



NGS perspective on IRanges package

- Representation of information on chromosomes/contigs
 - Intervals with or without associated data
 - Piecewise constant measures (e.g. coverage)
- Vector and interval operations for these representations
 - Interval overlap calculations
 - Coverage area within peak regions
- Metadata scheme for self-documenting objects



Two most important classes in IRanges

- *RangedData* - intervals and associated data on chromosomes/contigs. It can be conceptualized as a data table that is sorted by the chromosomes/contigs indicator column.
- *RleList* - coverage (or other piecewise constant measures) on chromosomes/contigs. RLE is an initialism for run length encoding, a standard compression method in signal processing.

RangedData class decomposition

- *RangedData*
 - *RangesList* - intervals on chromosomes/contigs. Extracted using the `ranges` function.
 - *Ranges* - intervals for a specific chromosome/contig. Most common subclass is *IRanges*.
 - *SplitDataFrameList* - data on chromosomes/contigs. Extracted using the `values` function.
 - *DataFrame* - data for a specific chromosome/contig.

Creating a new RangedData object

New object to use in interval operations

```
> ir <- IRanges(c(1, 8, 14, 15, 19, 34,  
+ 40), width = c(12, 6, 6, 15, 6,  
+ 2, 7))  
> strand <- rep(c("+", "-"), c(4, 3))  
> rd <- RangedData(ranges = ir, strand = strand,  
+ space = "chr1")
```



Low level data access

Accessors

```
> start(rd)
```

```
[1]  1  8 14 15 19 34 40
```

```
> end(rd)
```

```
[1] 12 13 19 29 24 35 46
```

```
> width(rd)
```

```
[1] 12  6  6 15  6  2  7
```

RangedData subsetting

```
> rd[1:5, ]
```

RangedData with 5 rows and 1 value column across 1 space

	space	ranges		strand
	<character>	<IRanges>		<character>
1	chr1	[1, 12]		+
2	chr1	[8, 13]		+
3	chr1	[14, 19]		+
4	chr1	[15, 29]		+
5	chr1	[19, 24]		-

Shifting intervals

- If your interval bounds are off by 1, you can shift them.

```
> rd2 <- rd
> ranges(rd2) <- shift(ranges(rd2), 1)
> rd2
```

RangedData with 7 rows and 1 value column across 1 space

	space	ranges		strand
	<character>	<IRanges>		<character>
1	chr1	[2, 13]		+
2	chr1	[9, 14]		+
3	chr1	[15, 20]		+
4	chr1	[16, 30]		+
5	chr1	[20, 25]		-
6	chr1	[35, 36]		-
7	chr1	[41, 47]		-

Resizing intervals (1/2)

- One common operation in ChIP-seq experiments is to “grow” and alignment interval to an estimated fragment length.

```
> rd3 <- rd
> pos <- values(rd3)[, "strand"] == "+"
> ranges(rd3)[pos] <- resize(ranges(rd)[pos],
+   120)
> ranges(rd3)[!pos] <- resize(ranges(rd)[!pos],
+   120, start = FALSE)
```

Resizing intervals (2/2)

```
> rd3
```

```
RangedData with 7 rows and 1 value column across 1 space
```

	space	ranges		strand
	<character>	<IRanges>		<character>
1	chr1	[1, 120]		+
2	chr1	[8, 127]		+
3	chr1	[14, 133]		+
4	chr1	[15, 134]		+
5	chr1	[-95, 24]		-
6	chr1	[-84, 35]		-
7	chr1	[-73, 46]		-

Restricting interval bounds

- The previous operation created some negative start values.
We can “clip” those negative values.

```
> ranges(rd3) <- restrict(ranges(rd3),
+   1)
> rd3
```

RangedData with 7 rows and 1 value column across 1 space

	space	ranges		strand
	<character>	<IRanges>		<character>
1	chr1	[1, 120]		+
2	chr1	[8, 127]		+
3	chr1	[14, 133]		+
4	chr1	[15, 134]		+
5	chr1	[1, 24]		-
6	chr1	[1, 35]		-
7	chr1	[1, 46]		-

Normalizing intervals

- *Ranges* can represent a set of integers
- *NormalIRanges* formalizes this, with a compact, normalized representation
- `reduce` normalizes ranges

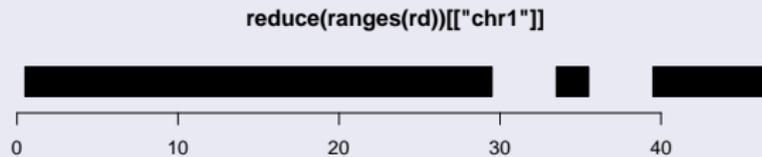
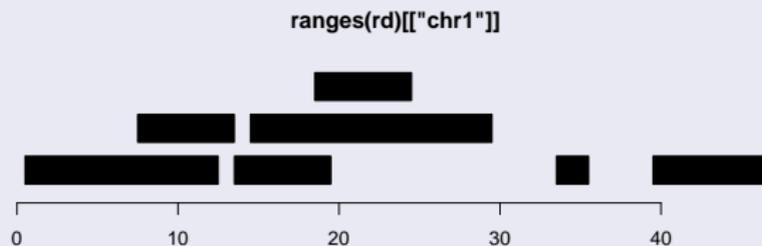
Code

```
> reduce(ranges(rd))
```

