

# ggbio: visualization tool kit for genomic data.

Tengfei Yin

Iowa State University

July 31, 2012

# Outline

## ① Introduction

- Background
- Motivation

## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

## ③ Conclusion

# Outline

## ① Introduction

- Background
- Motivation

## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

## ③ Conclusion

# Outline

## ① Introduction

Background

Motivation

## ② ggbio

Grammar of graphics

Extension

Overview and layout

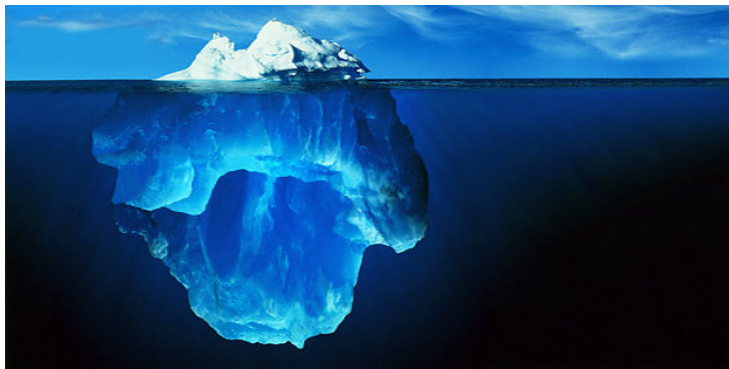
Specialized plots

## ③ Conclusion

# Tip of iceberg

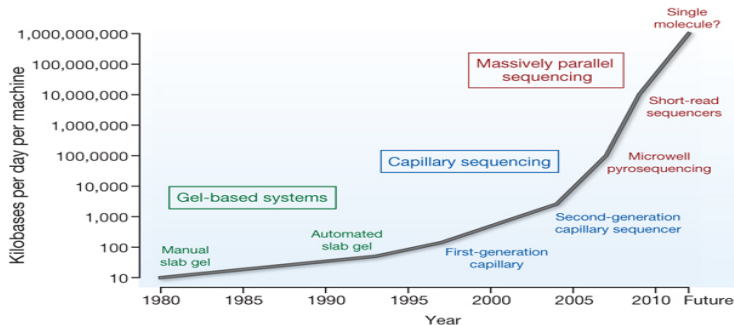
## Definition

**big data** is a loosely-defined term used to describe data sets so large and complex that they become awkward to work with using on-hand database management tools. Difficulties include capture, storage, search, sharing, analysis, and **visualization**(*wiki*).



# Sequencing data: the iceberg.

- Next generation sequencing data is huge, and keep growing.
  - Raw image over TB.
  - Over GB results data per run.
- Analysis to find *tip* of that “Ice berg”.
  - Analytical tool kits.
  - Visualization tool kits.

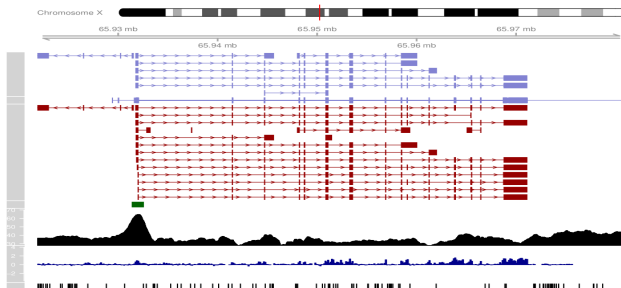


# Review of current tools

- GenomeGraphs, Gviz(in R).
- Interactive Desktop-app: IGB/IGV(Java).
- Web based: UCSC genome browser, DNAnexus.
- Specialized: Circos.

