# Visualizing Billions of Data Points: Doing It Right

Alexey Zaytsev,

Skoltech, CDISE

18 January

Slides are based on http://go.continuum.io/visualizing-billions-data-points/

**Skoltech**
Skolkovo Institute of Science and Technology

# Overview

1. Visualizing big data — what is the problem?

2. Datashading

3. Datasets:

    - 10 million points of NYC Taxi data

    - 3 billion points of OpenStreetMap data

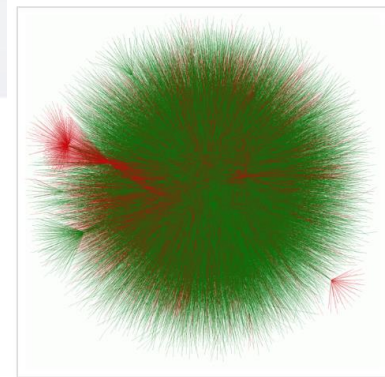    - 300 million points of US Census data

**Skoltech**
Skolkovo Institute of Science and Technology

# Billions and billions…
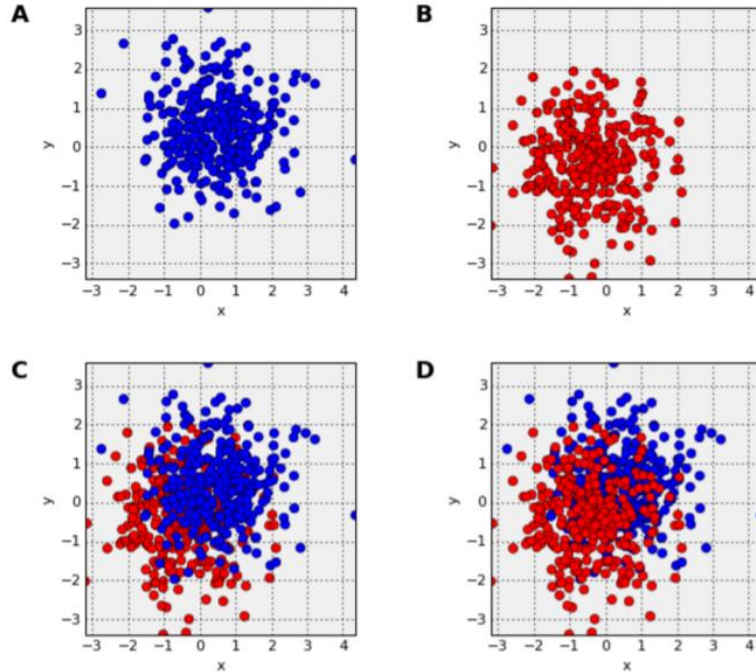
# Big data magnifies small problems

Big data presents storage and computation problems

More importantly, standard plotting tools have problems that are magnified

by big data:

- Overdrawing/Overplotting
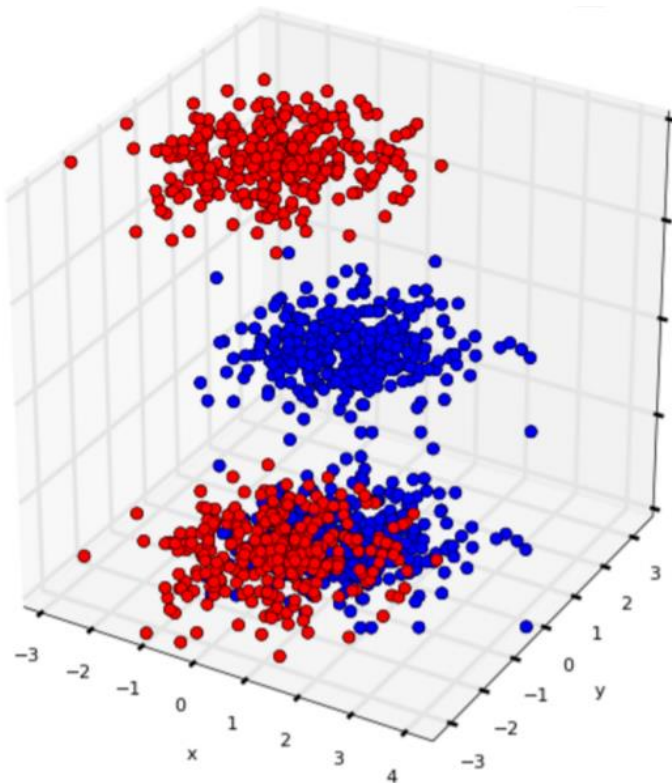
- Saturation

- Undersaturation

- Binning issues

We'll first explain these problems, and then present *datashading*

technique to address them.

