# Text Classification for Qualitative Analysis at Scale
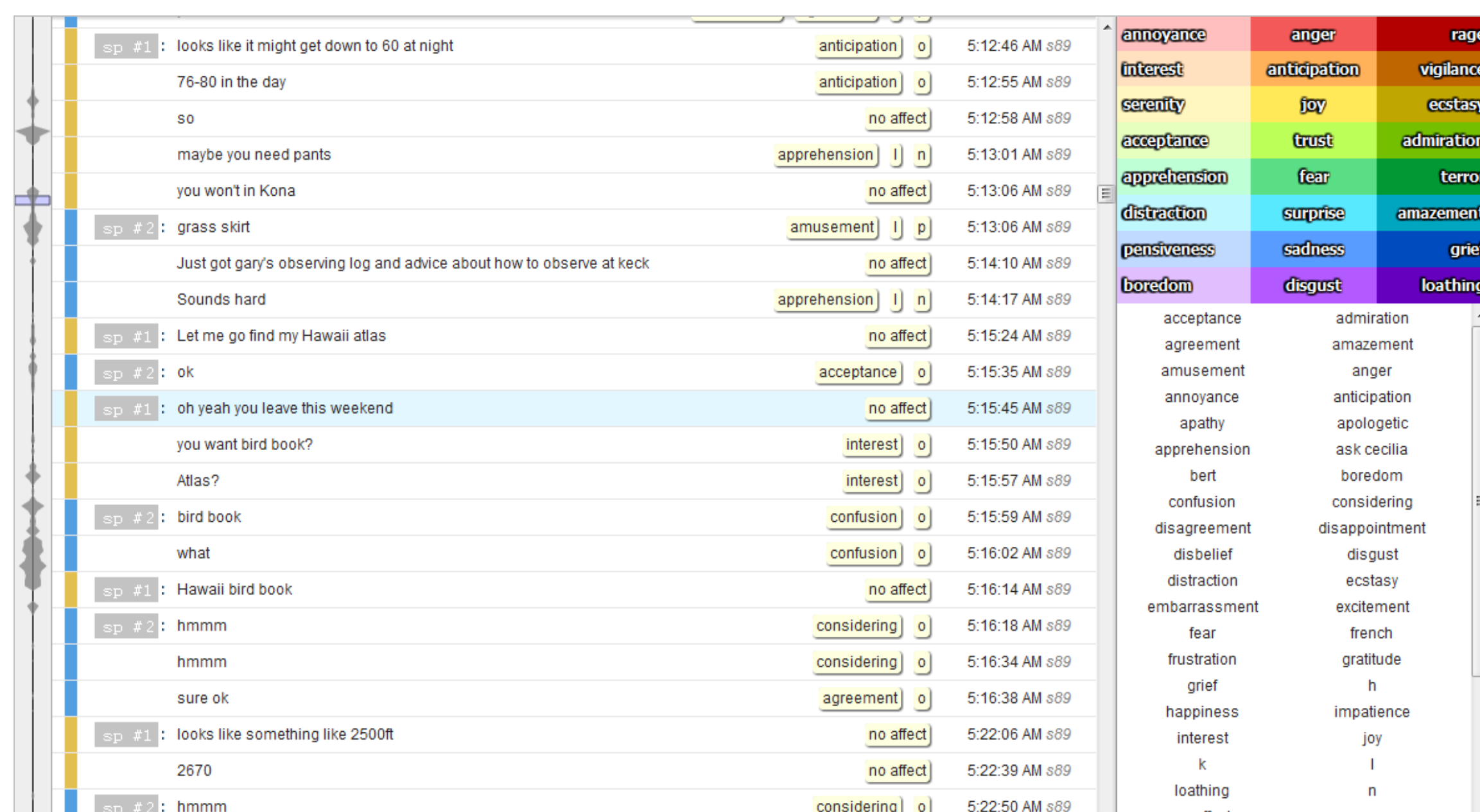
Michael Brooks, Katie Kuksenok, Megan Kelly Torkildson, John J. Robinson, Ray Hong, Taylor J. Scott, Daniel Perry, Cecilia Aragon

University of Washington

{mjbrooks, kuksenok, mtorkild, soco, rayhong, omni, dbperry, aragon}@uw.edu

## Human Centered Text Analysis

How can computational analysis tools like **machine learning** and **visualization** support **qualitative analysis** of large text datasets?