# Review of current tools

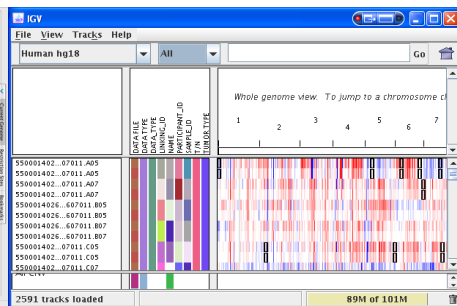
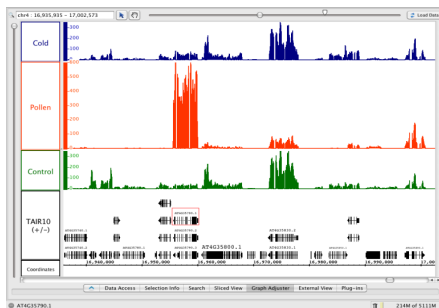
- GenomeGraphs, Gviz(in R).
- Interactive Desktop-app: IGB/IGV(Java).
- Web based: UCSC genome browser, DNAnexus.
- Specialized: Circos.





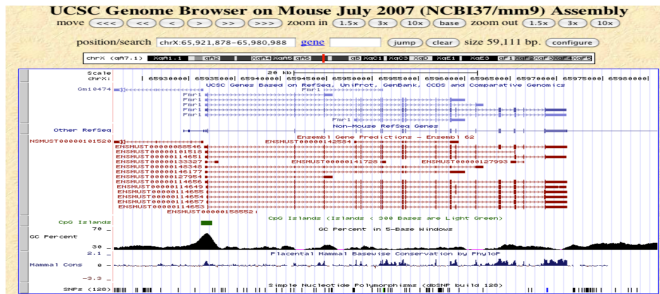
# Review of current tools

- GenomeGraphs, Gviz(in R).
- Interactive Desktop-app: IGB/IGV(Java).
- Web based: UCSC genome browser, DNAnexus.
- Specialized: Circos.



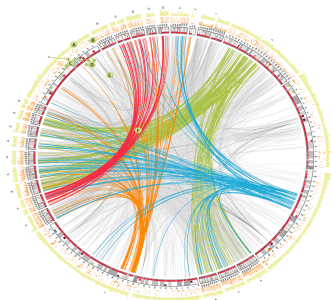
# Review of current tools

- GenomeGraphs, Gviz(in R).
- Interactive Desktop-app: IGB/IGV(Java).
- Web based: UCSC genome browser, DNAnexus.
- Specialized: Circos.



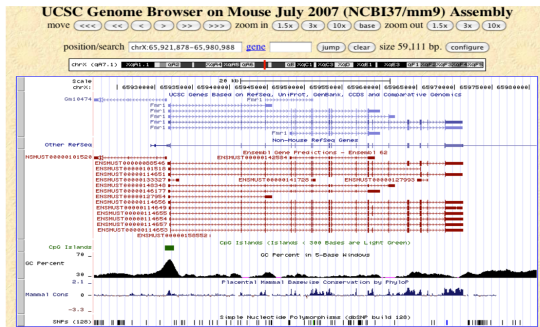
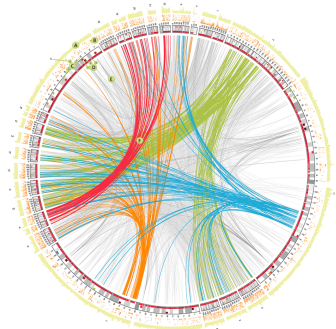
# Review of current tools

- GenomeGraphs, Gviz(in R).
- Interactive Desktop-app: IGB/IGV(Java).
- Web based: UCSC genome browser, DNAnexus.
- Specialized: Circos.



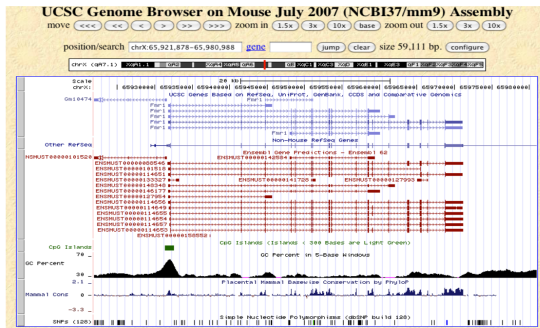
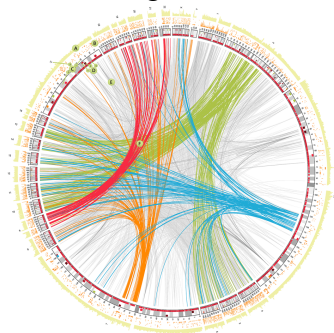
# Generalization possible?

