

# Arun Srinivasan

 @arun\_sriniv

Phd student, MPIPZ, Germany

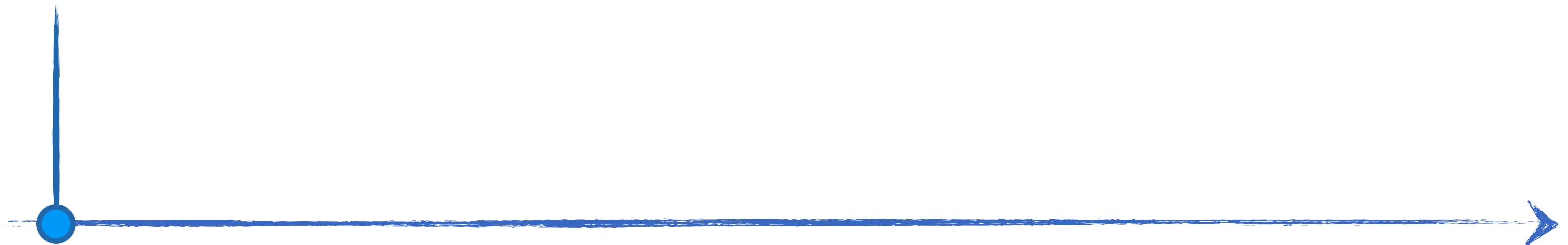
