

# Statement of Purpose of Max Klein

## For School of Information, UC Berkeley

### Research Agenda

Gender, race, sexuality, nationality, social-class, native-language, height and weight. Can we make an exhaustive list of social biases? I believe we already have, but haven't realized it.

Social biases are a main barrier to entry and expansion in sociotechnical systems<sup>1</sup>, but it's not certain we've accounted for them all. I'm tantalized by the idea that we could automatically detect biases we haven't yet identified – that we could uncover Rumsfeldian unknown unknowns. I propose a research agenda to classify the already-known social biases by as they appear in the collaborative technologies, and then to search for unidentified biases using those classifications. As an explanatory example, create a statistical model of how the known skewed distributions gender, race, and nationality exist in Wikidata (the free knowledge base that feeds Wikipedia), and then inspect *all* the property distributions for properties that match the biased patterns. The project grows more complex by allowing property-pairs (e.g. gender by race), different sociotechnical communities (e.g. Freebase, Git-hub), and different models of bias (e.g. editorship-measures). If successful we will find overlooked prejudices of people using a technology.

### Preparations

This research agenda is the direct consequence of the experience I have gathered in my academic, industry, and personal portfolio. They are the three strands which constitute the braid of my readiness: technical grokability, Open Culture social awareness, and self-driven curiosity.

In Spring of 2015, I will present my research on "The Virtuous Circle of Wikipedia"<sup>2</sup> at Computer-Supported Collaborative Work conference 2015. The paper provides a new definition and measure for the “collaborativeness” of a socio-technical community – the degree to which being a good user is correlated with editing with other good users. In this case I studied the non-profit Wikipedia, using economic insight and a variant of the Google PageRank algorithm. This will guide my proposed research agenda because by determining the connection between biases and collaborativeness, we address whether found biases reflect society at large, or only the biases of the editing community

Moreover my “Virtuous Circle” research has prepared me by stretching my technical strengths and weaknesses. The project was a collaboration with UC Berkeley iSchool PhD Thomas Maillart, who pushed me to grasp new mathematical models and higher methodological rigor. My BA in mathematics was useful in learning the required network science quickly demonstrated by the fact that an early stage of the research was presented as a poster to the NetSci 2014 conference<sup>3</sup>. Likewise, I was also pushed to come to speed with a base computational social science literature. An example of exceeding this goal is evident in the “Method of Reflections” technique which we borrowed and utilized in the paper – after which I created the first open-source implementation of the method<sup>4</sup>.

To open source that algorithm was instinctual because I am a part of the academic-hacker Wiki Research community. It was there that my blog posts on the gender biases in different Wikipedia languages was picked up by Hanyang University Sociology Professor Piotr Konieczny. We came to work together on creating "Wikipedia Gender Inequality Index" (WIGI<sup>5</sup>), an upcoming Open Dataset of extracted information from Biography articles across all Wikipedia Languages. My contribution was to put together all technical infrastructure to re-index, analyze and display the dataset each month. Moreover I sharpened my social research thinking from Piotr when we conjured the notion to aggregate our place of birth, citizenship, and

---

<sup>1</sup> Halfaker, Geiger, Morgan, and Reidl. [The Rise and Decline of an Open Collaboration System](#) (2013)

<sup>2</sup> Klein and Maillart [The Virtuous Circle of Wikipedia](#) (2014)

<sup>3</sup> Klein and Maillart [Poster - The Virtuous Circle of Wikipedia](#) (2014)

<sup>4</sup> Klein [Method of Reflections Explained and Exemplified](#) (2014)

<sup>5</sup> Klein and Konieczny [WIGI: Wikipedia Gender Inequality Index](#)

ethnicity data into “World Cultures” categories using Mechanical Turk. Currently in the submission process, WIGI is an existing prototype and first step in the stream of the research agenda put forth.

My preparedness is anchored by a final personal factor: the philosophy to steer my own course and navigating by natural curiosity. One telling fact about the above two projects is that I have not done them under the direction of an institution. Both have come voluntarily without pay and on my free time. The last major directives I received were 3 years ago at the beginning of my stint as Research Assistant at OCLC Inc. where I was hired to improve Wikipedia-Library integration, which I accomplished by way of writing Wikipedia bots to add content from Library databases. After publishing about the process<sup>6</sup> and finally amassing over 2 million edits, I came to see what was driving me in the position. Strangely, it was the *Open Notebook Science* ideology – having the source, data and analysis synced publicly online as it evolves in front of me, thereby submitting my ideas to the full scrutiny of the internet hivemind. More and more my ideas spilled over the scope of the OCLC Research blog<sup>7</sup>, and I started hosting them myself at *notconfusing.com*<sup>8</sup>. My natural inclinations to start investigating on my own started receiving attention and energizing feedback – particularly on posts exploring the data of gender<sup>9</sup> and language<sup>10 11</sup>. Although I didn't fully realize it then, this new networked, boss-free but still peer-reviewed world meant that I already started to pursue my the kind of research that I now propose to formalize.

## Faculty Interest

No blog is an island. As my research efforts intensify I have a growing need for guidance in framing and methodological focus. In a conversation with Coye Cheshire – the faculty member with whom I'm interested in working with, I was aroused by the discussion of his research around online effects gender<sup>12</sup>, race<sup>13</sup>. Most fascinatingly is how they can be hypothesized to comment on trust and inequality in our society<sup>14</sup>. Those are precisely the larger themes I am dancing around. I would ask that he could put his hands on my shoulders and move me to the exact location of the grander questions. For the agenda outlined here I may think my worry is “how to classify bias?” or “what is the extent of prejudice online?” The real question is likely more poignant and sociological, only I do not know how to ask it yet.

Earning a doctorate degree from iSchool in the grand scheme is a waypoint on the path towards becoming a professional researcher. The end product of my education is necessarily full time research, because the end goal of my education is knowing how we are effecting sociotechnical systems with outmoded social biases, which will forever be less than fully understood.

---

<sup>6</sup> Klein and Kyrios [VIAFbot and the Integration of Library Data on Wikipedia](#) (2013)

<sup>7</sup> Klein [HangingTogether.org Blog Posts](#) (2012-14)

<sup>8</sup> Klein [Notconfusing.com](#) Research Blog Posts (2012-14)

<sup>9</sup> Klein [Sex Ratios in Wikidata Part III](#) (2014)

<sup>10</sup> Klein [Actionable Metrics for Uganda and Côte D'Ivoire](#) (2014)

<sup>11</sup> Klein [The Most Unique Wikipedias According to Wikidata](#) (2013)

<sup>12</sup> Antin, Yee, and Cheshire [Gender Differences in Wikipedia Editing](#) (2011)

<sup>13</sup> Mendelson, Shaw, Fiore, and Cheshire [Black/White Dating Online](#) (2014)

<sup>14</sup> Cook and Cheshire [Social Exchange, Power, and Inequality in Networks](#) (2013)