The wealth of data generated online through social media, forums, and email is creating new areas and paradigms for the social sciences. However, obtaining and working with datasets at this scale requires particular skills and resources. We need **better tools for broader participation** across research communities. We are using a human-centered approach to design and build software supporting large-scale qualitative research.
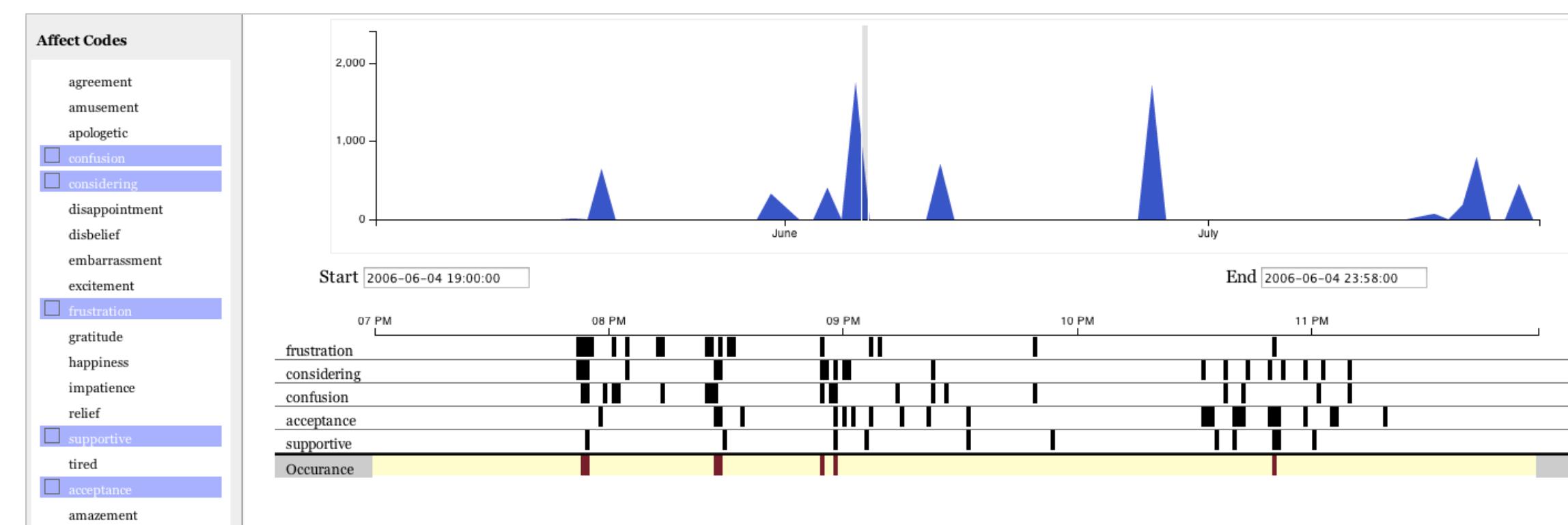


*A tool we created for coding chat messages more efficiently. Visualizing message frequency over time shows the pace of conversations. Color-coding speakers reveals the conversation thread. The grid of codes helps users remain cognizant of the range of available affect categories.*

## Qualitative Data Analysis

- Many approaches to qualitative research exist.
- Place emphasis on **context, meaning, and interpretation**.
- Explicitly acknowledge and reflect on the **role of the researcher**.

Researchers often combine qualitative and quantitative approaches, but qualitative methods are built on different foundations from quantitative methods. This has implications for design.

Qualitative analysis relies on reading, coding, and summarizing text, **time-consuming** but critical meaning-making activities. Qualitative researchers are often limited to small, focused datasets. We believe that machine learning and visualization can magnify or extend the power of qualitative analysis.



*A prototype visualization showing frequency of affect codes over time within our chat dataset. Visualizations such as this can show patterns in the categories which researchers could explore further.*

## Computational Social Science

- Social media and online communication tools produce vast amounts of text data: Twitter, email, online forums, and instant messaging.
- New research bridging social science and computer science is thriving on this public record of human social activity.
- How is activity online linked with offline phenomena such as community well-being and depression?
- How does information spread and propagates online?

While the research questions in this area are social, the methods are usually computational or statistical. We are building **analytics software grounded in qualitative methods**, to make large-scale social data more useful for research where a qualitative/interpretive approach is appropriate.