**Skoltech**
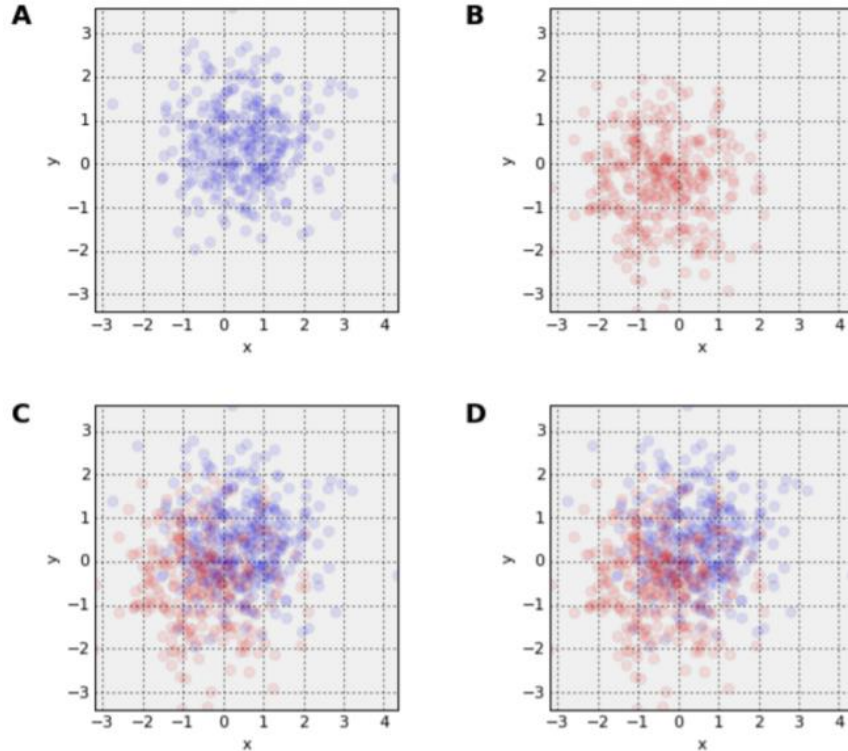Skolkovo Institute of Science and Technology

# Overdrawning



- For a scatterplot, the order in which points are drawn is very important
- The same distribution can look entirely different depending on plotting order
- Last data plotted overplots

Skoltech
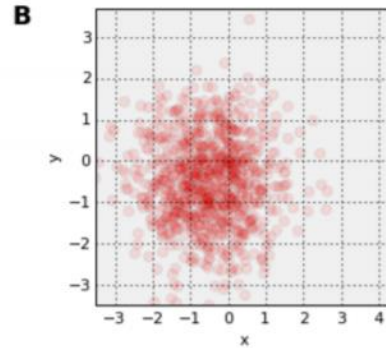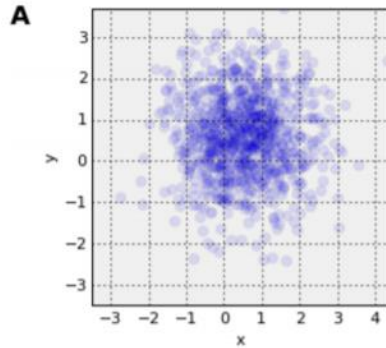Skolkovo Institute of Science and Technology

