# The DataBridge: A Social Network for Long Tail Science Data

**Howard Lander**

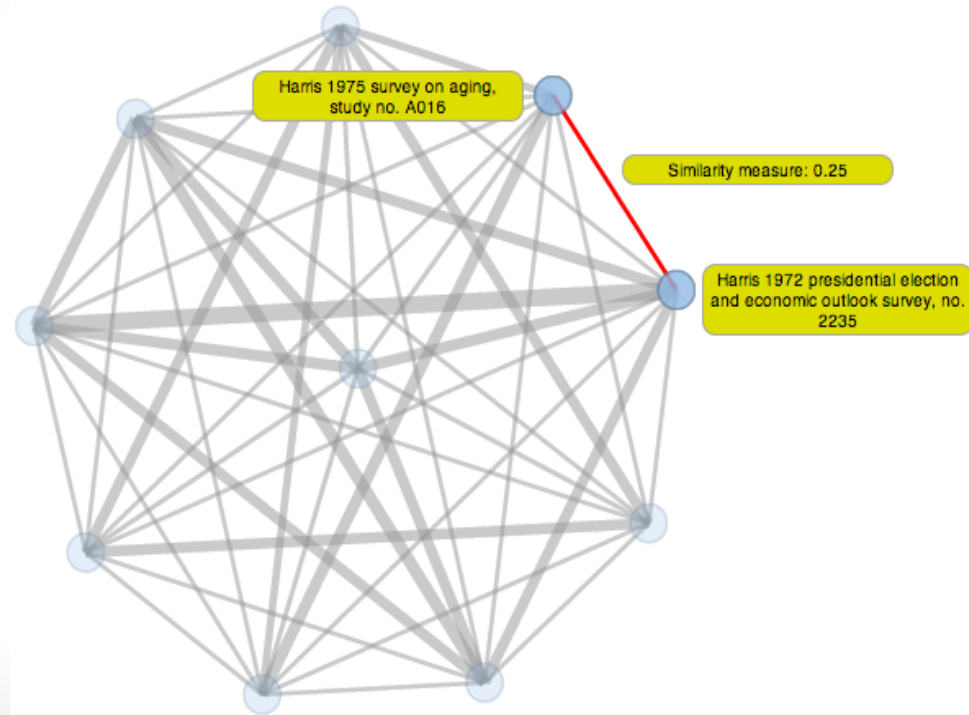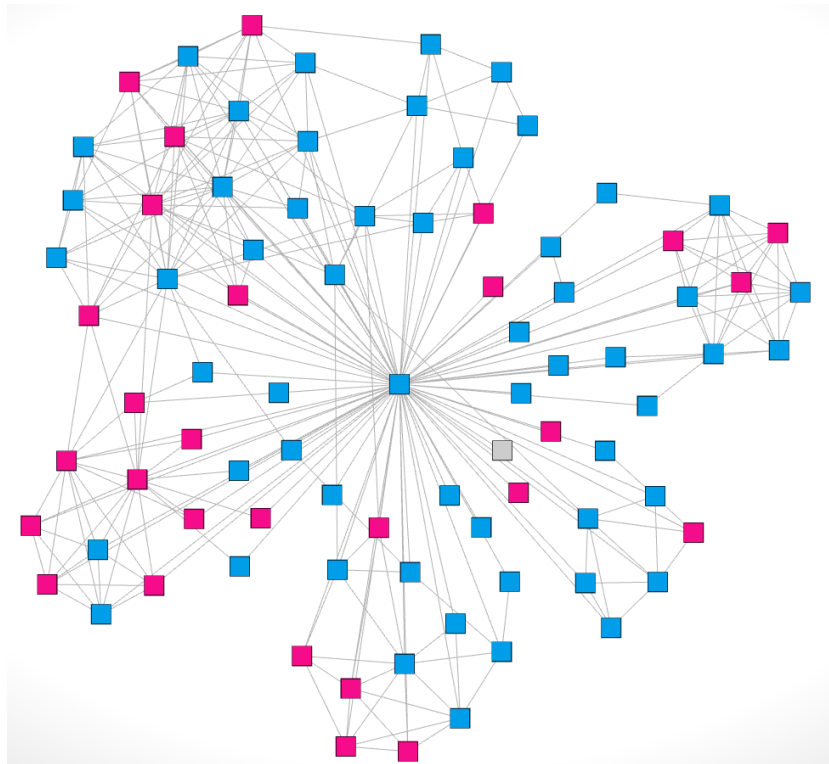**howard@renci.org**

