

# GenomeSpy Workshop

At Oslo Bioinformatics Workshop Week 2024

11.12.2024

# About Me

- Kari Lavikka
  - PhD student at Hautaniemi Lab,  
Medical Faculty, University of Helsinki
  - MSc in Data Science, BSc in Biology
  - Author of GenomeSpy
- 
- Email: [kari.lavikka@helsinki.fi](mailto:kari.lavikka@helsinki.fi)
  - Website: <https://karilavikka.fi/>

# Hautaniemi Lab



Prof.  
Hautaniemi

That's Me!

# Deciphering cancer genomes with GenomeSpy: a grammar-based visualization toolkit

Kari Lavikka <sup>1,\*</sup>, Jaana Oikkonen <sup>1</sup>, Yilin Li <sup>1</sup>, Taru Muranen <sup>1</sup>, Giulia Micoli <sup>1</sup>, Giovanni Marchi <sup>1</sup>, Alexandra Lahtinen <sup>1</sup>, Kaisa Huhtinen <sup>1,2</sup>, Rainer Lehtonen <sup>3</sup>, Sakari Hietanen <sup>4</sup>, Johanna Hynninen <sup>4</sup>, Anni Virtanen <sup>5</sup>, and Sampsia Hautaniemi <sup>1,\*</sup>

<sup>1</sup>Research Program in Systems Oncology, Research Programs Unit, Faculty of Medicine, University of Helsinki, 00014 Helsinki, Finland

<sup>2</sup>Cancer Research Unit, Institute of Biomedicine and FICAN West Cancer Centre, University of Turku, 20521 Turku, Finland

<sup>3</sup>Applied Tumor Genomics Research Program, Research Programs Unit, Faculty of Medicine, University of Helsinki, 00014 Helsinki, Finland

<sup>4</sup>Department of Obstetrics and Gynecology, University of Turku and Turku University Hospital, 20521 Turku, Finland

<sup>5</sup>Department of Pathology, University of Helsinki and HUS Diagnostic Center, Helsinki University Hospital, 00260 Helsinki, Finland

\*Correspondence address. Kari Lavikka, Research Program in Systems Oncology, Research Programs Unit, Faculty of Medicine, University of Helsinki, 00014 Helsinki, Finland, E-mail: [kari.lavikka@helsinki.fi](mailto:kari.lavikka@helsinki.fi); Sampsia Hautaniemi, Research Program in Systems Oncology, Research Programs Unit, Faculty of Medicine, University of Helsinki, 00014 Helsinki, Finland, E-mail: [sampsia.hautaniemi@helsinki.fi](mailto:sampsia.hautaniemi@helsinki.fi)

## Abstract

**Background:** Visualization is an indispensable facet of genomic data analysis. Despite the abundance of specialized visualization tools, there remains a distinct need for tailored solutions. However, their implementation typically requires extensive programming expertise from bioinformaticians and software developers, especially when building interactive applications. Toolkits based on visualization grammars offer a more accessible, declarative way to author new visualizations. Yet, current grammar-based solutions fall short in adequately supporting the interactive analysis of large datasets with extensive sample collections, a pivotal task often encountered in cancer research.

**Findings:** We present GenomeSpy, a grammar-based toolkit for authoring tailored, interactive visualizations for genomic data analysis. By using combinatorial building blocks and a declarative language, users can implement new visualization designs easily and embed

# Program

- Background: DECIDER project
- Background: Why yet another tool
- SegmentModel Spy demo
- DECIDER visualization demo and exercise
- Visualization Grammars
- Code-Along: Let's make a ClinVar visualization
- Using in Observable notebooks
- Local development, remote deployment
- Questions, discussion, etc.

Survey

# Background

Why was GenomeSpy developed?

