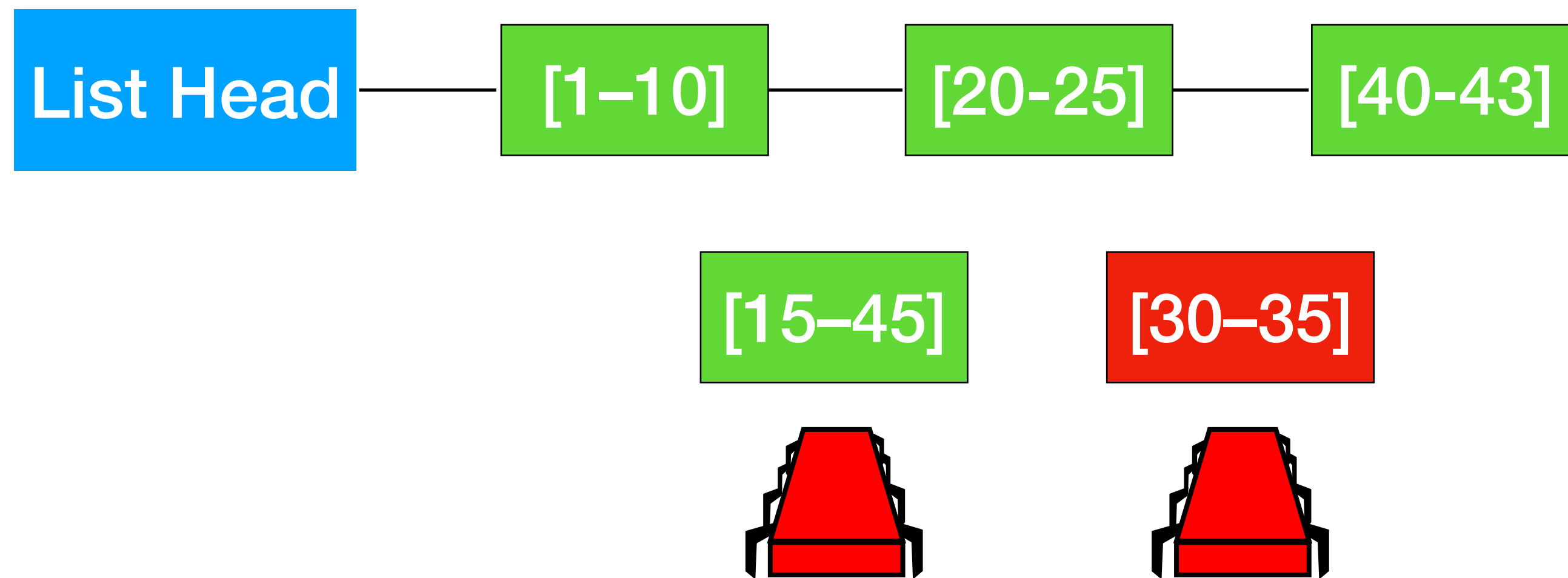  - Sorted by their starting addresses

# List-based Range Locks

- A range lock is acquired once a range is inserted into a list

  - Sorted by their starting addresses

# List-based Range Locks

- A range lock is acquired once a range is inserted into a list

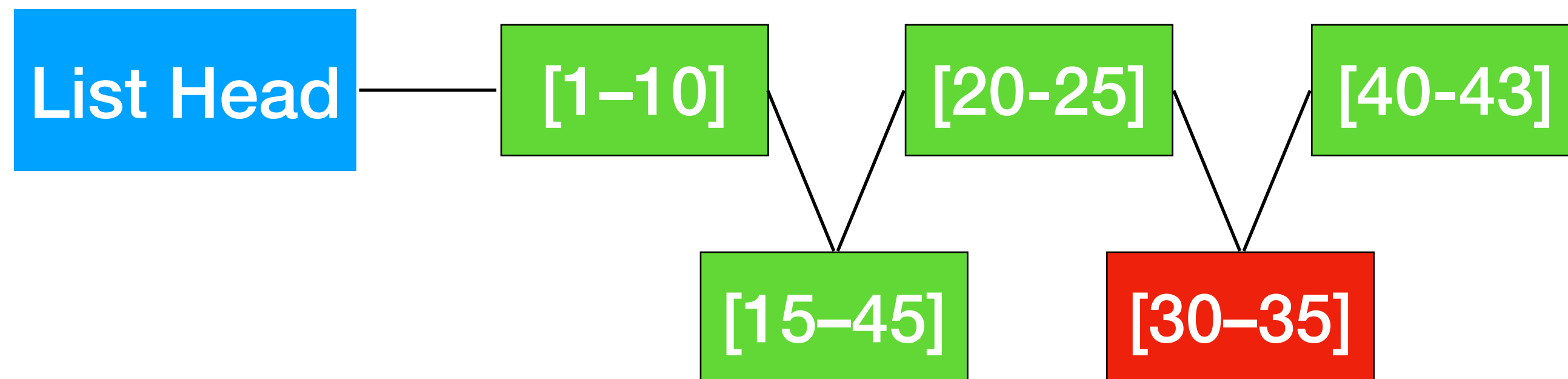  - Sorted by their starting addresses

We only need an extra validation step for Read-Write semantics

List Head — [1–10]    [20-25]    [40-43]

[15–45]    [30–35]

# VM Management in the Kernel

`0xffffffffff`

Virtual
Address
Space

`0x00000000`

# VM Management in the Kernel

**0xffffffff**

**VMA 1**

*start:*
*length:*
*Access rights:*
**READ|WRITE**

**VMA 2**

*start:*
*length:*
*Access rights:*
**READ|WRITE**

**VMA 3**

*start:*
*length:*
*Access rights:*
**NONE**

**Virtual Address Space**

**0x00000000**

14

# VM Management in the Kernel

**VMA 1**
*start:*
*length:*
*Access rights:*
**READ|WRITE**

**VMA 2**
*start:*
*length:*
*Access rights:*
**READ|WRITE**

**VMA 3**
*start:*
*length:*
*Access rights:*
**NONE**

`0xffffffff`

**Virtual Address Space**

`0x00000000`

`mm_rb`

VMA 1

VMA 3

VMA 2

15

# VM Management in the Kernel

```
mprotect(1000, 100, READ|WRITE)
```

```
mprotect(3000, 100, NONE)
```

**VMA**
*start:* 1000
*length:* 5000
*Access rights:*
**READ|WRITE**

0xffffffffff

Virtual
Address
Space

0x00000000

# VM Management in the Kernel

`0xffffffff`

`mprotect(1000, 100, READ|WRITE)`

**VMA**
*start:* 1000
*length:* 5000
*Access rights:*
**READ|WRITE**

`mprotect(3000, 100, NONE)`

