# Reverse Engineering Chart Data with WebPlotDigitizer
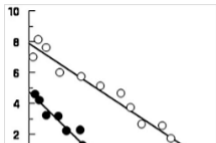
**Ankit Rohatgi** <ankitrohatgi@hotmail.com>

May 3, 2017

🌐 http://arohatgi.info/WebPlotDigitizer
🐙 http://github.com/ankitrohatgi/WebPlotDigitizer

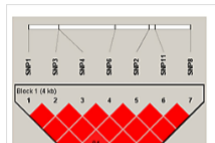# Raw Data?



**Representative Scatchard plot of HNECA saturation binding data**
Laura Bazzichi ⌄                30/12/2011
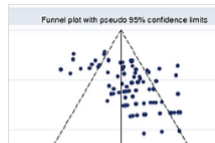


**Forest plot of sensitivity**
Steve Goodacre ⌄                30/12/2011
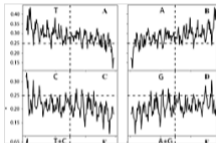


**Linkage disequilibrium (D') plot of IPF1 gene in Caucasians**
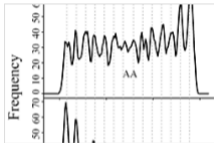Mohammad A Karim ⌄                30/12/2011



**Funnel plot for sensitivity**
Steve Goodacre ⌄                30/12/2011



**Base composition plot in core region of mixture alignment**
Ji-Ping Z. Wang ⌄                30/12/2011
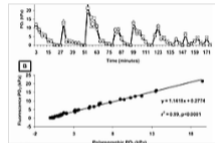


**Frequency plot of TT and AA signals in the alignment presented in**
Ji-Ping Z. Wang ⌄                30/12/2011
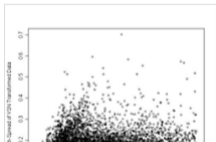


**L'Abbe plot of risk of leak in single-layer vs**
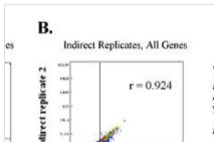Satoru Shikata ⌄                30/12/2011



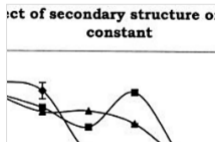**Plot of fluorescence, polarographic and predicted partial oxygen tensio...**
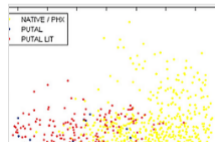Andrew D Shaw ⌄                30/12/2011