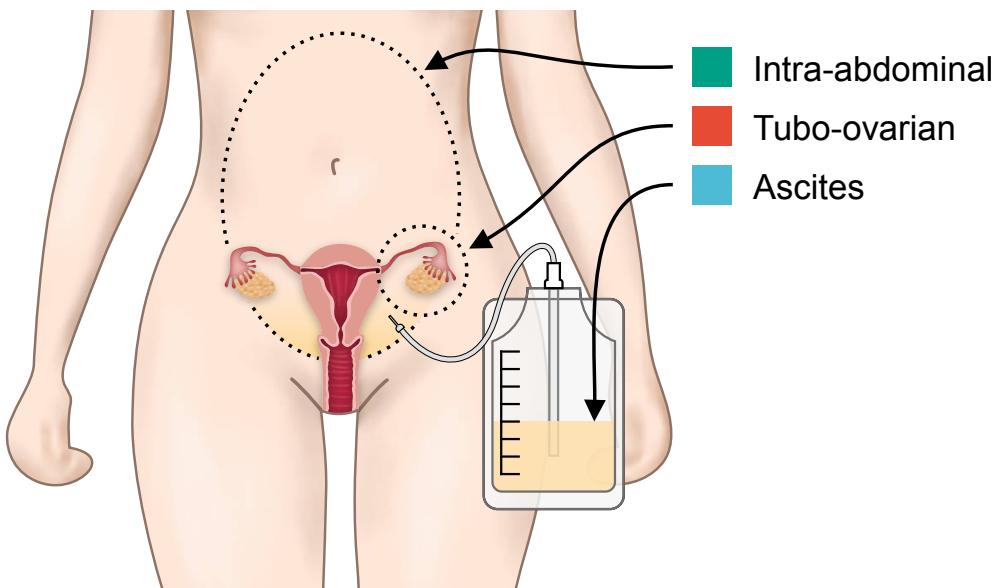


Improving clinical decisions in cancer

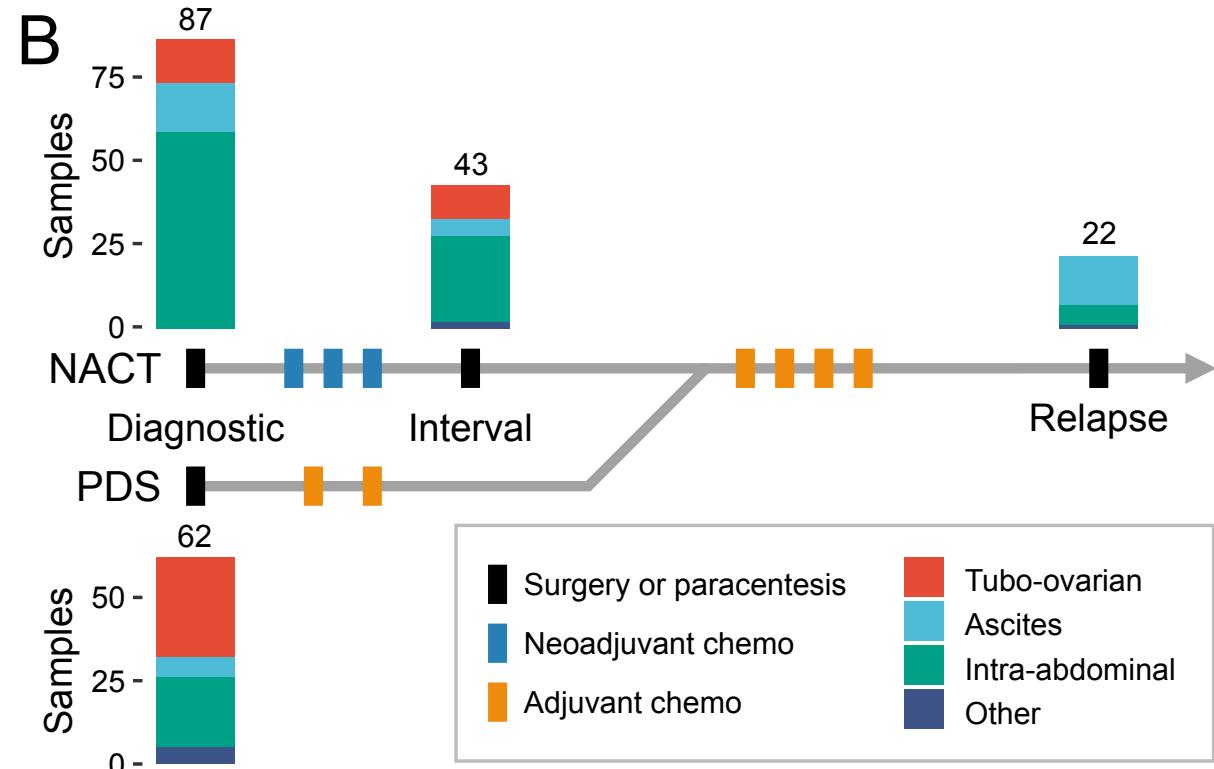
- Ovarian High-Grade Serous Carcinoma (HGSC)
- Diagnostic tools and treatments
- Overcoming Drug Resistance
- Funded by Horizon 2020

# DECIDER HGSC Cohort

A



B



A. Lahtinen et al., "Evolutionary states and trajectories characterized by distinct pathways stratify patients with ovarian high grade serous carcinoma," *Cancer Cell*, May 2023, doi: 10.1016/j.ccr.2023.04.017.

# Geneticists and others work with...

- ~1500 WGS samples
- Copy numbers
- SNVs / Indels
- Structural Variants
- RRBS Methylation
- Expression data

# Geneticists and others need to...

- Explore the whole cohort for patterns
- Stratify the cohort in various ways
- Study individual patients
- Rapidly check hypotheses
- Etc...

The State of the Art

# awesome-genome-visualization

- <https://cmdcolin.github.io/awesome-genome-visualization/>
- Compiled by Colin Diesh  
(the lead JBrowse 2 author)



cmdcolin

# UCSC XENA



UNIVERSITY OF CALIFORNIA  
SANTA CRUZ

DATA SETS

VISUALIZATION

TRANSCRIPTS

DATA HUBS

VIEW MY DATA

BOOKMARK

HELP

MORE TOOLS

TCGA Prostate Cancer (PRAD)