Do they really look so different?



# Generalization possible?

Possible to generalize!



# Outline

## ① Introduction

Background

Motivation

## ② ggbio

Grammar of graphics

Extension

Overview and layout

Specialized plots

## ③ Conclusion

# Motivation

- Explore the data in different ways.
- Generalize genomic visualization frame work.
- Elegant graphics for publications.
- Modularize components to facilitate construction of high level graphics.

# Design basic

- A powerful computational platform.
  - R: statistical platform, numerous model.
- A general data model.
  - Bioconductor: For I/O, data model, and some analytical.
- A general graphic model and grammar.
  - *ggplot2*: A grammar of graphics in R.



# Design basic

- A powerful computational platform.
  - R: statistical platform, numerous model.
- A general data model.
  - Bioconductor: For I/O, data model, and some analytical.
- A general graphic model and grammar.
  - *ggplot2*: A grammar of graphics in R.



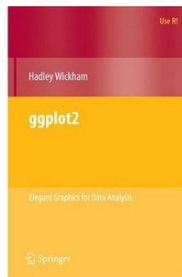
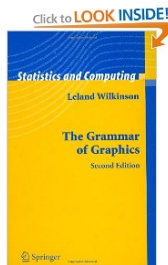
# Design basic

- A powerful computational platform.
  - R: statistical platform, numerous model.
- A general data model.
  - Bioconductor: For I/O, data model, and some analytical.
- A general graphic model and grammar.
  - *ggplot2*: A grammar of graphics in R.



# Design basic

- A powerful computational platform.
  - R: statistical platform, numerous model.
- A general data model.
  - Bioconductor: For I/O, data model, and some analytical.
- A general graphic model and grammar.
  - *ggplot2*: A grammar of graphics in R.



# Outline

## ① Introduction

- Background
- Motivation

## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

## ③ Conclusion

# Outline

## ① Introduction

Background

Motivation

## ② ggbio

Grammar of graphics

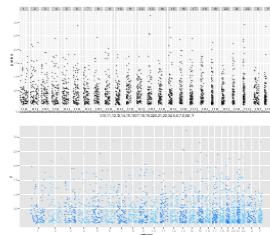
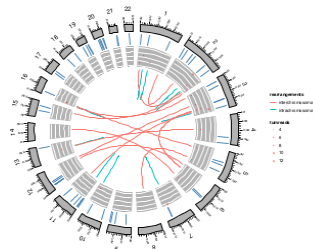
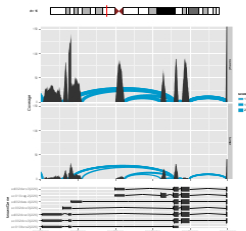
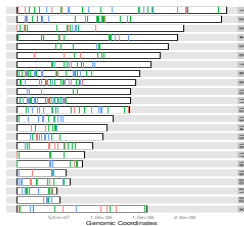
Extension