**Plot of the rank of the median probe**
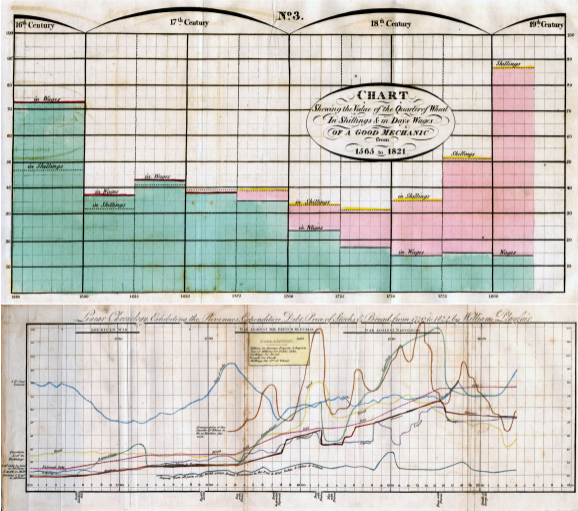


**(A) Scatter plot of direct labeling**



**A plot of the second-order rate**
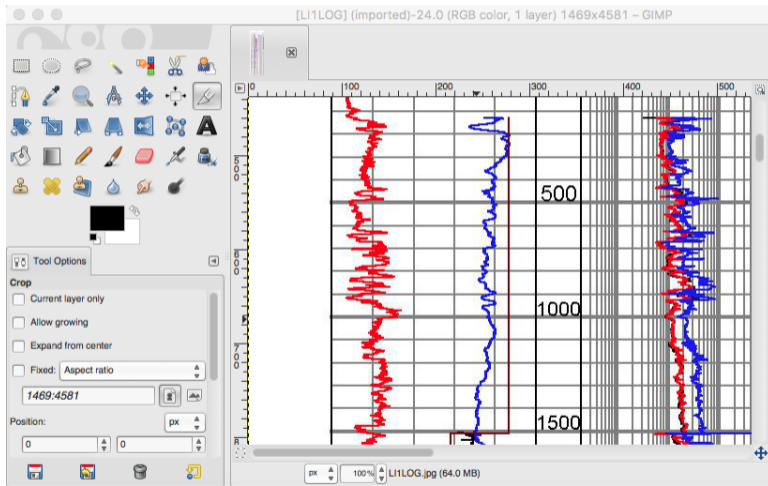


**Kernel-based scatter plot**
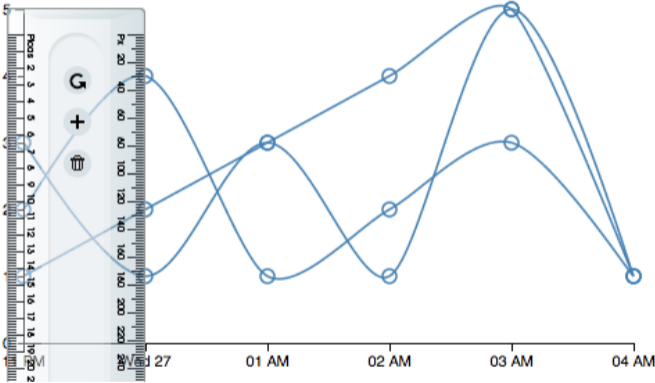
# Fetching Raw Data

Contact Authors





William Playfair, 1759-1823

# Fetching Raw Data



Pixel Counting?

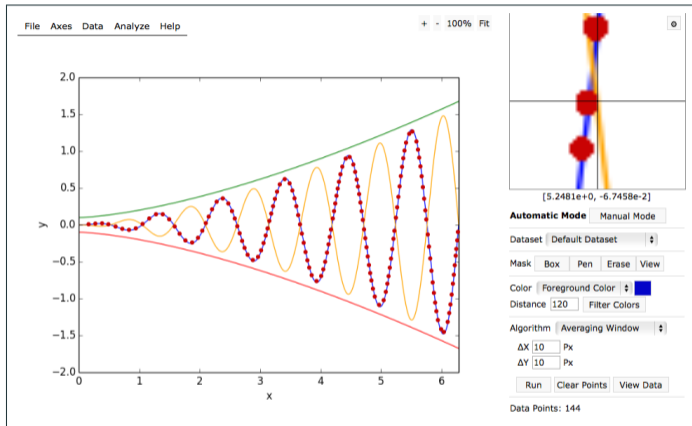# Fetching Raw Data



Geometry?

## Existing tools

- A few tools are available[1–6], but with many issues:
  - Difficult to access or incompatible with the operating system
  - Support only XY charts
  - Complicated interface
  - Accuracy concerns
  - Some are expensive, but not significantly better
  - Minimal automation

## Existing tools

- A few tools are available[1–6], but with many issues:
    - Difficult to access or incompatible with the operating system
    - Support only XY charts
    - Complicated interface
    - Accuracy concerns
    - Some are expensive, but not significantly better
    - Minimal automation
- Complete automation is still an area of active research[7–10].

# WebPlotDigitizer



- Free, opensource, web based tool
- Works with a wide variety of charts
- Partial automation with sub-pixel resolution algorithms

# Demonstration

Also available on http://arohatgi.info/WebPlotDigitizer

# Workflow

# XY Charts



## Affine Transformation



DIFFERENTIAL SCALING

SKEW

ROTATION

TRANSLATION

### X and Y Axes Calibration

Enter X-values of the two points clicked on X-axis and
Y-values of the two points clicked on Y-axes

| | Point 1 | Point 2 | Log Scale |
|---|---|---|---|
| X-Axis: | 1 | 6 | |
| Y-Axis: | -2 | 2 | |

*For dates, use yyyy/mm/dd format (e.g. 2013/10/23 or
2013/10). For exponents, enter values as 1e-3 for 10^-3.