Filtered to 492 Samples

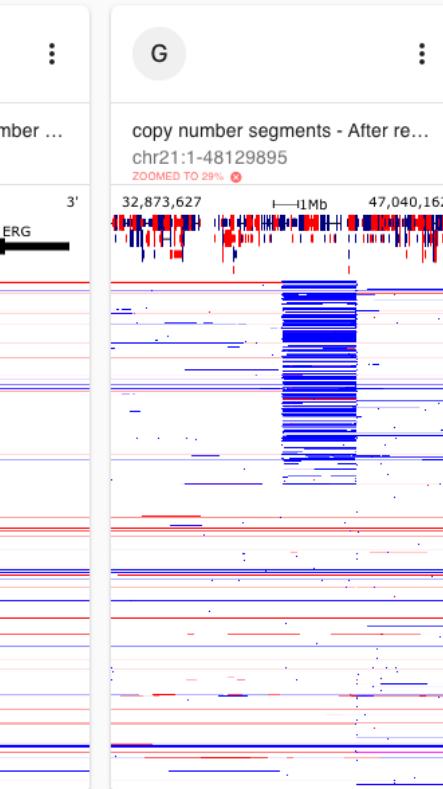
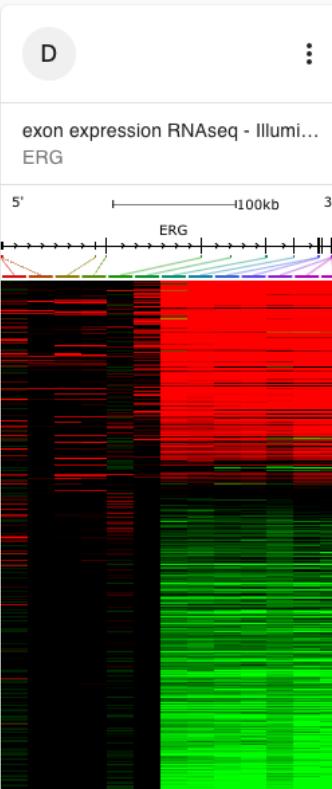
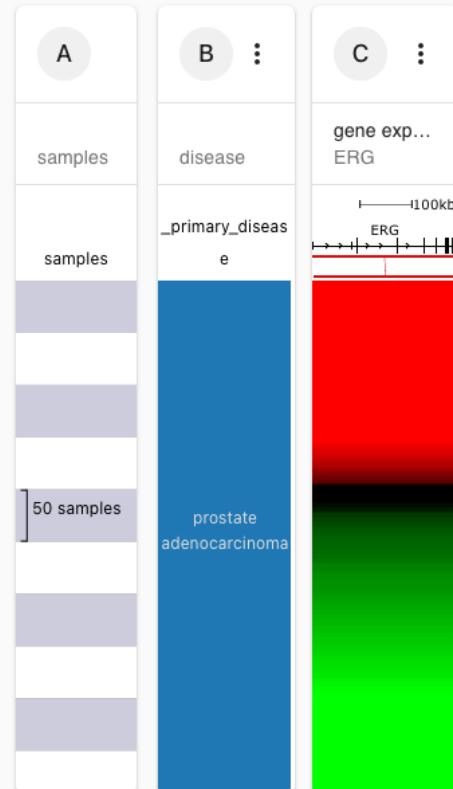


Type here or use dropper to select samples



ZOOM: None

User Guide • Zoom Help • Tooltip Help



Click to Add Column

7.3

log2(norm\_ count+1)

12

low high

log2(RPKM+1)

-0.93

0.93

log2(tumor/normal)

-0.39

0.39

log2(tumor/normal)

-0.39

0.39

log2(tumor/normal)

IGV



# GenomeSpy

How was it born?