# Overdrawning



- Underlying issue is just occlusion
- Same problem happens with one category, but less obvious
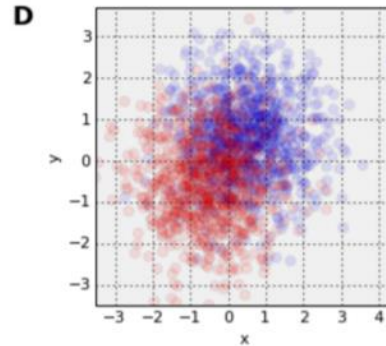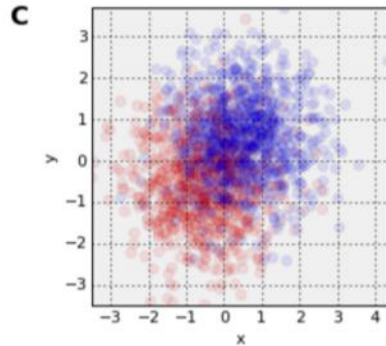- a Can prevent occlusion using transparency

# Saturation



- For alpha = 0.1, up to 10 points can overlap before saturating the available brightness
- Now the order of plotting matters less
- After 10 points, first-plotted data still lost
- For one category, 10, 20, or 2000 points overlapping will look identical

**Skoltech**
Skolkovo Institute of Science and Technology
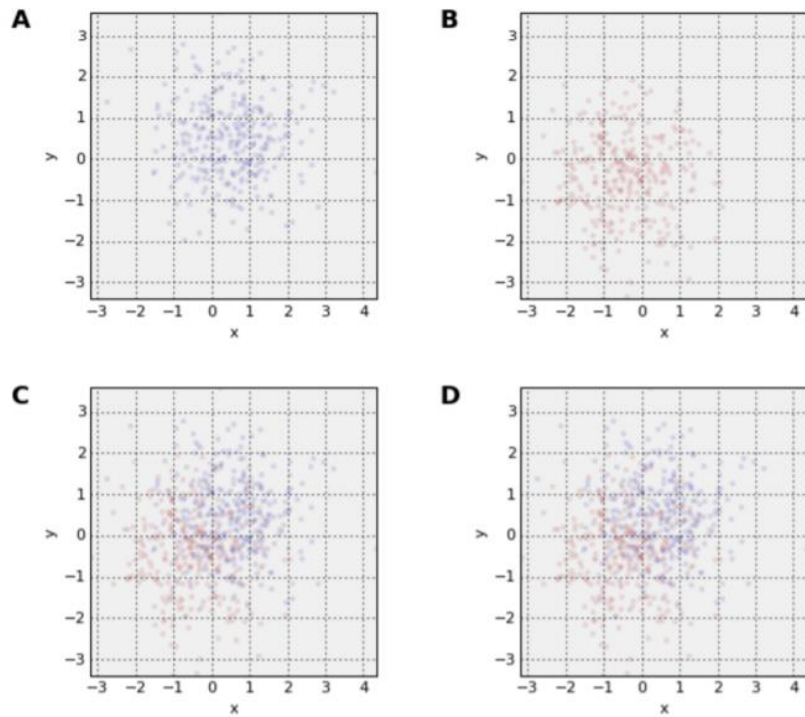
# Saturation



Same alpha value, more points:

- Now is highly misleading

- alpha value depends on size, overlap of dataset

- Difficult-to-set parameter, hard to know when data is misrepresented

**Skoltech**
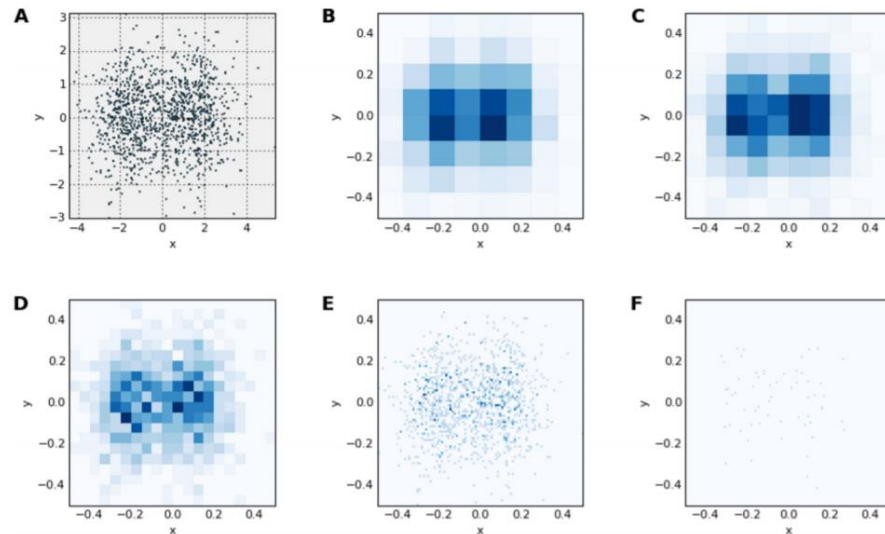Skolkovo Institute of Science and Technology

# Saturation



Can try to reduce point size to reduce overplotting and saturation

- Now points are hard to see, with no guarantee of avoiding problems

- Another difficult-to-set parameter

- For really big data, scatterplots start to become very inefficient: many points per pixel — may as well be binning by pixel

Skoltech
Skolkovo Institute of Science and Technology

# Binning

- Can use heatmap instead of scatter
- Avoids saturation by auto-ranging on bins
- Result independent of data size
- Here two merged normal distributions look very different at different binning
- Another difficult-to-set parameter

**Skoltech**
Skolkovo Institute of Science and Technology

# Main challenges while plotting big data

- When exploring really big data, the visualization is all you have —
  there's no way to look at each of the individual data points
- Common plotting problems can lead to completely incorrect
  conclusions based on misleading visualizations
- Slow processing makes trial and error approach ineffective

When data is large, you don't know when the plot is
lying.

**Skoltech**
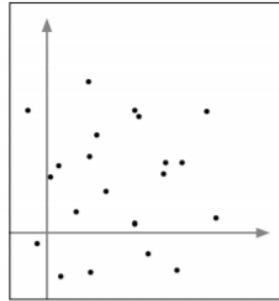Skolkovo Institute of Science and Technology

## Datashading

- Flexible, configurable pipeline for automatic plotting

- Prevents overplotting, saturation, and undersaturation •

- Mitigates binning issues by interactivity, even of very large datasets on ordinary machines

- Allows rapid iteration of visual styles & configs, interactive selections and filtering, to support data exploration

**Skoltech**
Skolkovo Institute of Science and Technology

# Datashading pipeline starts with projection



Data → Project / Synthesize → Scene

- Stage 1: select variables (columns) to project onto the screen
- Data often filtered at this stage

Skoltech
Skolkovo Institute of Science and Technology

# Datashading pipeline continues with aggregation



Data → Project / Synthesize → Scene → Sample / Raster → Aggregates

- Stage 2: Aggregate data into a fixed set of bins

- Each bin yields one or more scalars (total count, mean, stddev, etc.)

**Skoltech**
Skolkovo Institute of Science and Technology

# Datashading pipeline ends with transfer



- Stage 3: Transform data using one or more transfer functions, culminating in a function that yields a visible image
- Each stage can be replaced and configured separately

