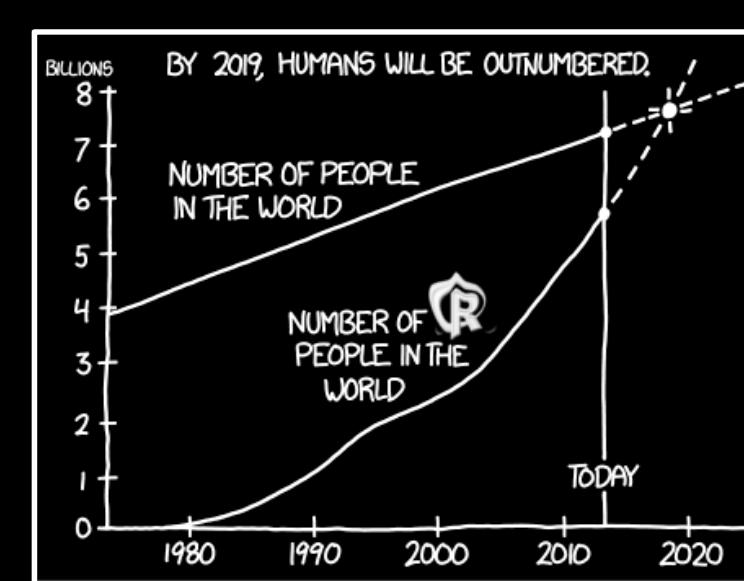


OUR USERS ALL AROUND THE WORLD

GERGELY DAROCZI

RAPPORTER.NET



OVERVIEW

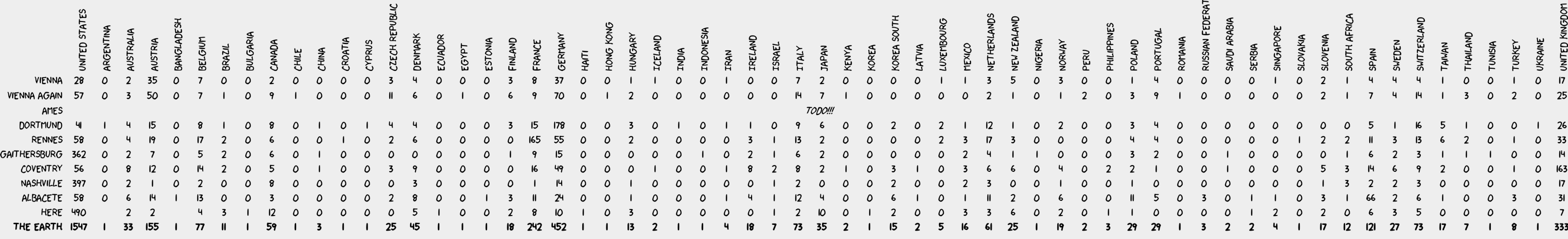
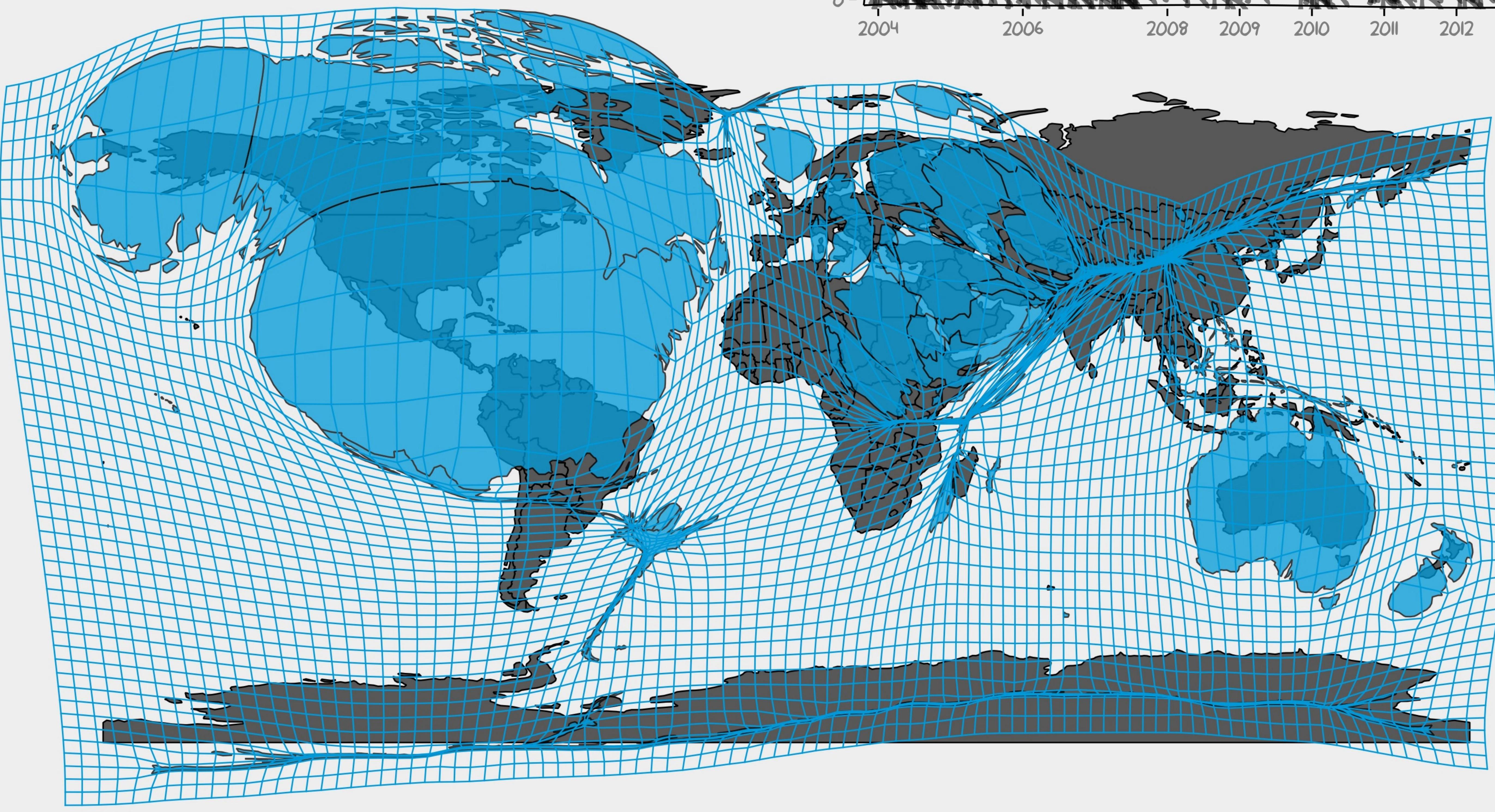
This poster was inspired by a recent blog post on „Where is the R activity?” and a few follow-up articles, where the authors tried to determine the **number of R users all around the world**.