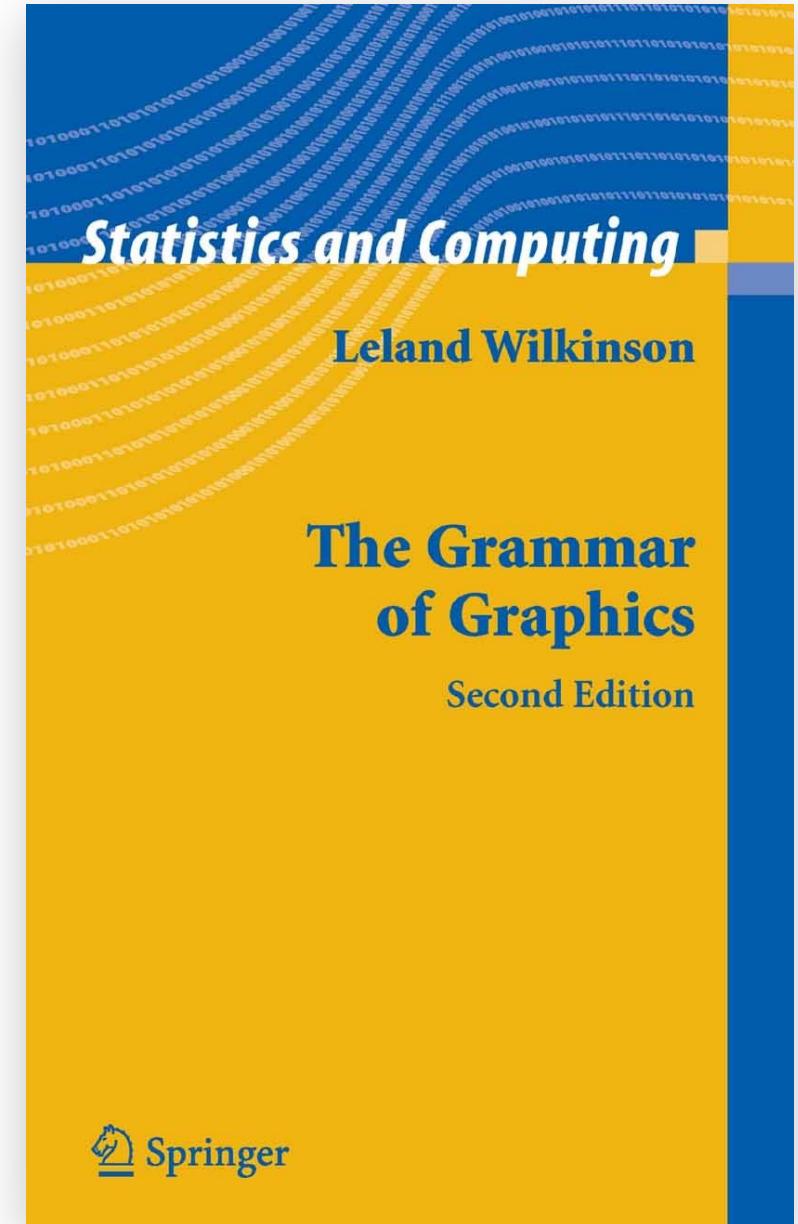
极客湾



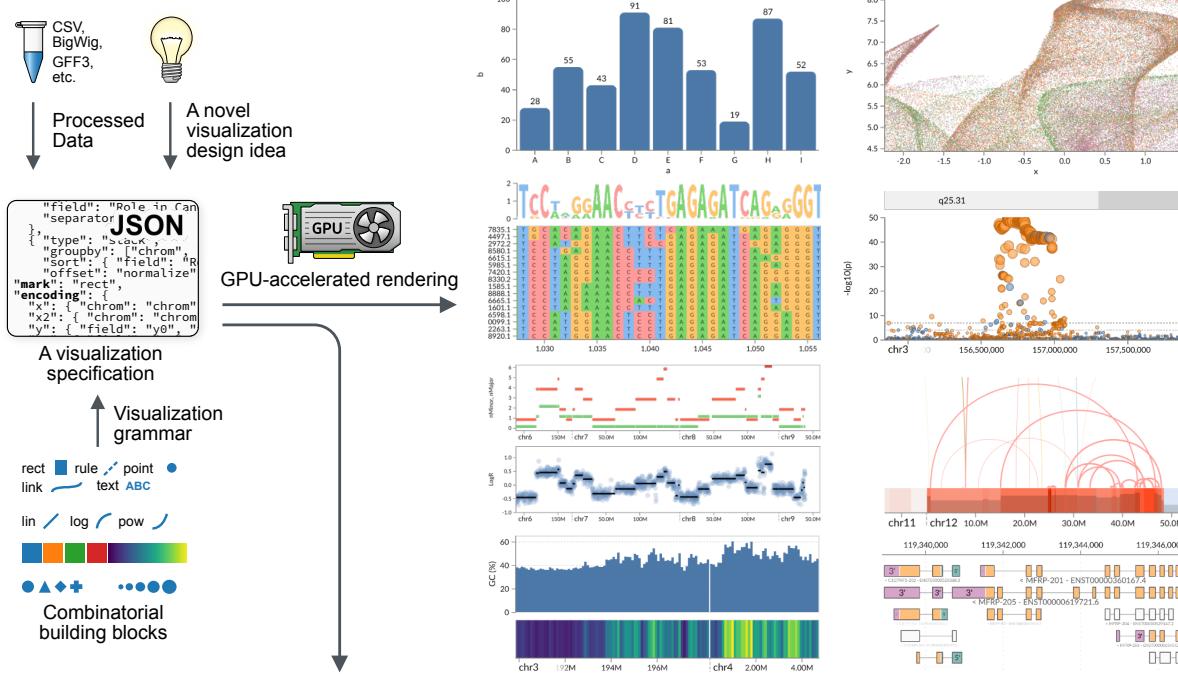
By 极客湾Geekerwan, CC BY 3.0, <https://commons.wikimedia.org/w/index.php?curid=151470847>

# Genomics Game Engine?

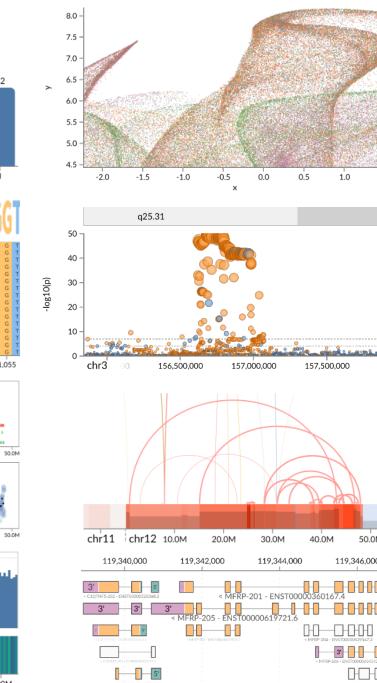
- GPU programing is hard
- What would be a suitable abstraction?
- A Visualization Grammar!
  - Examples: `ggplot2`, Vega-Lite
  - More on that later...
- Developing tailored GPU-accelerated visualizations is now easy
  - At least for me 😎



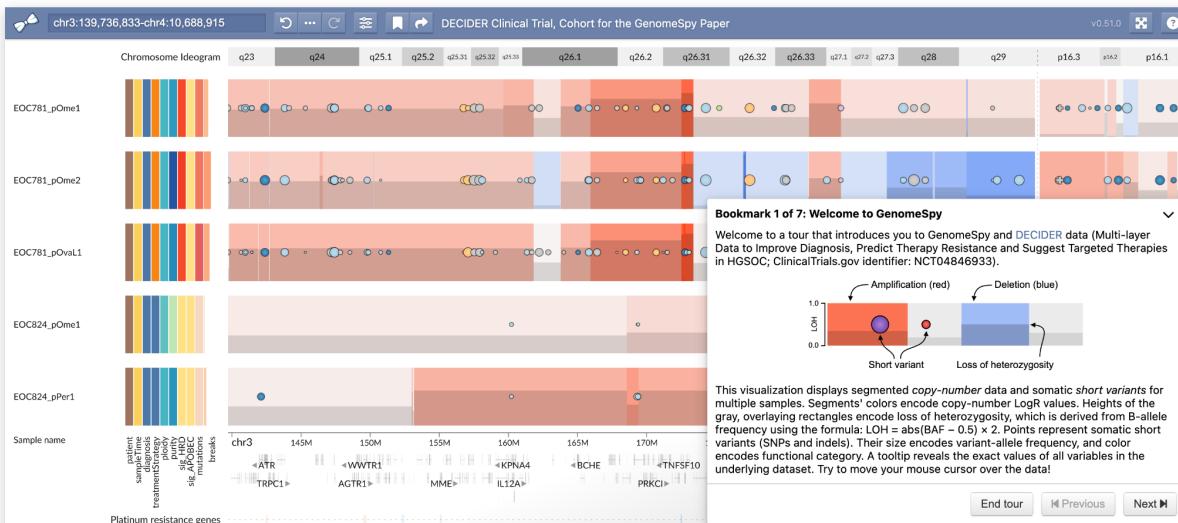
### A Authoring tailored visualizations



### B Tailored, interactive visualizations rendered by GenomeSpy Core



### C Analyzing sample collections using GenomeSpy App



# Demo: SegmentModel Spy

- GATK has tons of parameters with unclear documentation.  
How to study their effect?
- A simple app for copy-number segmentation assessment
- Try it yourself: <https://genomespy.app/segmentmodel/>

# Let's Explore the DECIDER Cohort

- Hautaniemi Lab geneticists uses this every day
- Navigate to: <https://tinyurl.com/mr4d9tev>
- Exercise: HGSC has a unique copy-number landscape. Do the copy-number peaks correlate with other peaks?