co-developer, data.table

# TIMELINE

melt  
dcast



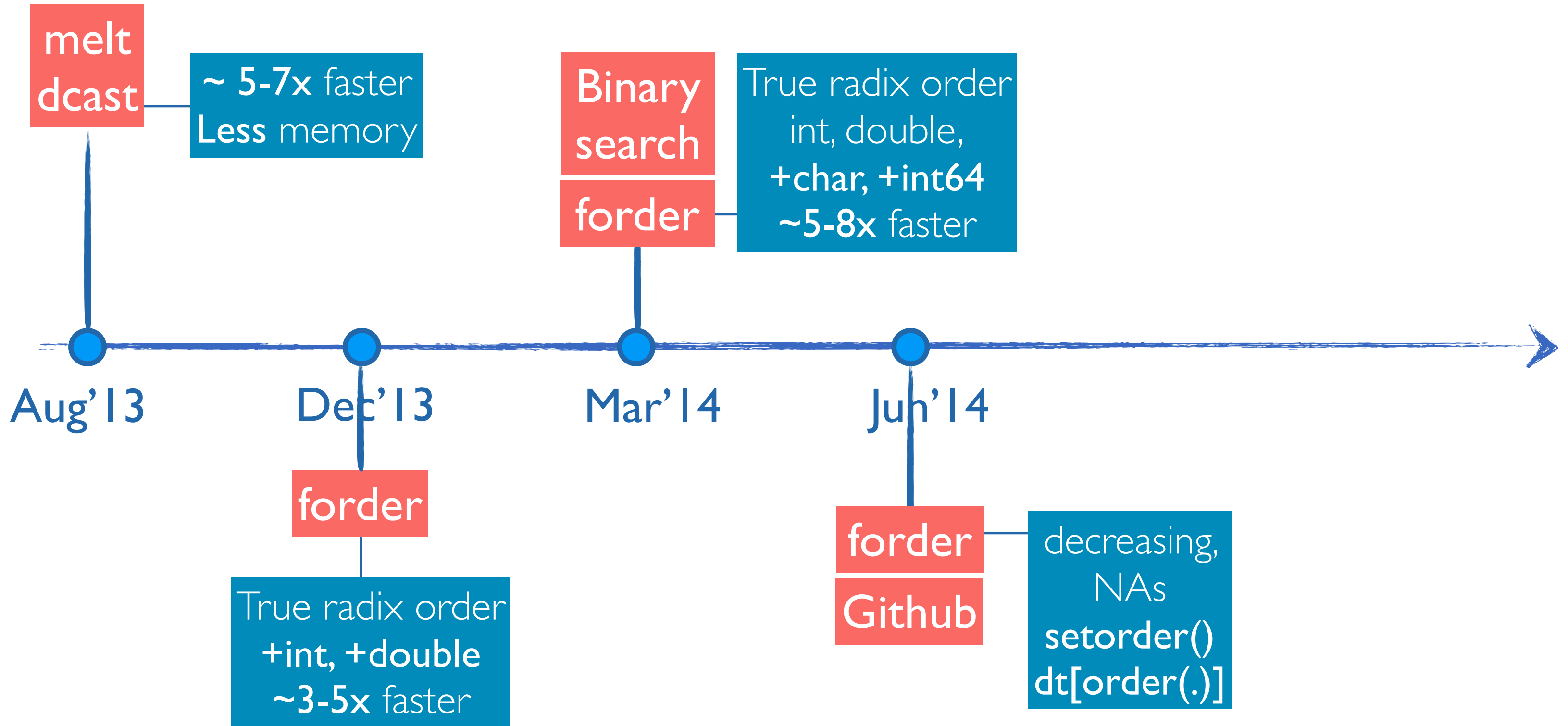