Normalizing intervals

Code

```
> reduce(ranges(rd))
```



Set operations

- *Ranges* as set of integers: `intersect`, `union`, `gaps`, `setdiff`
- Each range as integer set, in parallel: `pintersect`, `punion`, `pgap`, `psetdiff`

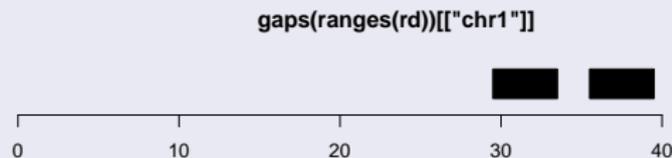
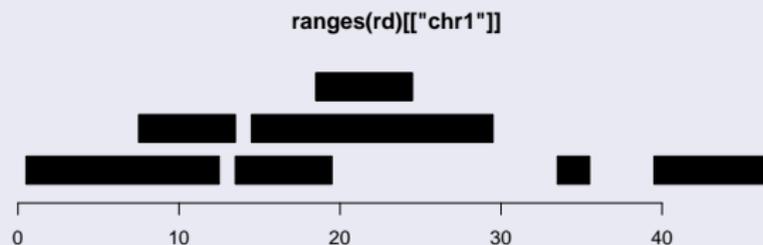
Example: `gaps`

```
> gaps(ranges(rd))
```

Set operations

Example: gaps

```
> gaps(ranges(rd))
```



Disjoining intervals

- Disjoint ranges are non-overlapping
- `disjoin` returns the widest ranges where the overlapping ranges are the same

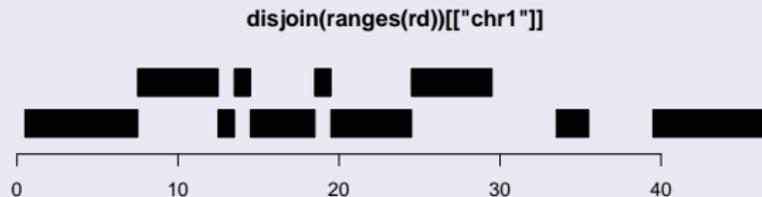
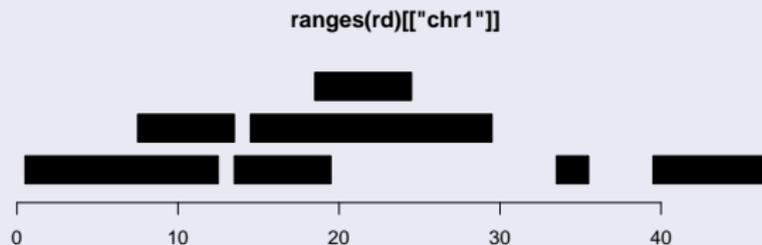
Code

```
> disjoin(ranges(rd))
```

Disjoining intervals

Code

```
> disjoint(ranges(rd))
```



Overlap detection

- `overlap` detects overlaps between two *Ranges* objects
- Uses interval tree for efficiency

Code

```
> ol <- findOverlaps(ranges(rd), reduce(ranges(rd)))  
> as.matrix(ol)
```

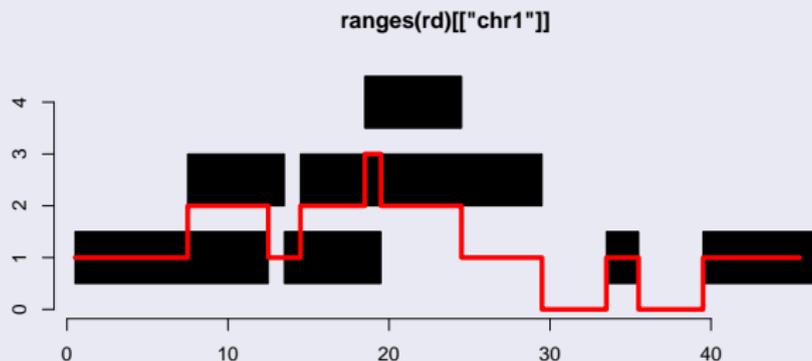
	query	subject
[1,]	1	1
[2,]	2	1
[3,]	3	1
[4,]	4	1
[5,]	5	1
[6,]	6	2
[7,]	7	3

Counting overlapping intervals

coverage counts number of ranges over each position

Code

```
> cover <- coverage(ranges(rd))
```



Finding nearest neighbors

- nearest finds the nearest neighbor ranges (overlapping is zero distance)
- precede, follow find non-overlapping nearest neighbors on specific side

Positional piecewise constant measures

- The number of genomic positions in a genome is often in the billions for higher organisms, making it challenging to represent in memory.
- Some data across a genome tend to be sparse (i.e. large stretches of “no information”)
- The *IRanges* package solves the set of problems for positional measures that tend to have consecutively repeating values.
- The *IRanges* package *does not* address the more general problem of positional measures that constantly fluxuate, such as conservation scores.