# Visualization Grammars

A Brief Introduction

# What Are Visualization Grammars?

- Visualization grammars are frameworks that describe how data is **transformed** into visualizations through a set of **composable rules**.
- They provide a consistent and systematic approach to creating complex visualizations by breaking them into fundamental **components**.

# Wilkinson's Grammar

ELEMENT: *line(position(smooth.linear(age\*bp)))*

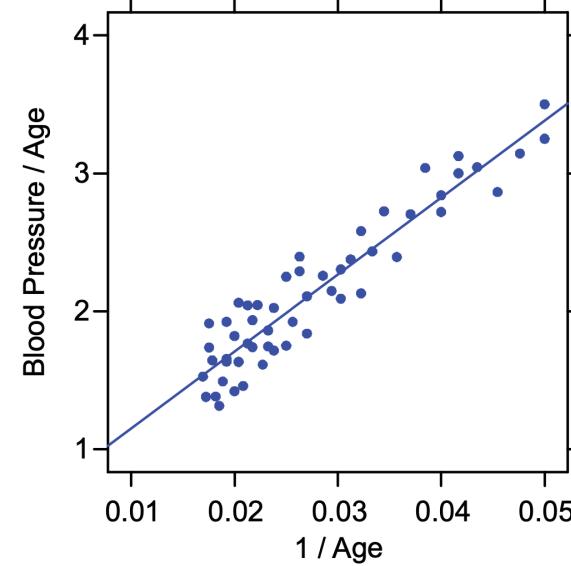
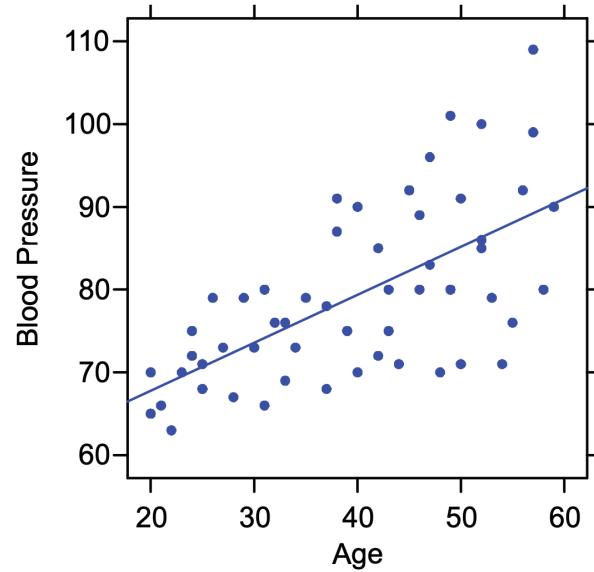
ELEMENT: *point(position(age\*bp))*