**Skoltech**
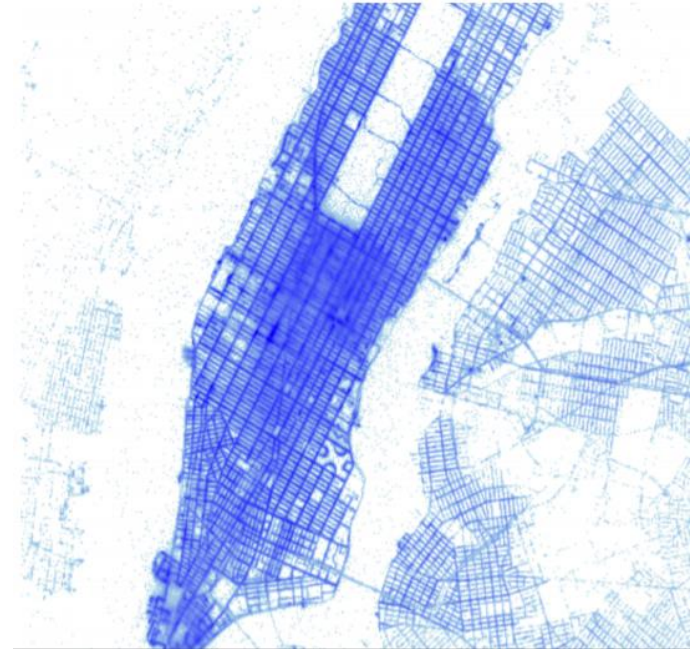Skolkovo Institute of Science and Technology

# NYC taxi data

- Data for 10 million New York City taxi trips

- Even 100,000 points gets slow for scatterplot

- Parameters usually need adjusting for every zoom

- True relationships within data not visible in std plot

Datashading automatically reveals the entire dataset, including outliers, hot spots, and missing data

**Skoltech**
Skolkovo Institute of Science and Technology

# Billion OSM points

Open Street Map data:

- About 3 billion GPS coordinates

- https://blog.openstreetmap.org/
  2012/04/01/bulk-gps-point-data/

- This image was rendered in one
  minute on a standard MacBook with 16 GB RAM

- Renders in 7 seconds on a 128GB Amazon EC2
  instance

**Skoltech**
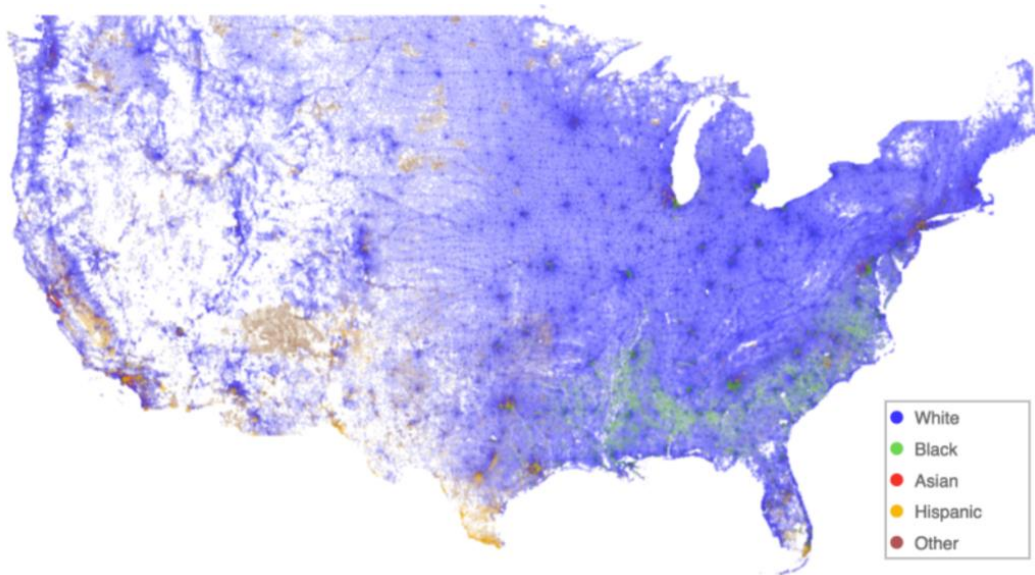Skolkovo Institute of Science and Technology

# Billion OSM points

- OSM database is 40GB, larger than the 16GB RAM on this machine
- Fast out-of-core operation powered by: •
  - Numba: generates fast C and GPU code from Python source
  - Dask: Parallelizes tasks
- datashader source code is all Python

**Skoltech**
Skolkovo Institute of Science and Technology

# 300 million points Census data

- One point per person

- 300 million total

- Categorized by race

- Datashading shows faithful

  distribution per pixel



Legend:
- White
- Black
- Asian
- Hispanic
- Other

**Skoltech**
Skolkovo Institute of Science and Technology

# Example of data exploration with datashader

http://datashader.org/topics/nyc_taxi.html

Skoltech
Skolkovo Institute of Science and Technology