Protecting ranges naively can create data races

**Virtual Address Space**

`0x00000000`

# Refined Ranges for VM

```
VM_Operation(start, length, args..){
  Acquire_mm_sem();
  VMA = find_vma(start);
  // operation logic
  …
  read_only operations
  Decide if structural modification is required
  …
  Release_mm_sem();
}
```

# Refined Ranges for VM

```
VM_Operation(start, length, args..){
  Acquire_mm_sem();

  VMA = find_vma(start);

  // operation logic

  …

  read_only operations

  Decide if structural modification is required

  …

  Release_mm_sem();
}
```

# Refined Ranges for VM

Traverses the red-black tree `mm_rb`

```
VM_Operation(start, length, args..){
    Acquire_mm_sem();
    VMA = find_vma(start);
    // operation logic
    …
    read_only operations
    Decide if structural modification is required
    …
    Release_mm_sem();
}
```

# Refined Ranges for VM

Protect with range lock of input range

```
VM_Operation(start, length, args..){
    Acquire_mm_sem();
    Acquire_RL_Read(start, start+length);
    VMA = find_vma(start);
    Release_RL();
    // operation logic
    …
    read_only operations
    Decide if structural modification is required
    …
    Release_mm_sem();
}
```

# Refined Ranges for VM

```
VM_Operation(start, length, args..){
    Acquire_mm_sem();
    Acquire_RL_Read(start, start+length);
    VMA = find_vma(start);
    Release_RL();
    Acquire_RL_Write(VMA.start-x, VMA.end+x);
    // operation logic
    …
    read_only operations
    Decide if structural modification is required
    …
    Release_RL();
    Release_mm_sem();
}
```