Overview and layout

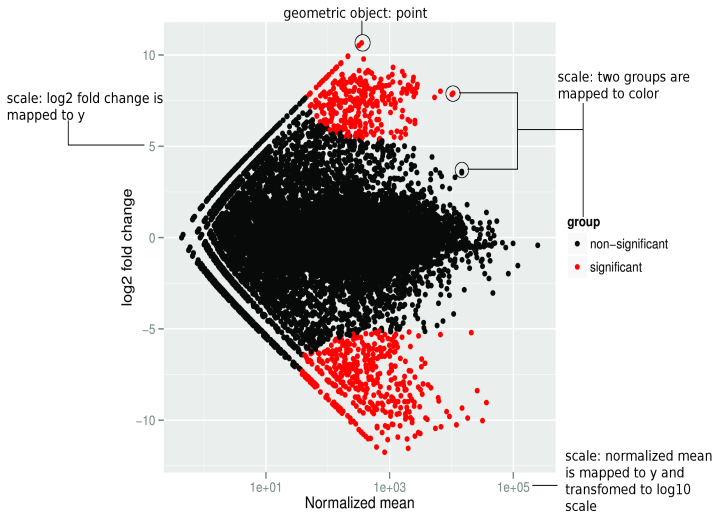
Specialized plots

## ③ Conclusion

# What can ggbio do?



# What is grammar of graphics(GoG)



# Outline

## ① Introduction

Background

Motivation

## ② ggbio

Grammar of graphics

Extension

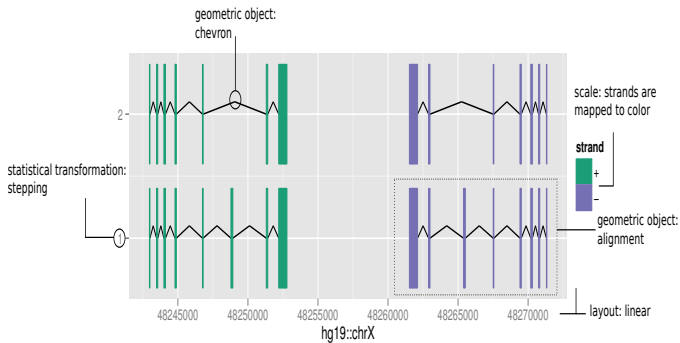
Overview and layout

Specialized plots

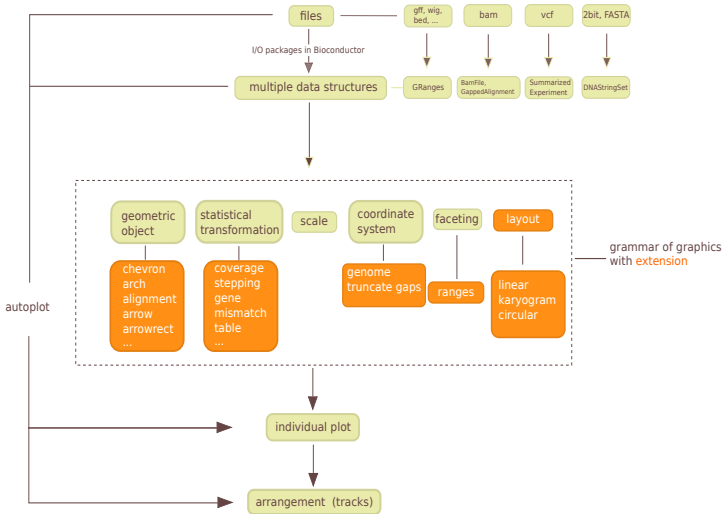
## ③ Conclusion






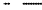










# GoG in Genomic World.











# Pipeline in ggbio



# Extended grammar of graphics table I

Comp	name	usage	icon
<b>geom</b>	geom_rect	rectangle	
	geom_segment	segment	
	geom_chevron	chevron	
	geom_arrow	arrow	
	geom_arch	arches	
	geom_bar	bar	
	geom_alignment	alignment (gene)	
<b>stat</b>	stat_coverage	coverage (of reads)	
	stat_mismatch	mismatch pileup for alignments	
	stat_aggregate	aggregate in sliding window	
	stat_stepping	avoid overplotting	
	stat_gene	consider gene structure	
	stat_table	tabulate ranges	
	stat_identity	no change	

# Extended grammar of graphics table II

<b>coord</b>	linear	ggplot2 linear but facet by chromosome	
	genome	put everything on genomic coordinates	
	truncate gaps	compact view by shrinking gaps	
<b>layout</b>	track	stacked tracks	
	karyogram	karyogram display	
	circle	circular	
<b>faceting</b>	formula	facet by formula	
	ranges	facet by ranges	
<b>scale</b>	not extended	<i>ggplot2</i> default	

# Extended grammar of graphics table III

# Outline

## ① Introduction

- Background
- Motivation

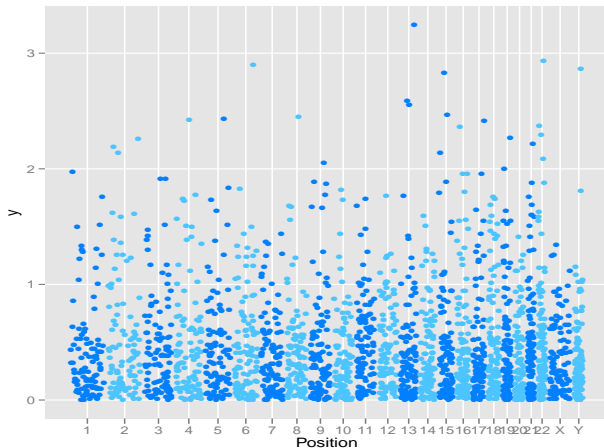
## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