TRANS: *inverseage* = *inverse(age)*

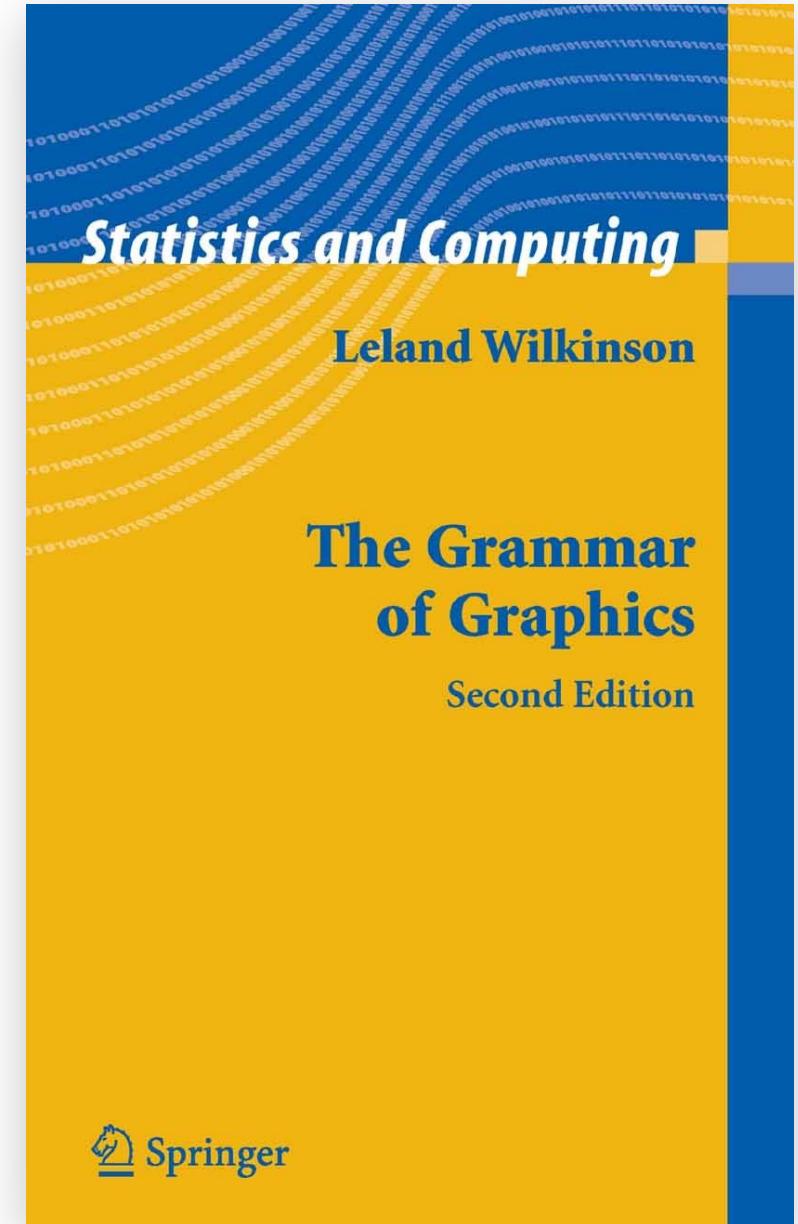
TRANS: *bpbyage* = *ratio(bp, age)*

ELEMENT: *line(position(smooth.linear(inverseage\*bpbyage)))*

ELEMENT: *point(position(inverseage\*bpbyage))*

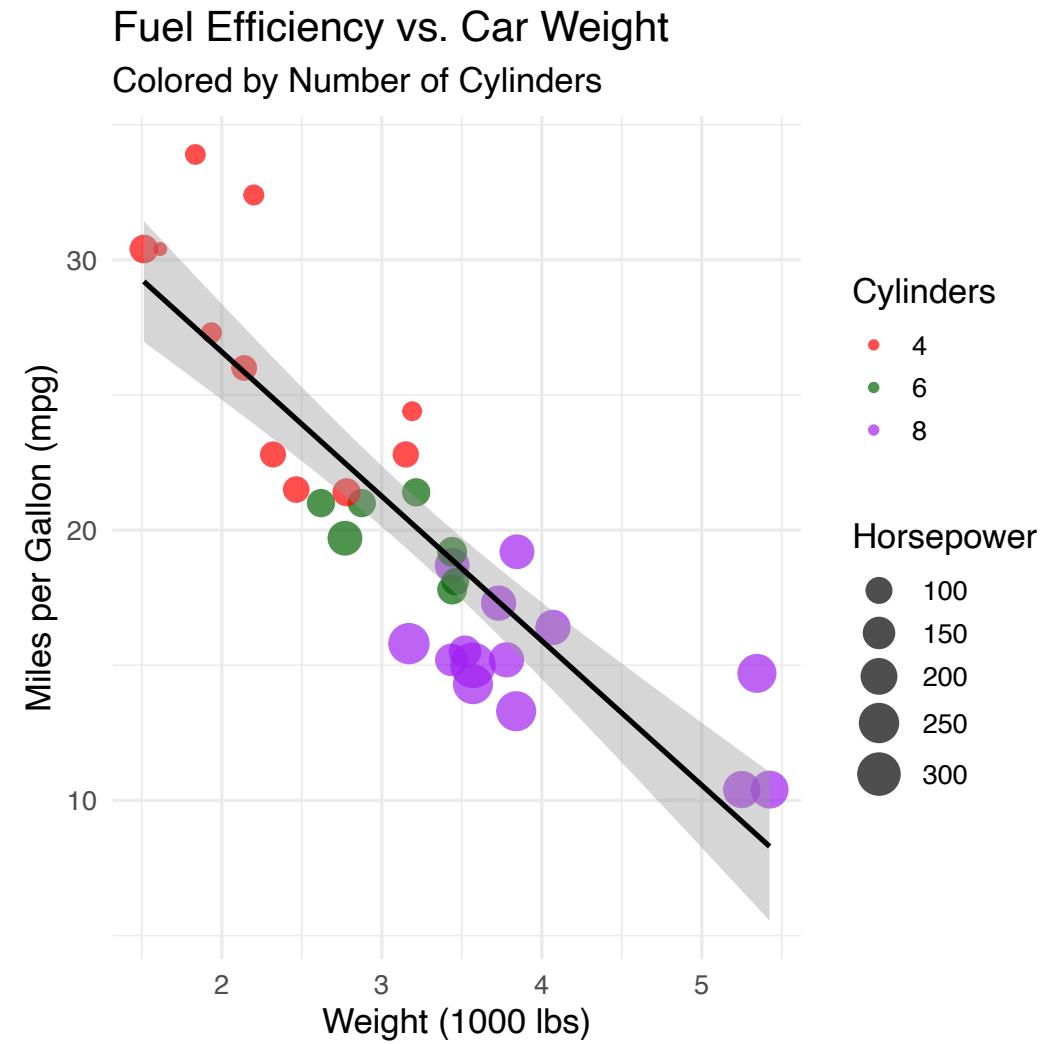


**Figure 9.14** Weighted least squares via transformation of variables



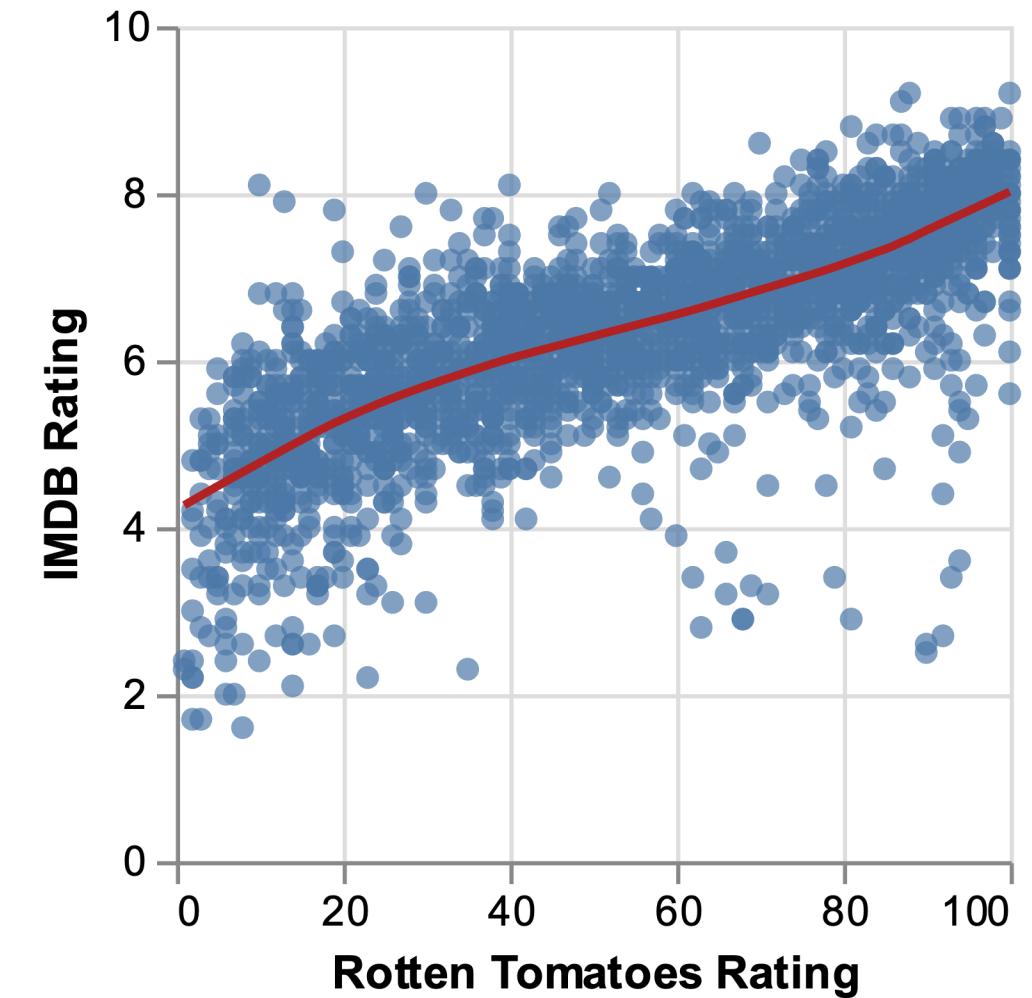
# ggplot2

```
ggplot(data = mtcars, aes(x = wt, y = mpg)) +
  geom_point(aes(color = factor(cyl), size = hp),
             alpha = 0.7) +
  geom_smooth(method = "lm", color = "black") +
  scale_color_manual(values = c(
    "4" = "red",
    "6" = "darkgreen",
    "8" = "purple")) +
  scale_size_continuous(range = c(2, 8)) +
  labs(
    title = "Fuel Efficiency vs. Car Weight",
    subtitle = "Colored by Number of Cylinders",
    x = "Weight (1000 lbs)",
    y = "Miles per Gallon (mpg)",
    color = "Cylinders",
