**Working toward effective anonymization for surveillance data:**

**Innovation at South Africa's Agincourt Health and Demographic Surveillance Site**

- Lori M. Hunter, Wayne Twine and Catherine Talbot

**Extended Abstract**

Research on population-environment connections is often hampered by lack of the necessary geographic coding required to link people to place.

To facilitate research on socio-ecological intersections in low-income settings, this project offers innovations in data integration through a focus on one of 50+ "Health and Demographic Surveillance Systems" providing unparalleled social science data in hard-to-reach world regions. This paper reports on 3 goals.
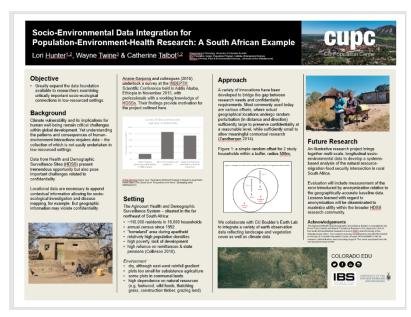
First, we **anonymize** household-scale data in our rural South African study setting, the Agincourt Health and Demographic Surveillance site. Integration of local environmental information with household data requires household locational information – but such information could breach confidentiality. Anonymization involves small random perturbations of locational data to protect household confidentiality.

Using clusters of household defined by differing rules, we then **integrate** earth observational information to measure local environmental conditions and temporal change.

Third, the **empirical effects** of different approaches to clustering will be explored to determine those which best balance research and confidentiality requirements. We will examine the intersections between migration, natural resource availability and food security using natural resource measures based on different clustering approaches.

In the long-run, the goal is to develop effective and efficient options for the provision of geocodes to researchers for socioecological data integration within Agincourt and potentially within other research settings.

We have already done a poster presentation of an early version of this project at CU Population Center's mini-conference on Africa Population-Health-Environment in April 2018. Since then, we have generated different clustering approaches and are implementing those now. The clusters will then be linked with earth observational data, measuring natural resources, which we already have in-hand. Results across models of migration, natural resources and household well-being (namely, food security) will be contrasted and provide the focus of this presentation.

Examples of three different clustering techniques:

1) clustering by village sections,



2) with the addition of roads as boundaries,



3) clusters based on distance from village edge.