## ③ Conclusion

# Manhattan plot(Grandlinear)

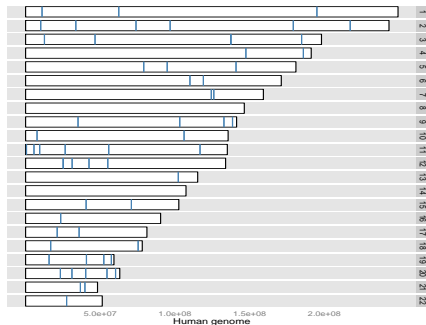
```
> autoplot(obj, coord = "genome", ...)
```



# Karyogram layout

@

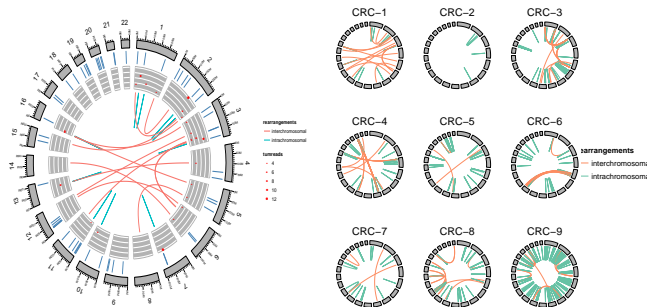
```
> autoplot(obj, layout = "karyogram", ...)
```





# Circular layout

```
> autoplot(obj, layout = "circle", ...)
```



# Outline

## ① Introduction

- Background
- Motivation

## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

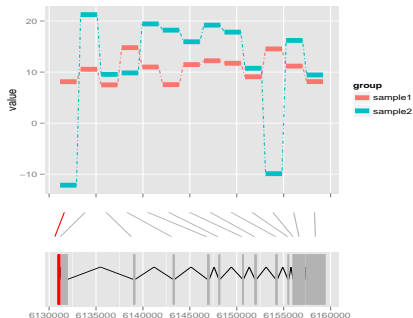
## ③ Conclusion

# Specialized plots.

- Higher level: for convenience usage.
- Complex than prototypes, usually integrated with tracks.

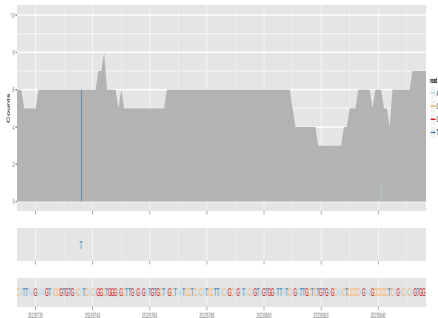
## 'ranges-linked-to-data' plot.

```
> plotRangesLinkedToData(obj, ...)
```



## Mismatch summary plot

```
> autoplot(obj, stat = 'mismatch', ...)
```

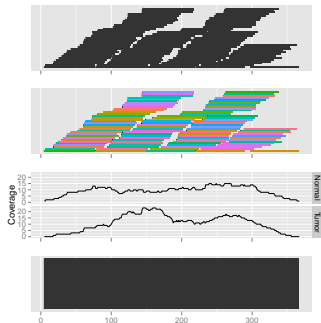


# Object supported(overview).

<i>GRanges</i>				

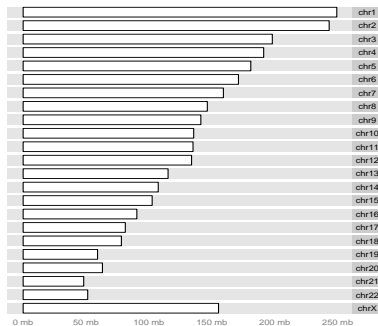
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>			