**Protect with range lock of VMA range+Δ**

**check if the `mm_rb` changed meanwhile**
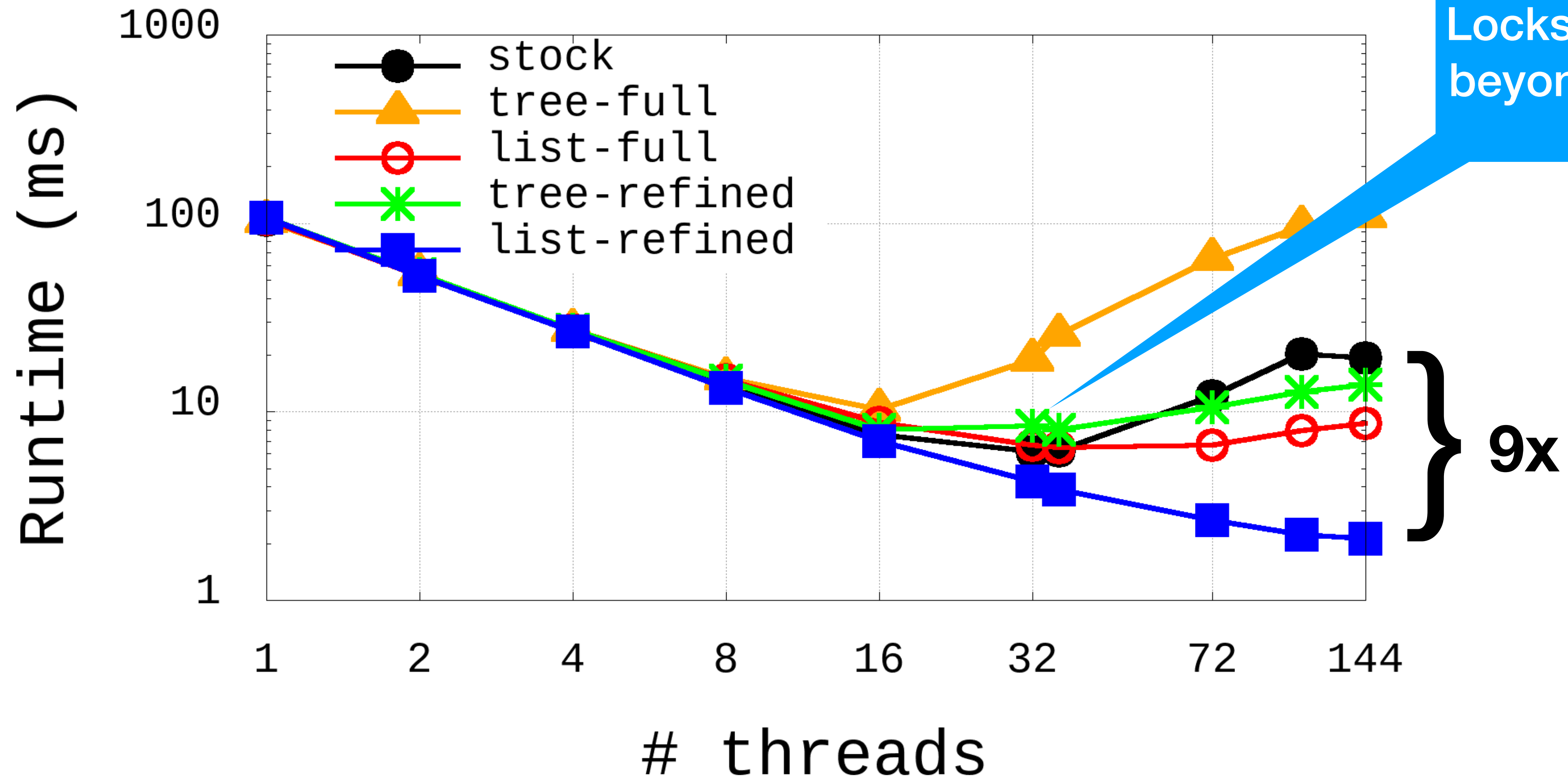
# Refined Ranges for VM

```
VM_Operation(start, length, args..){
    Acquire_mm_sem();

    Acquire_RL_Read(start, start+length);

    VMA = find_vma(start);

    Release_RL();

    Acquire_RL_Write(VMA.start-x, VMA.end)+x;

    // operation logic

    …

    read_only operations

    if structural modification is required{
        Release_RL();

        Acquire_RL_Write(0,2^63-1);

        retry();
    }

    …

    Release_RL();

    Release_mm_sem();

}
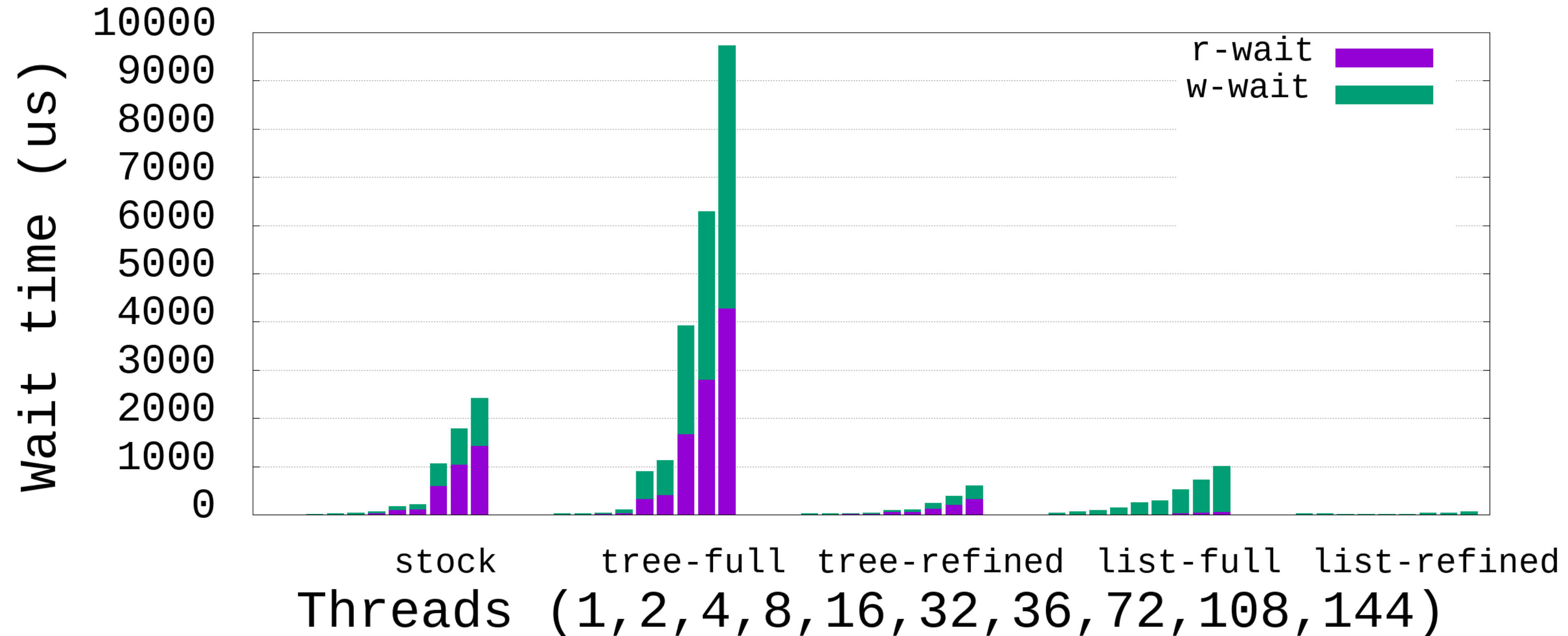```

Acquire full range lock and retry

# Evaluation

- Linux kernel 4.16.0-rc2

- 4 Intel Xeon E7-8895 v3 (144 threads)

- Metis benchmark (wrmem)

- Baselines:

  - Stock

  - Tree-based RL (with and w/out speculation)

  - List-based RL (with and w/out speculation)

# Evaluation

# Evaluation



Wait time (us)

Threads (1,2,4,8,16,32,36,72,108,144)

stock    tree-full    tree-refined    list-full    list-refined

r-wait
w-wait

Collected using
`lock_stats`

26

# More in the paper…

- Evaluation:

  - More workloads

  - User-space applications

- Range Locks design

  - Fast path, avoiding starvation, memory reclamation

  - Range locks for skip lists

# Conclusion

- Scalable linked list-based Range Locks

- New speculative approach for the Linux VM

- Using Range Locks for concurrent data structures