The first version of my related experiment included various data sources on the number of **R Foundation members**, attendee-lists of the **useR! conferences**, the number and size of **R User Groups** registered on meetup.com, **R package downloads** from Rstudio cloud-CRAN servers, R-related **search queries via Google** and the geographical distribution of **GitHub users** with at least one R repository. The results and a combined R score were published on an interactive D3.js application.

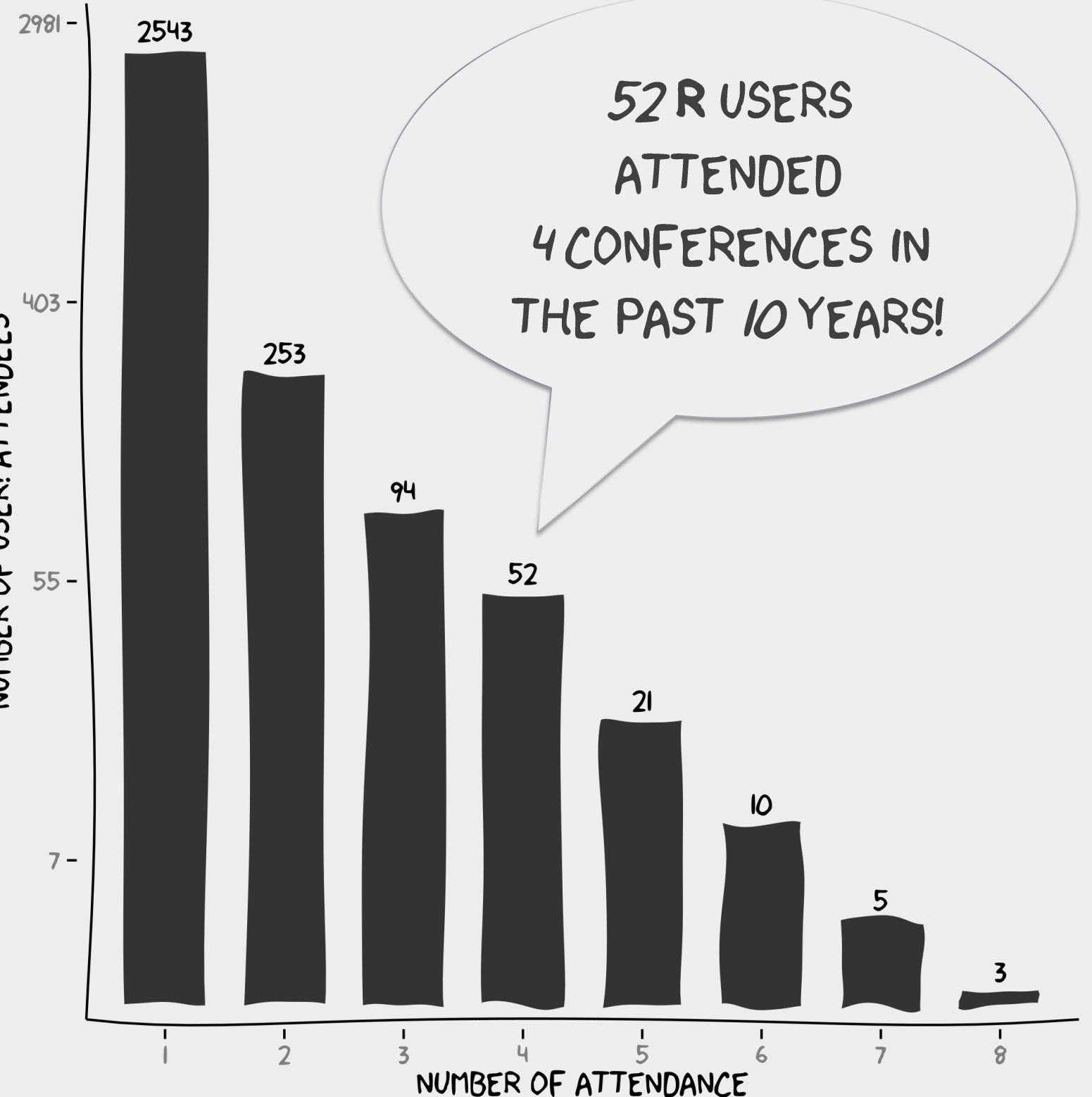
Now we focus on the attendees of the useR! conferences along with the fresh dataset of **600+ registrants** here (aka *useR! 2014*). First, we visualized the number of participants in time and space, then identified the number of attendance on a useR basis, so that we could **estimate** the size of the unkown population with **loglinear models**. How many guys and girls are using R? 2 or 3M? Even more?

Although this poster was not be able to answer this question after all, it still features some related statistical methods along with a few **XKCD** webcomic-styled plots of the results. Further analysis was also done on the proportion of genders, professional titles and headlines of the attendees etc.