# Object supported(overview).

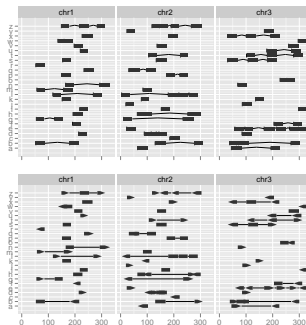
<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>		





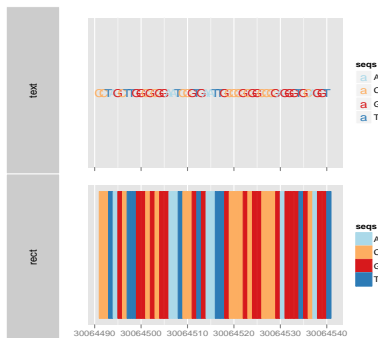
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	



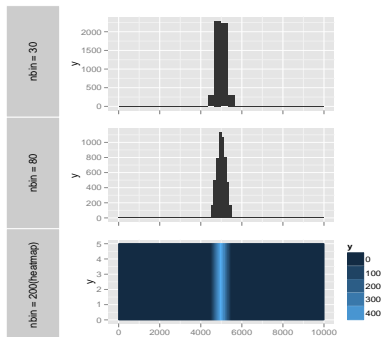
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>



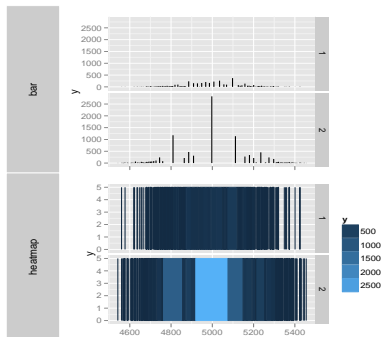
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>				



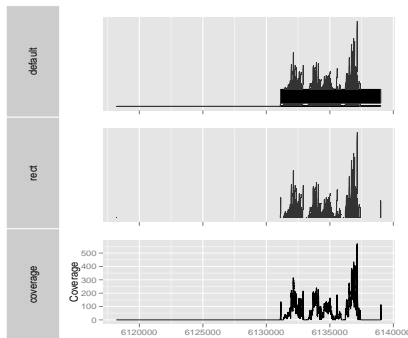
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>			



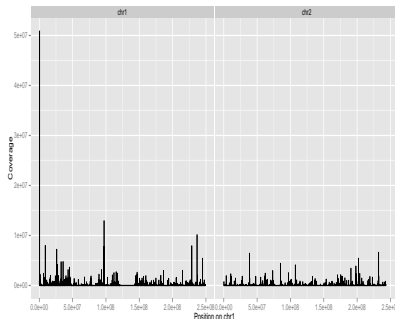
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>		



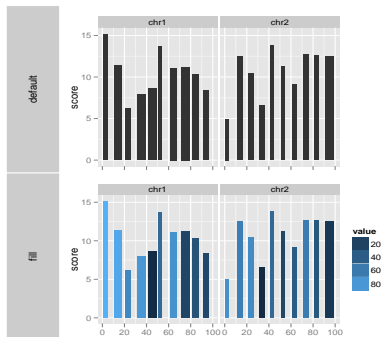
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	



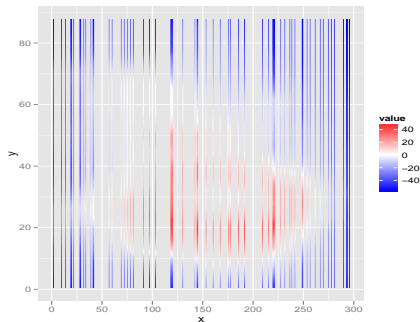
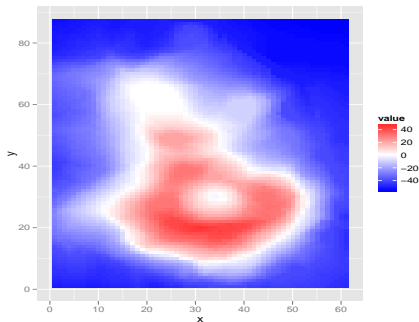
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>