OK

# Bar Charts and Histograms



## Bar Chart

## Histogram*

*Calibrate as a 2D XY plot

**Bar Chart Calibration**

Enter the values at the two points selected on the continuous axes along the bars

| Point 1 | Point 2 | Log Scale |
|---------|---------|-----------|
| 0 | 100 | ☐ |

OK

# Polar Diagrams

# Ternary Diagrams

# Scaled Images (Maps, Microscope, etc.)



Apple A7 (ifixit.com)
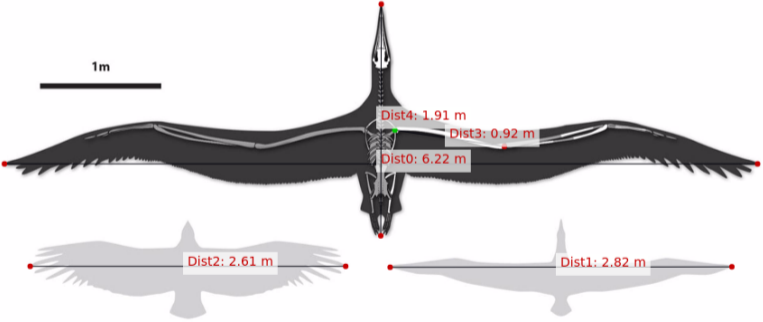


Arado AR 65F (the-blueprints.com)
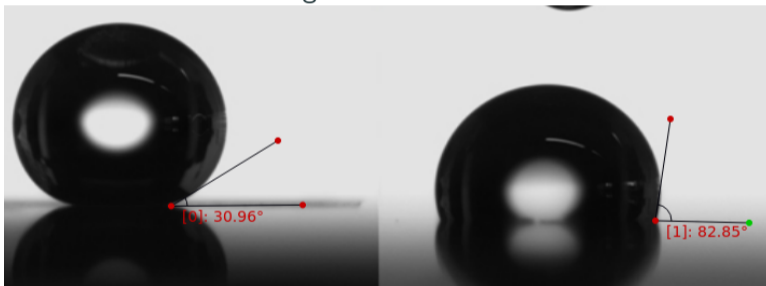
# Measurement Tools

File    Axes    Data    **Measure**    Help

**Distances**

Angles

2.0



## Distance Measurement

1m

Dist4: 1.91 m

Dist3: 0.92 m

Dist0: 6.22 m

Dist2: 2.61 m

Dist1: 2.82 m

*Pelagornis sandersi* (National Geographic)

# Measurement Tools



Angle Measurement

# Auto-Extraction



Challenges:

- Color and shape based image segmentation
- Region of interest identification
- Sub-pixel thinning, centroid estimation

# Averaging Window Algorithm



Suited for continuous curves and data points

# X Step with Interpolation Algorithm



Suited for continuous or discontinuous curves, data points and noisy data

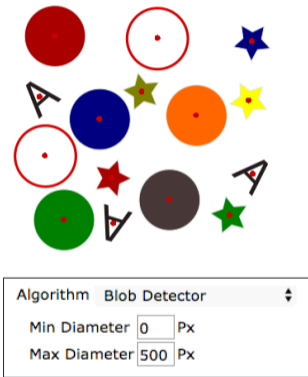# X Step with Interpolation Algorithm



Dashed Lines

Noisy Data

# Bar Charts and Histogram Algorithms



Suited for vertical or horizontal bar charts and histograms

# Blob Extraction



- Connected component labeling
- Computes:
  - Centroid
  - Area
  - Moment Invariant
- Shape based extraction

# Grid Removal

# Accuracy: Summary

## Accuracy: Independent Studies



[11]



[12]



[13]

# Interesting Use Cases



**Mark Brandon** ✔
@icey_mark

**Follow**

Said it before Web Plot Digitizer is the just
the most amazingly useful tool for extracting
data from papers
arohatgi.info/WebPlotDigitiz…

# Interesting Use Cases



Cindy Harnett
@CindyHarnett

Liking webplotdigitizer for grabbing data from pics in #lab, thanks @CousinAmygdala for the tip.
arohatgi.info/WebPlotDigitiz…

# Interesting Use Cases



Thurston Sexton. "Optimal Modeling of Knots in Wood". Arizona State University, 2015

# R Package (Under Development)



https://github.com/ankitrohatgi/digitizeR

# Native Desktop App (Under Development)



Qt/C++

# Public Issue Tracker

- http://arohatgi.info/WebPlotDigitizer
- http://www.github.com/ankitrohatgi/WebPlotDigitizer
- ankitrohatgi@hotmail.com
- ankit_rohatgi

# References

[1]     *PlotDigitizer*. URL: http://plotdigitizer.sourceforge.net/.