## Text Classification

- *Text classification* has been applied to important problems including sentiment analysis, topic detection, and spam filtering.
- Decomposes raw text documents into numeric **features** such as the presence or absence of specific words and phrases.
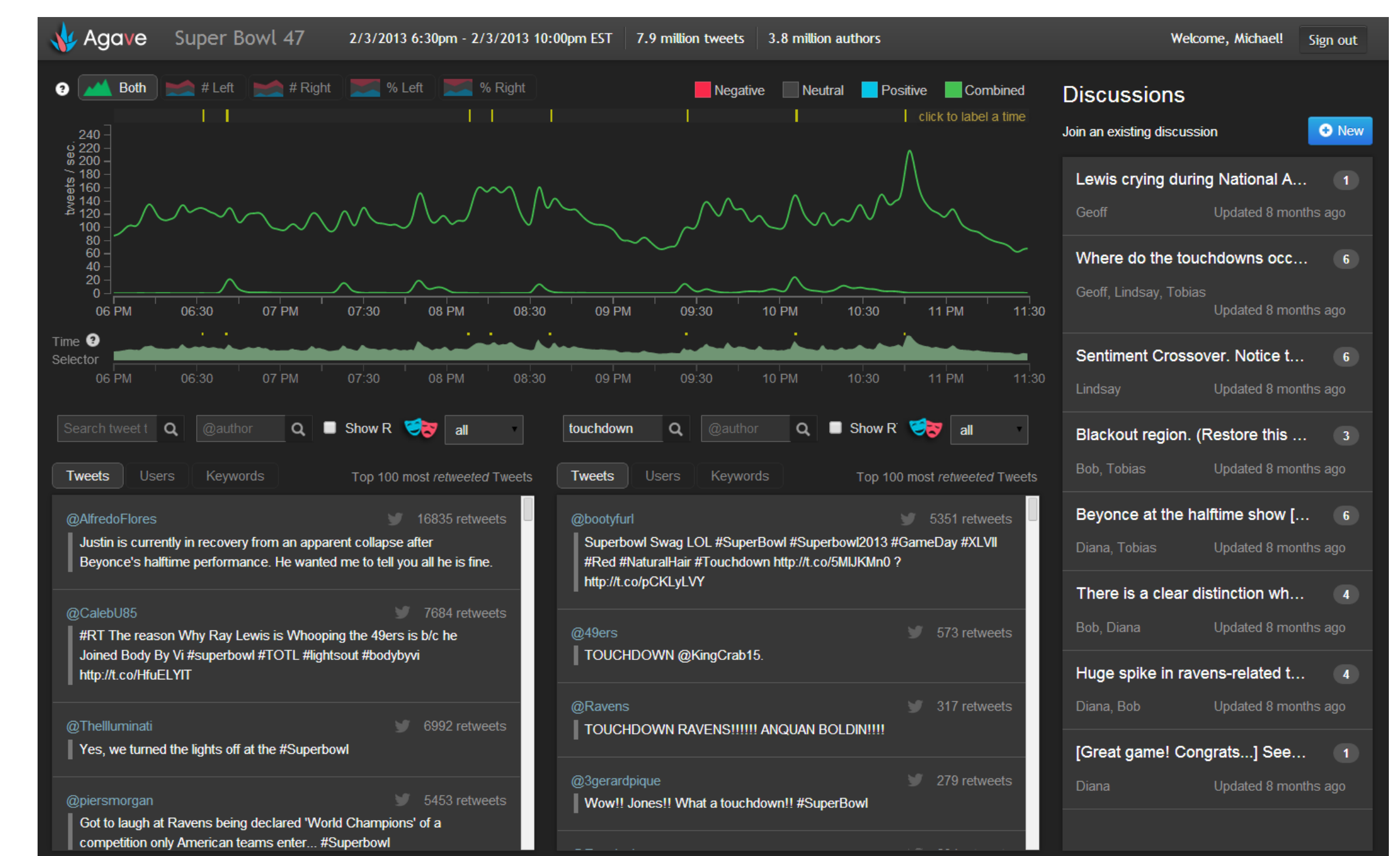- **Classifiers** can be trained to map these feature values onto a set of document classes.

Trained classifiers are typically evaluated and optimized based on **quantitative metrics** such as accuracy, precision, or recall. This requires comparing the classifier's output to some pre-classified example documents, or ground **truth data**. It is not clear whether these approaches to evaluation are appropriate and sufficient when applied as part of qualitative analysis.

## Visual Analytics

- Visual analytics is the use of interactive visualizations and computation for **analytical reasoning and data exploration**.
- Leverage the strengths of both humans and computers in a complementary fashion.
- Interactive visualizations could help researchers understand text-based datasets, feature spaces, and classifier performance.

Large, text communication datasets are typically complex, multifaceted, and time-varying. Interactive visualizations can be helpful for exploring the relationships between people, the mood and content of conversation, as well as dynamics over time.

From the initial open-coding phase of qualitative research, users may wish to begin experimenting with machine learning, whether to find new data to examine or to extend their coding work over more data. Designing for an interactive and exploratory style of use will enable non-experts to effectively apply machine learning to their own research problems. Guided by their intuition and domain knowledge, users will be able to rapidly create and evaluate classifiers with the support of visualizations and example results.



*Agave, a system we built for collaborative exploration of large Twitter datasets. The screenshot shows a user inspecting tweets related to touchdowns in the 2013 Super Bowl, a 7 million tweet dataset. Tools that facilitate collaboration could be especially useful for research on large social data sets, which often involve multidisciplinary teams.*

Scientific Collaboration & Creativity Lab

HCDE Human Centered Design & Engineering

dub

UNIVERSITY of WASHINGTON