# Object supported(overview).

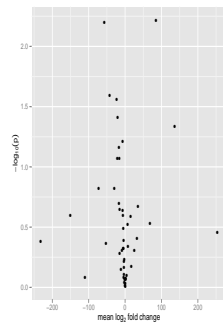
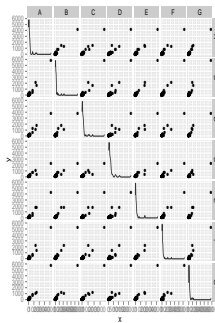
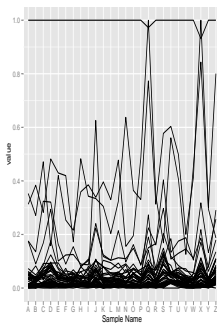
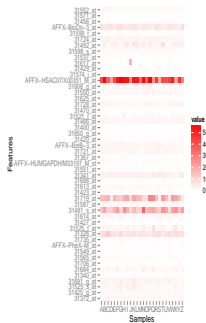
<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>
<i>matrix</i>				





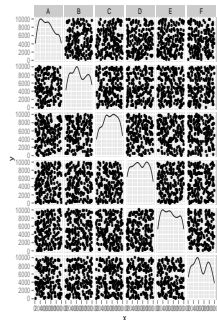
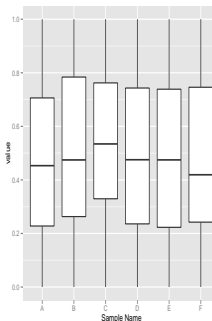
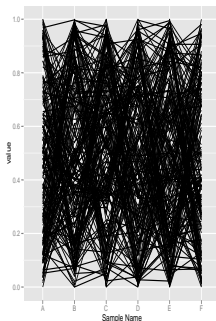
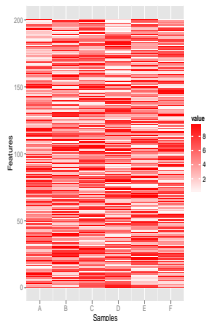
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>
<i>matrix</i>	<i>ExpressionSet</i>			



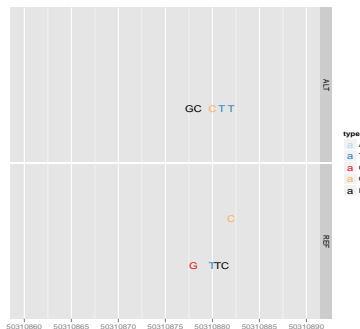
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>
<i>matrix</i>	<i>ExpressionSet</i>	<i>SummarizedExperiment</i>		



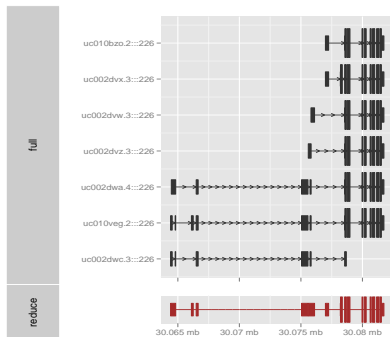
# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>
<i>matrix</i>	<i>ExpressionSet</i>	<i>SummarizedExperiment</i>	<i>VCF</i>	



# Object supported(overview).

<i>GRanges</i>	<i>IRanges</i>	<i>Seqinfo</i>	<i>GRangesList</i>	<i>BSgenome</i>
<i>Rle</i>	<i>RleList</i>	<i>GappedAlignment</i>	<i>BamFile</i>	<i>character</i>
<i>matrix</i>	<i>ExpressionSet</i>	<i>SummarizedExperiment</i>	<i>VCF</i>	<i>TranscriptDb</i>



# Outline

## ① Introduction

- Background
- Motivation

## ② ggbio

- Grammar of graphics
- Extension
- Overview and layout
- Specialized plots

## ③ Conclusion

# Future study

- Support more core data model in Bioconductor.
- More elegant theme for tracks.
- Keep improving with new ggplot2 development
- More powerful tracks function, may accept lattice graphics.

# Acknowledgment

- Michael Lawrence, Dianne Cook
- Genentech

Thank you !!!