Aug' 13



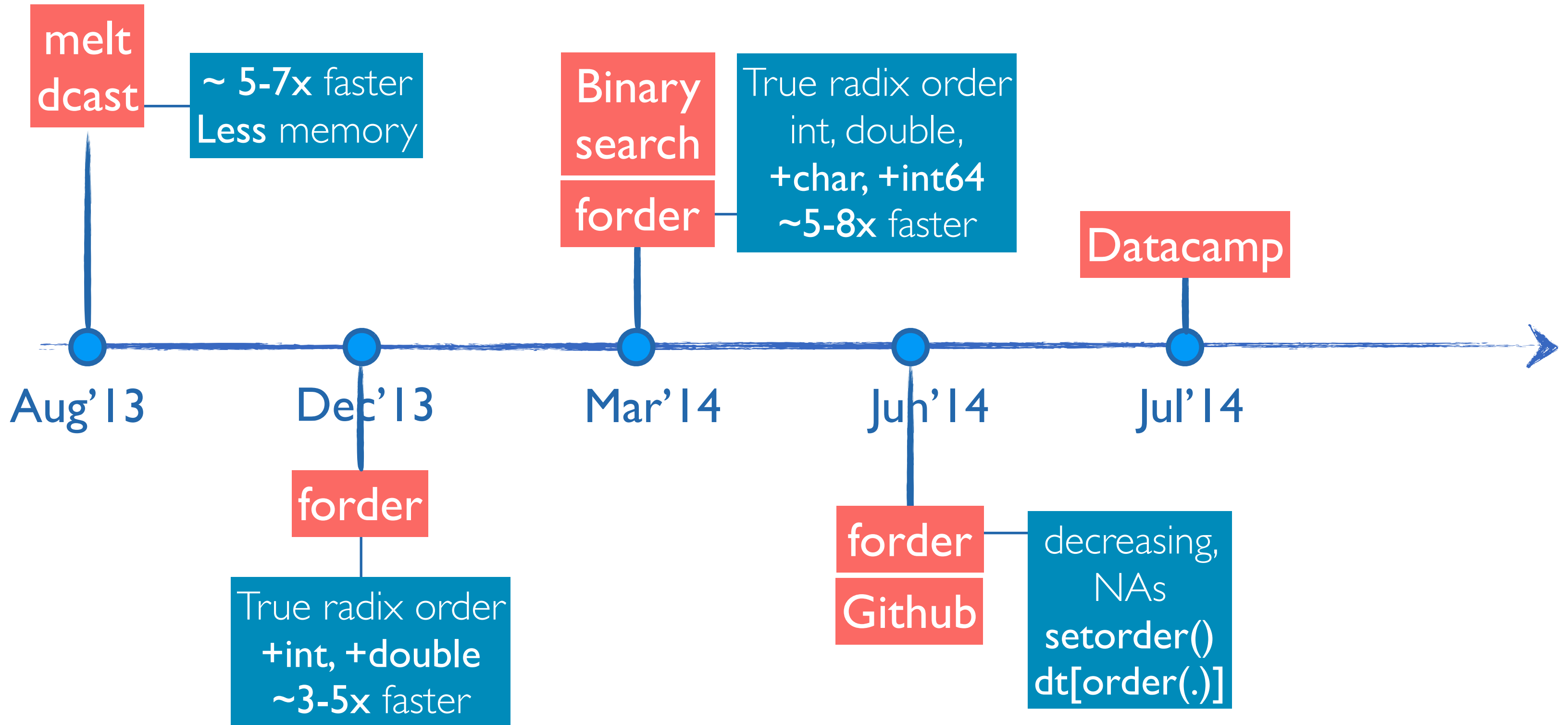
# TIMELINE



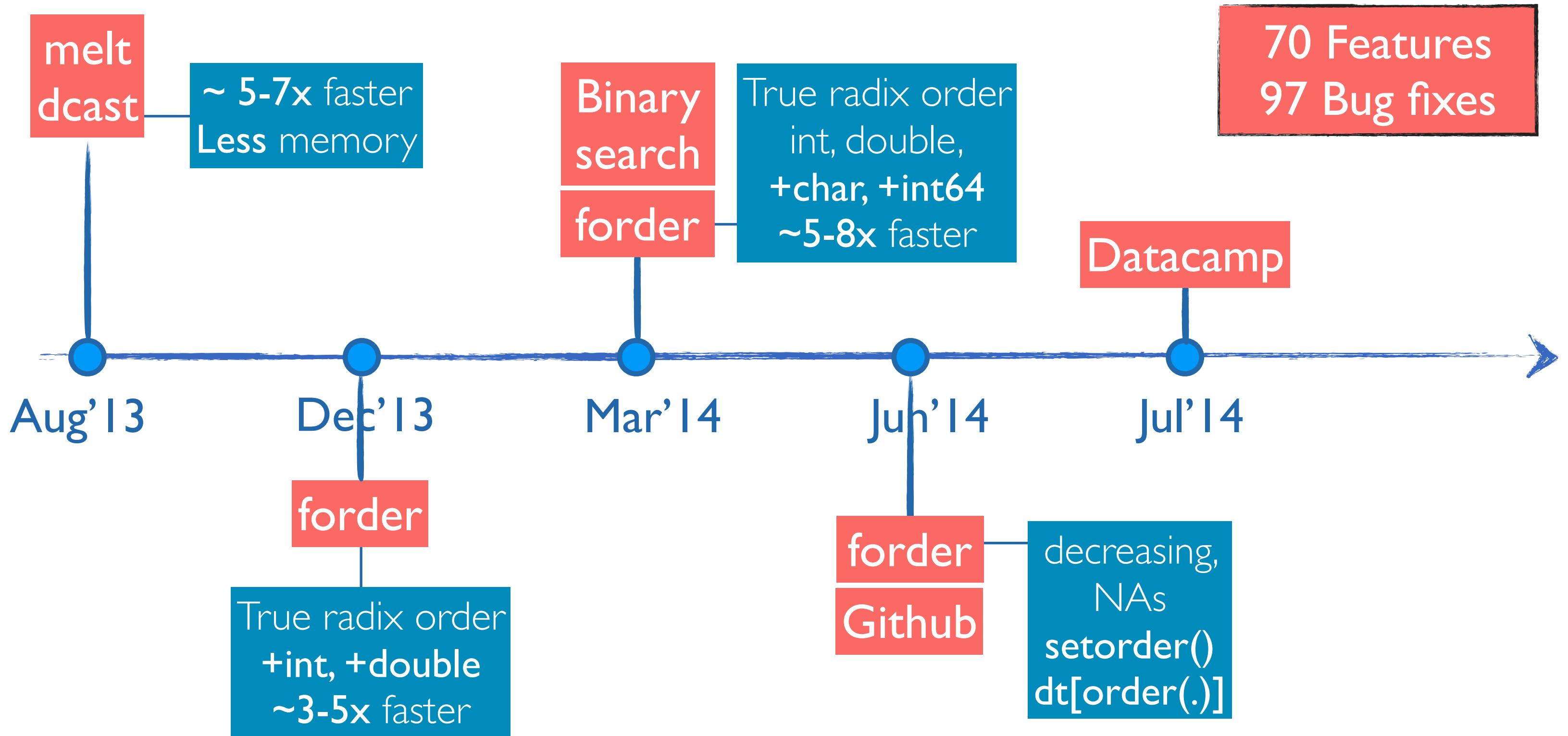