    size = "Horsepower"
  ) +
  theme_minimal(base_size = 14)
```



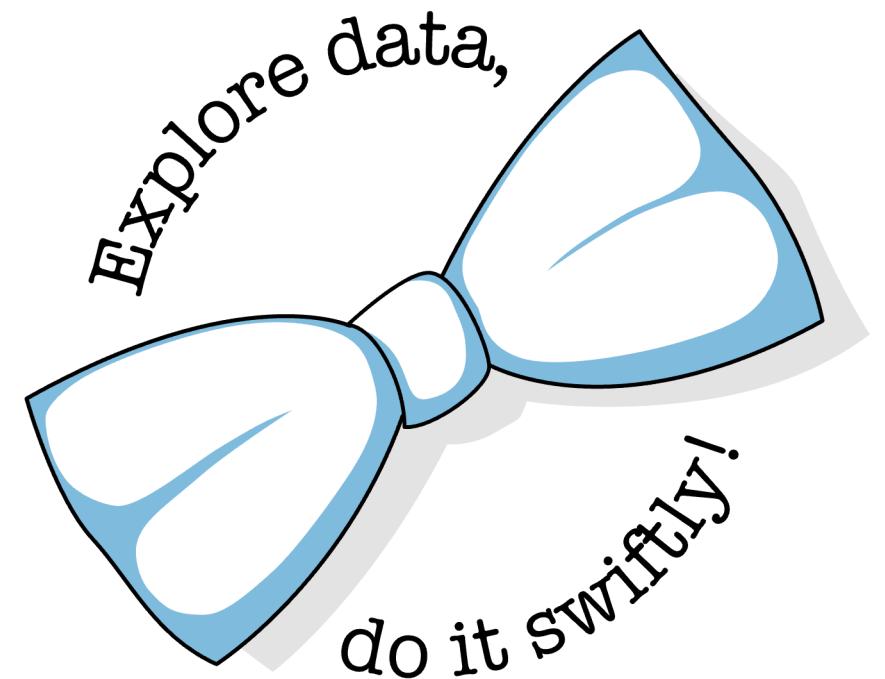
# Vega-Lite

```
{  
  "data": { "url": "data/movies.json" },  
  "encoding": {  
    "x": {  
      "field": "Rotten Tomatoes Rating",  
      "type": "quantitative"  
    },  
    "y": {  
      "field": "IMDB Rating",  
      "type": "quantitative"  
    }  
  },  
  "layer": [  
    {  
      "mark": { "type": "point", "filled": true }  
    },  
    {  
      "transform": [{  
        "loess": "IMDB Rating",  
        "on": "Rotten Tomatoes Rating"  
      }],  
      "mark": { "type": "line", "color": "firebrick" }  
    }  
  ]  
}
```



# GenomeSpy

- **HEAVILY** inspired by Vega-Lite
- Provides partial compatibility
- Genome-related enhancements
- Independent implementation
  - Scale transformations and rendering using GPU shaders
- <https://genomespy.app/>