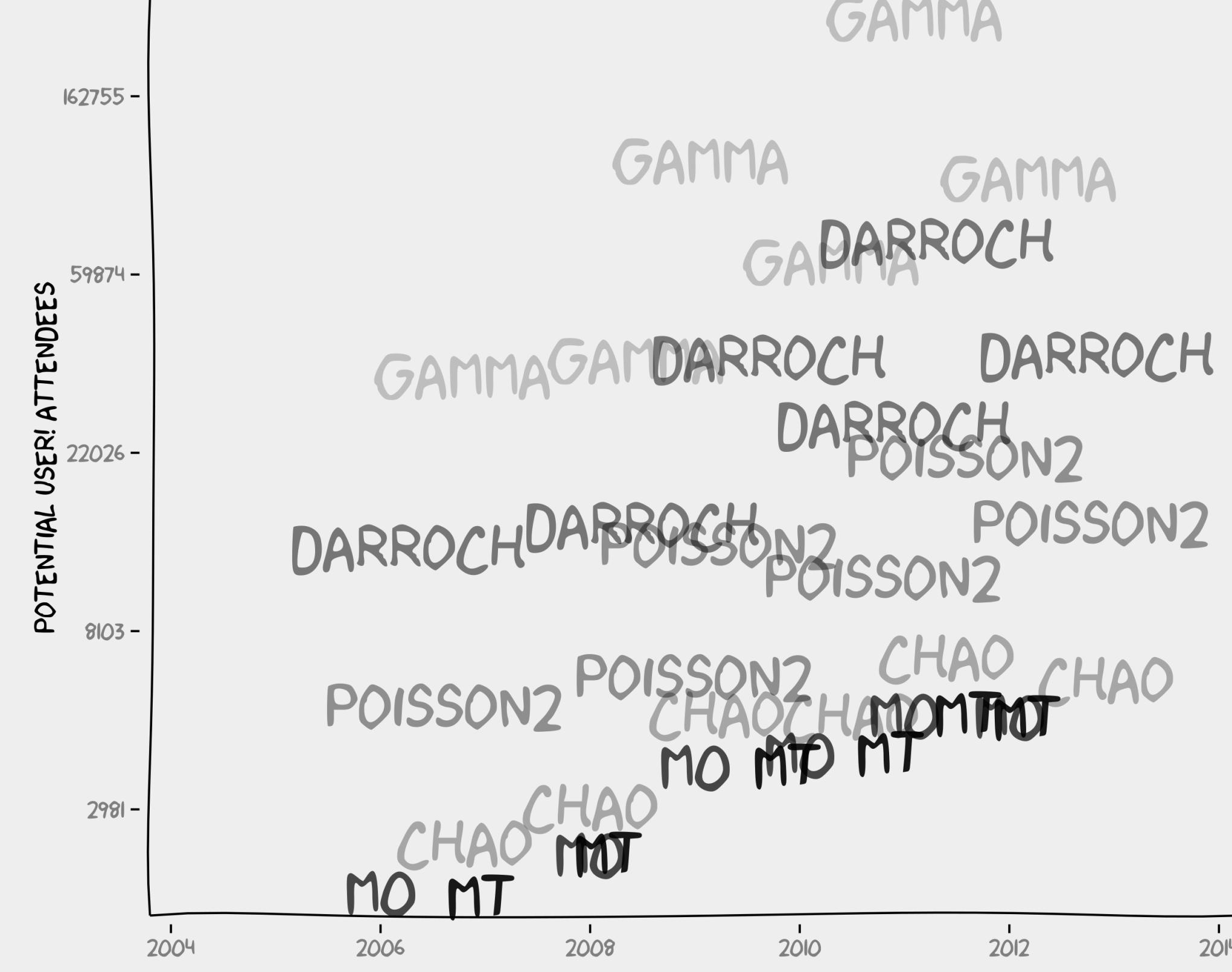
OVERALL NUMBER OF ATTENDEES FOR ALL USER! CONFERENCES BETWEEN 2004 AND 2014 – 10 YEARS!



MOST STATISTICIANS CAN READ, ANALYSE AND INTERPRET HUGE TABLES EASIER AND QUICKER COMPARED TO REALIZING WHAT'S GOING ON IN A CARTOGRAM!



1. The name of all **useR! 2004-2014** attendees was compiled in a list, where some (left figure) **recurring participants** were identified after manual data cleaning.
 2. These results can be used as the "number of units captured *i* time" in a "mark and recapture" research msethod for **estimating the unknown population size**.
 3. Building open-population models by estimating the **observation probability**, the **survival rate** and the **number of new arrivals** based on historical data.
 4. Building 6x6 **closed population** models based on 3-years long time intervals (right figure): the *lower bound* estimate of Chao showed an increase from 2,500 to 7,000 even with and without a time effect, while Daroch's model estimates tend to spread between 15,000 and 70,000. The Gamma and Poisson models of Rivest and Baillargeon (2012) returns some restrained (10,000-20,000) and some over-



REFERENCES

- James Cheshire (2013). Where is the R Activity? <http://spatial.ly>
 - Gergely Daroczi (2013). The attendants of useR! 2013 around the world. <http://blog.rapporter.net>
 - Gergely Daroczi (2014). R activity around the world. <http://blog.rapporter.net>
 - Dominique Andrieu, Christian Kaiser and André Ourednik: ScapeToad. v1.1. <http://scapetoad.choros.ch>
 - Gastner, M.T. and Newman, M. (2004). Diffusion-based method for producing density equalizing maps. PNAS 101(20): 7499-7504.
 - Sophie Baillargeon, Louis-Paul Rivest (2012). Rcapture: Loglinear Models for Capture-Recapture Experiments. R pkg v1.3-1. CRAN.
 - Emilio Torres Manzanera (2014). xkcd: Plotting ggplot2 graphics in a XKCD style. R pkg v0.0.3. CRAN

LIMITATIONS

- UseR! attendees do not represent R users.
 - Further data cleaning is required.
 - Missing data for 2007.
 - Chinese R Conference is not included.
 - Closed population model is not appropriate.
 - Improved model selection is desired.
 - We need more & independent data sources