# TIMELINE



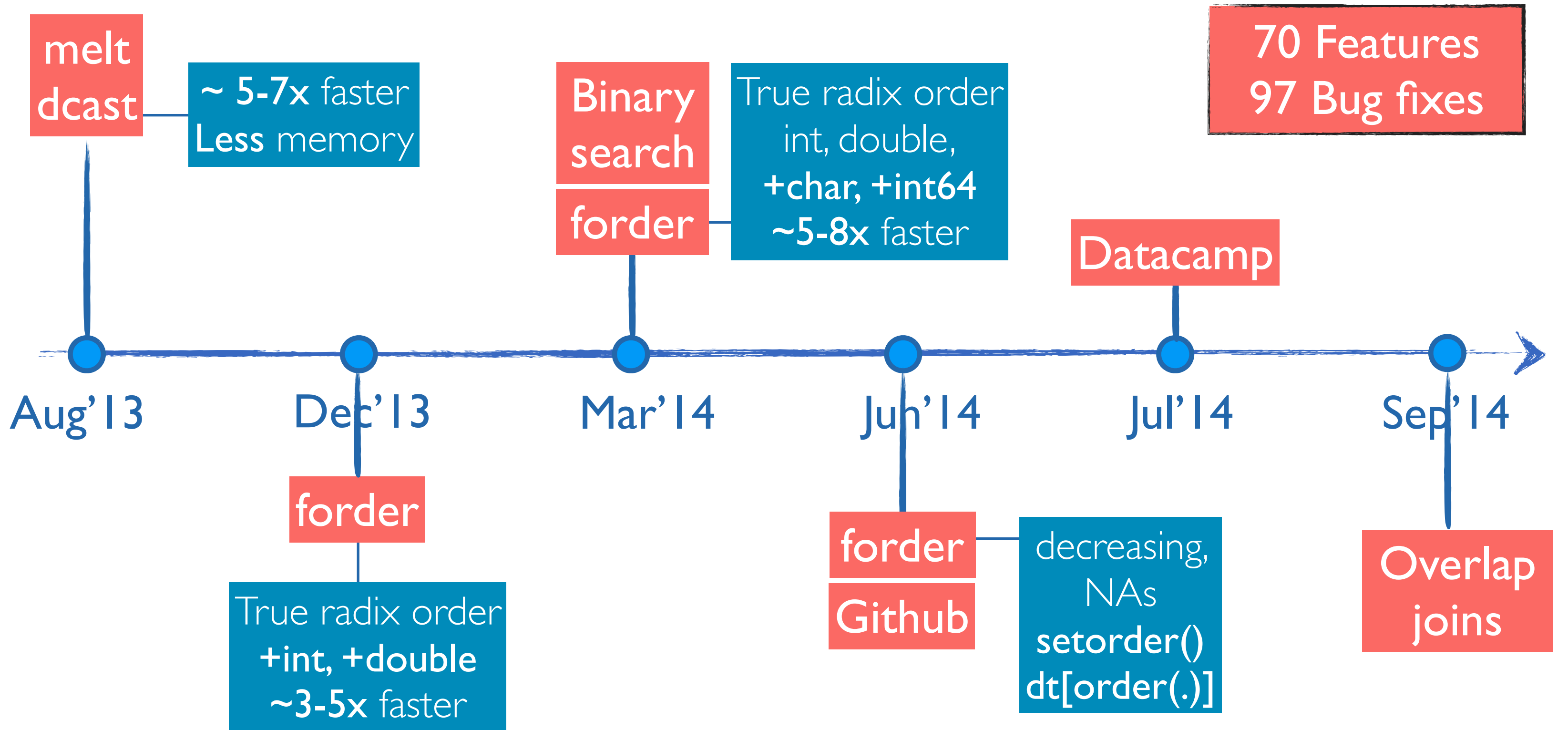
# TIMELINE



# TIMELINE