[2]     *DigitizeIt*. URL: http://www.digitizeit.de/.

[3]     *GetData graph digitizer*. URL: http://www.getdata-graph-digitizer.com/.
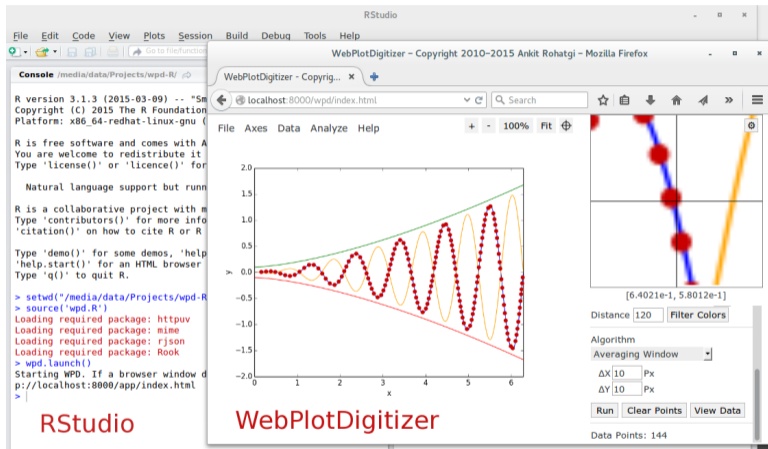
[4]     *Engauge Digitizer*. URL: http://markummitchell.github.io/engauge-digitizer/.

[5]     Biosoft. *Ungraph*. URL: http://www.biosoft.com/w/ungraph.htm.

[6]     Geomatix. *XYit*. URL: http://www.geomatix.net/xyit/.

[7]     Gonzalo Gabriel Méndez, Miguel A. Nacenta, and Sebastien Vandenheste. "iVoLVER". In: *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems - CHI 2016*. Association for Computing Machinery (ACM), 2016. DOI: 10.1145/2858036.2858435.

[8]     Noah Siegel et al. "FigureSeer: Parsing Result-Figures in Research Papers". In: *Computer Vision – ECCV 2016*. Springer Nature, 2016, pp. 664–680. DOI: 10.1007/978-3-319-46478-7_41.

# References

[9] Daekyoung Jung et al. "ChartSense: Interactive Data Extraction from Chart Images". In: ACM, May 2017. URL: https://www.microsoft.com/en-us/research/publication/chartsense-interactive-data-extraction-chart-images/.

[10] Sagnik Ray Choudhury and Clyde Lee Giles. "An Architecture for Information Extraction from Figures in Digital Libraries". In: *Proceedings of the 24th International Conference on World Wide Web - WWW 2015 Companion*. Association for Computing Machinery (ACM), 2015. DOI: 10.1145/2740908.2741712.

[11] Brittany U. Burda et al. "Estimating data from figures with a Web-based program: Considerations for a systematic review". In: *Research Synthesis Methods* (2017). DOI: 10.1002/jrsm.1232.

[12] Daniel Drevon, Sophie R. Fursa, and Allura L. Malcolm. "Intercoder Reliability and Validity of WebPlotDigitizer in Extracting Graphed Data". In: *Behav. Modif.* 41.2 (Mar. 2017), pp. 323–339. DOI: 10.1177/0145445516673998.

[13] M. Moeyaert, D. Maggin, and J. Verkuilen. "Reliability, Validity, and Usability of Data Extraction Programs for Single-Case Research Designs". In: *Behav. Modif.* 40.6 (Apr. 2016), pp. 874–900. DOI: 10.1177/0145445516645763.

[14] Thurston Sexton. "Optimal Modeling of Knots in Wood". Arizona State University, 2015.