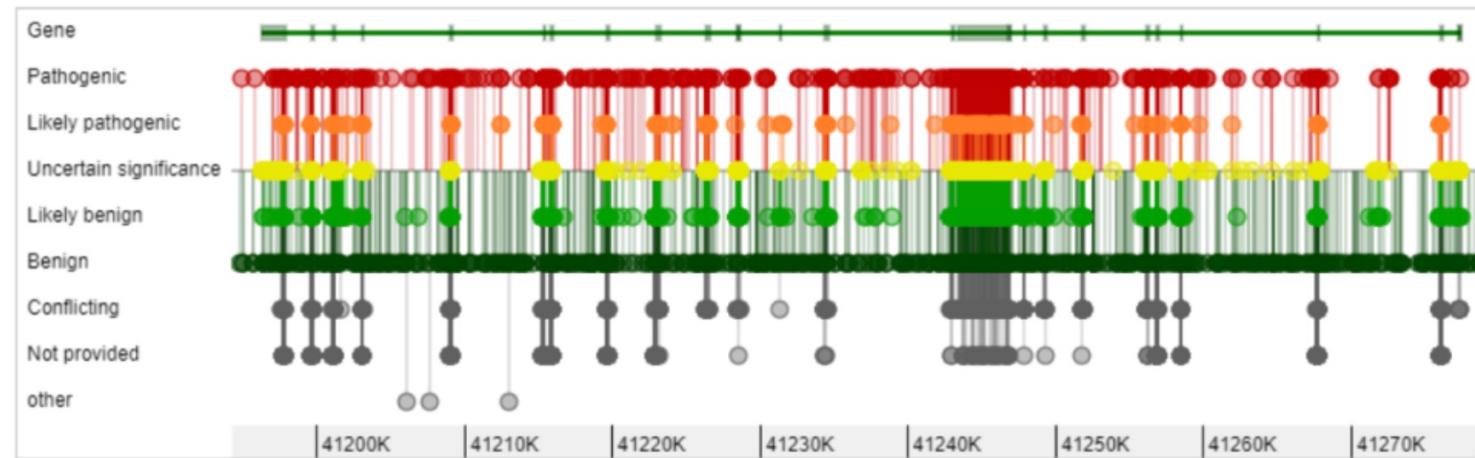
# Code-Along

Let's implement a visualization!

NEW

# ClinVar

## *New Visualization Tool*



## New ClinVar graphical display

*Maps clinically significant variants by gene and position!*

<https://ncbiinsights.ncbi.nlm.nih.gov/2022/08/30/clinvar-graphical-view/>

# GenomeSpy Playground

- A code editor and an interactive visualization
- For small-scale testing and prototyping
- <https://genomespy.app/playground/>

# ObservableHQ Notebooks

Embedding Visualizations

<https://observablehq.com/collection/@tuner/genomespy>

# Local Development, Remote Deployment

Developing more complex visualizations and deploying them  
on remote web servers

# Local Development

1. Open VS Code or another JSON-schema aware editor
2. Create an HTML file
3. Write a visualization spec
4. Start a local web server
  - `python -m http.server`
  - <https://github.com/danvk/RangeHTTPServer>
  - <https://www.npmjs.com/package/http-server>

# Remote Deployment

- No server-side logic needed
- Just put your files on a web server
  - ... and ensure that file permissions are ok

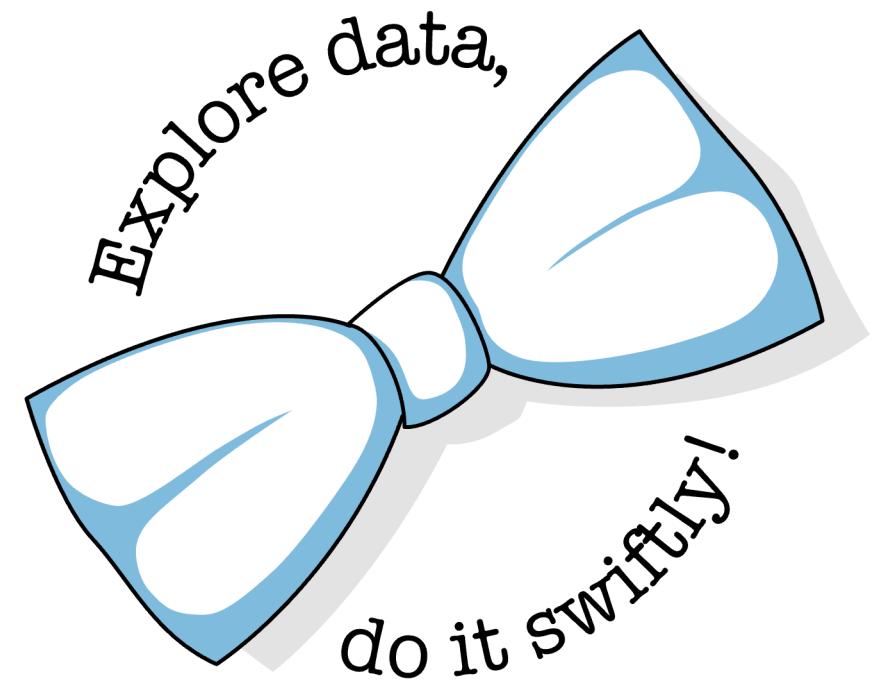
# Conclusions

# The Types of Visualization Challenges

- Tailored visualizations
- Large sample collections
- Finding patterns in large datasets

# Future Directions

- Better developer experience
- Growing the community
- Extending the grammar
- Analysis app improvements
  - RNA expression data
  - Offer simple statistics (ANOVA, etc.)



Questions?

# Follow or Contact Me for Updates

- Bluesky: @karilavikka.fi
- X: @KariLavikka
- Linkedin: <https://www.linkedin.com/in/karilavikka/>
- E-mail: kari.lavikka@helsinki.fi

# Evaluation form