# TIMELINE



# OVERLAP JOINS



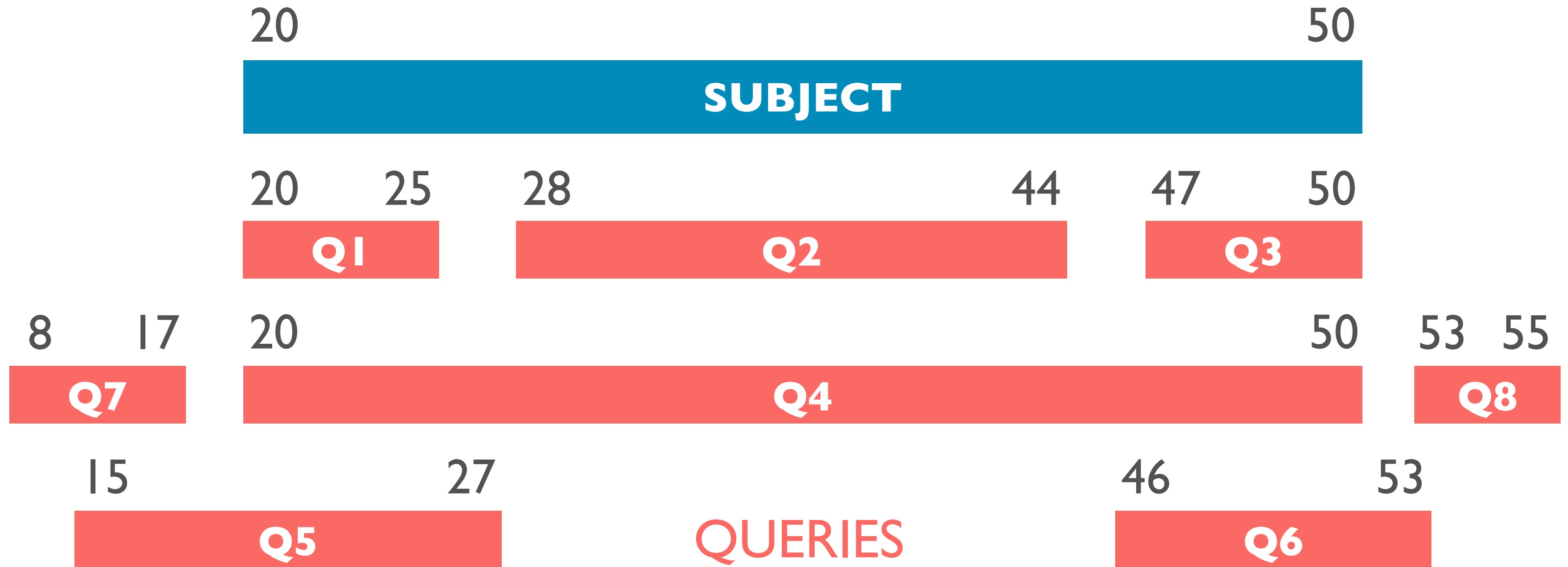
# WHAT ARE OVERLAPS?

20

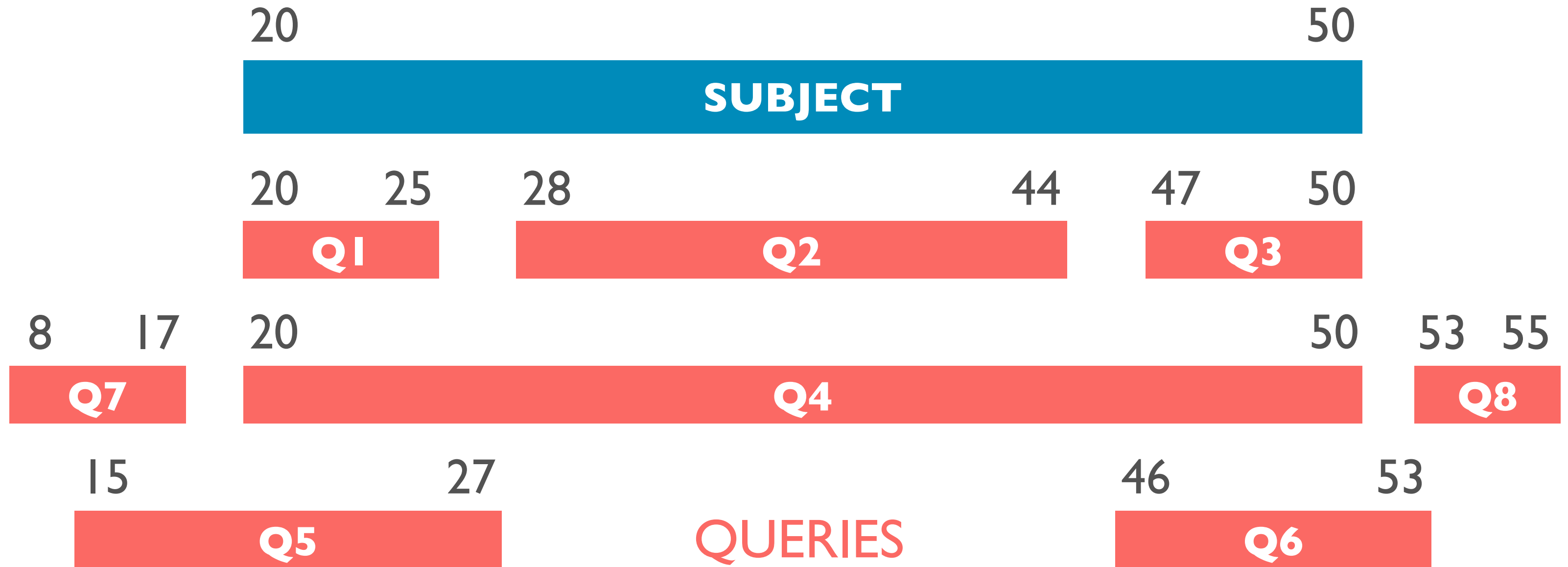
50

**SUBJECT**

# WHAT ARE OVERLAPS?

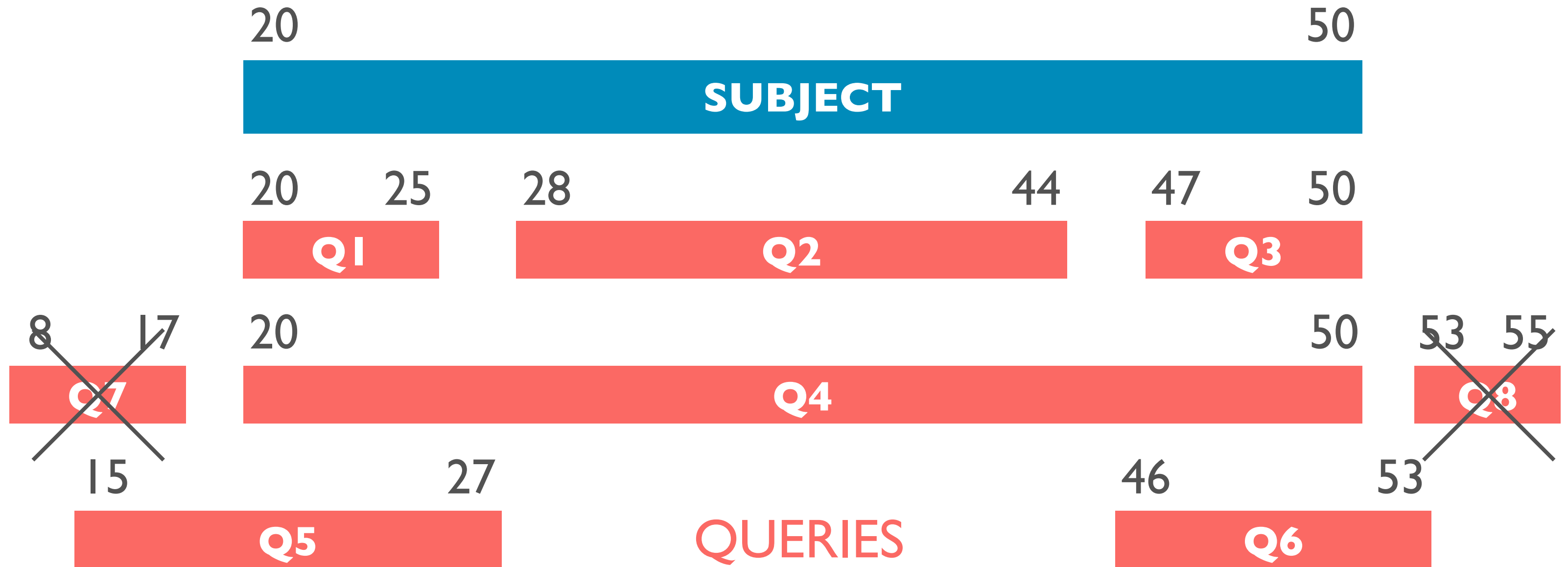


# WHAT ARE OVERLAPS?



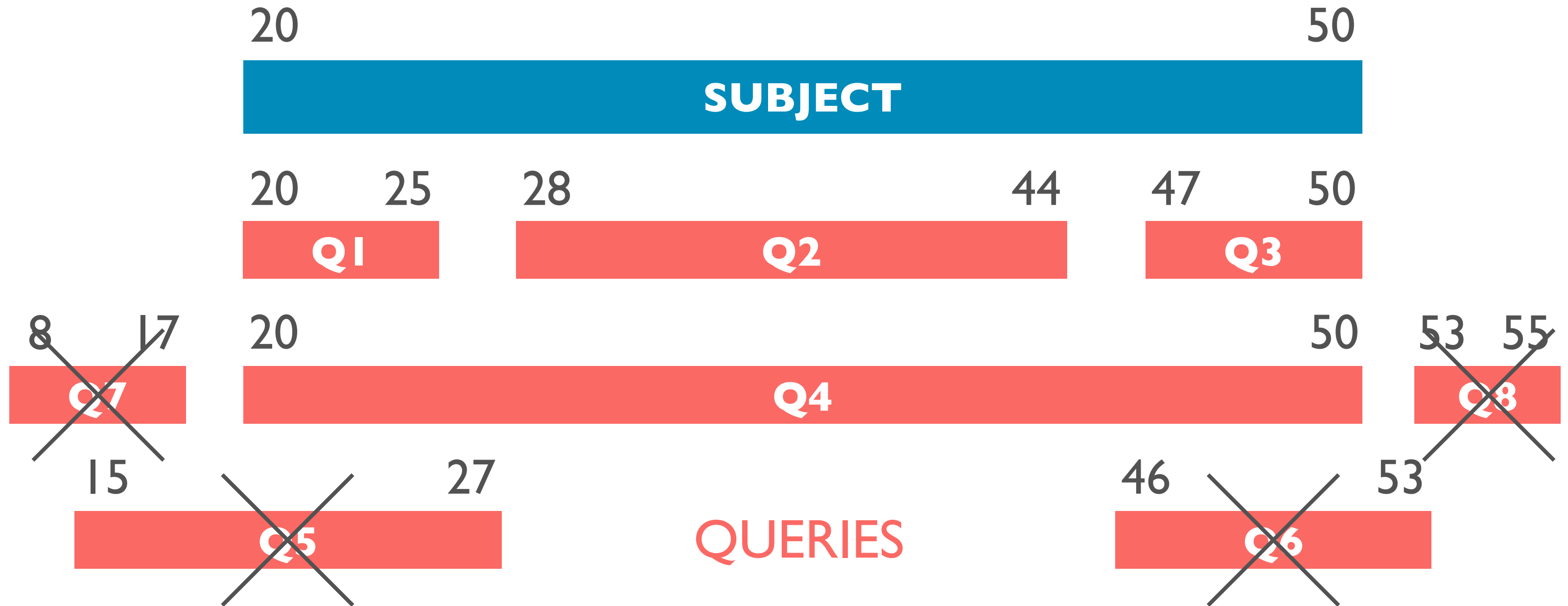
Which of the **red** ranges fall completely within the **blue** range?

# WHAT ARE OVERLAPS?



Which of the **red** ranges fall completely within the **blue** range?

# WHAT ARE OVERLAPS?



Which of the **red** ranges fall completely within the **blue** range?

# HOW IT WORKS: STEP 1

**SUBJECT** is *usually* much smaller than **QUERY**

	start	end
1:	10	16
2:	20	35
3:	30	45

**SUBJECT**

2D-form

# HOW IT WORKS: STEP 1

**SUBJECT** is *usually* much smaller than **QUERY**

	start	end
1:	10	16
2:	20	35
3:	30	45

**SUBJECT**

2D-form



1. Sort **SUBJECT** by **start**, **end**
2. Add one to **end**
3. collapse and sort *again*

# HOW IT WORKS: STEP 1

**SUBJECT** is *usually* much smaller than **QUERY**

	start	end
1:	10	16
2:	20	35
3:	30	45

**SUBJECT**

2D-form



1. Sort **SUBJECT** by **start, end**
2. Add one to **end**
3. collapse and sort *again*

pos
10
17
20
30
36
46



# HOW IT WORKS: STEP 1

**SUBJECT** is *usually* much smaller than **QUERY**

	start	end
1:	10	16
2:	20	35
3:	30	45

**SUBJECT**  
2D-form



1. Sort **SUBJECT** by **start, end**
2. Add one to **end**
3. collapse and sort *again*
4. Get **row numbers**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

**ID** -form

# HOW IT WORKS: STEP 1

**SUBJECT** is *usually* much smaller than **QUERY**

	start	end
1:	10	16
2:	20	35
3:	30	45

**SUBJECT**

2D-form



1. Sort **SUBJECT** by **start, end**
2. Add one to **end**
3. collapse and sort *again*
4. Get **row numbers**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

ID -form

Thanks to **Matt** for describing the technique based on conversation with **@corone**

# HOW IT WORKS: STEP 2

Use *rolling joins* on **QUERY**'s **start** and **end** separately

start	end
12	15
41	50
7	9
33	34

**QUERY**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

**ID -form**

→  
roll=TRUE  
LOCF

# HOW IT WORKS: STEP 2

Use *rolling joins* on **QUERY**'s start and end separately

start	end
12	15
41	50
7	9
33	34

**QUERY**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

**ID -form**

→  
roll=TRUE  
LOCF

# HOW IT WORKS: STEP 2

Use *rolling joins* on **QUERY**'s start and end separately

start	end
12	15
41	50
7	9
33	34

**QUERY**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

**ID -form**

→  
roll=TRUE  
LOCF

# HOW IT WORKS: STEP 2

Use *rolling joins* on **QUERY**'s start and end separately

start	end
12	15
41	50
7	9
33	34

**QUERY**

pos	row number
10	1
17	-
20	2
30	2,3
36	3
46	-

**ID -form**

roll=TRUE  
LOCF

start	end
10	16
NA	NA
NA	NA
20	35
30	45

**RESULT**

# CODE

```
foverlaps(QUERY,  
          SUBJECT,  
          type = "within")
```

# WHAT ELSE?



# WHAT ELSE?

- Intervals can be **integer**, **numeric**, **integer64**, **Date**, **POSIXct** etc

# WHAT ELSE?

- Intervals can be **integer**, **numeric**, **integer64**, **Date**, **POSIXct** etc
- **QUERY** need not be sorted; its order is preserved in **RESULT**

# WHAT ELSE?

- Intervals can be **integer**, **numeric**, **integer64**, **Date**, **POSIXct** etc
- **QUERY** need not be sorted; its order is preserved in **RESULT**
- Oh... and it's **FAST!**

# BENCHMARKS

QUERY	SUBJECT	Data dimensions	data.table <b>foverlaps</b>	GenomicRanges (bioconductor) <b>findOverlaps</b>
80M rows	33K rows	5 chromosomes, start, end N = 6 samples	2min	16min

*“GenomicRanges builds upon IRanges to add biological semantics to metadata”*

# BENCHMARKS

QUERY	SUBJECT	Data dimensions	data.table <b>foverlaps</b>	GenomicRanges (bioconductor) <b>findOverlaps</b>
80M rows	33K rows	5 chromosomes, start, end N = 6 samples	2min	16min
65M rows	35K rows	~7500 scaffolds, start, end N = 14 samples (time-course)	4min	>28hrs ~2hrs (parallel)

*“GenomicRanges builds upon IRanges to add biological semantics to metadata”*