**Renaissance Computing Institute**

**The University of North Carolina at Chapel Hill**

# The DataBridge: A Social Network for Data



Harris 1975 survey on aging, study no. A016

Similarity measure: 0.25

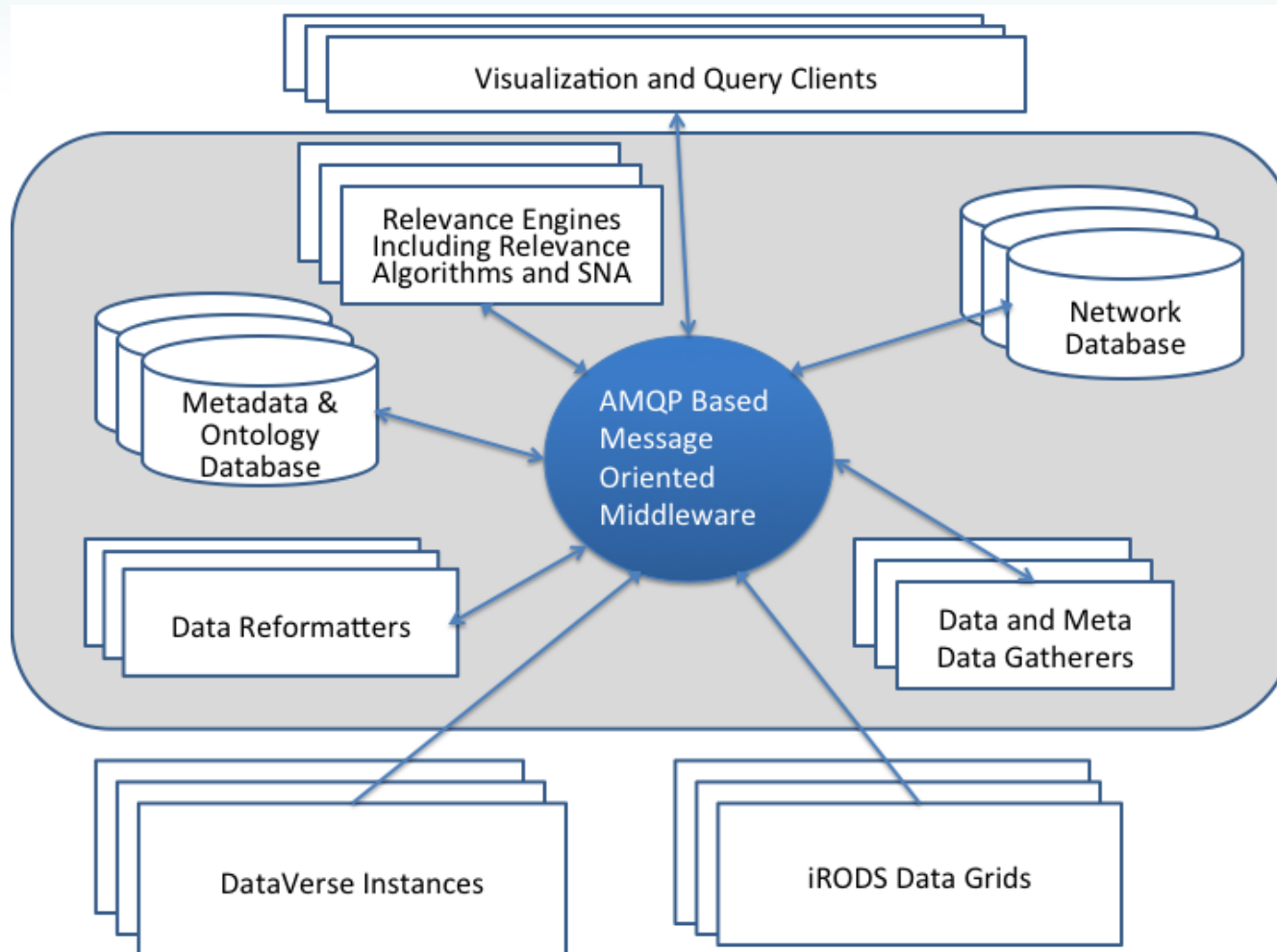Harris 1972 presidential election and economic outlook survey, no. 2235

# Dark Data from The Long Tail of Science

- **Long tail data is the small data sets produced by numerous investigators**
- **From Brahe to Mendel discovery has come from relatively small data sets**
- **Much long tail data is dark data, data "not easily found by potential users" (Heidorn)**
- **Long tail data sets lack structural advantages of "classic" Big Data.**
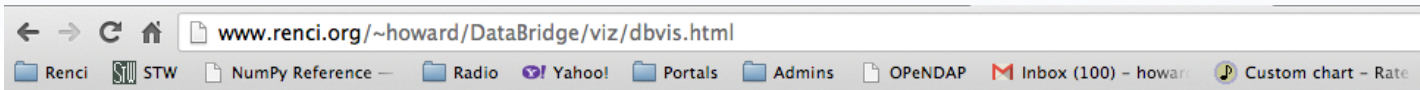
# The DataBridge Strategy: Building a Social Network for Scientific Data

- **Construct a multi-dimensional sociometric network for data. Three challenges:**
  - **Evaluate the similarity/relevancy of data sets**
  - **Perform community detection on the resulting set of similarities**
  - **Provide query interfaces on resulting multi-dimensional network**
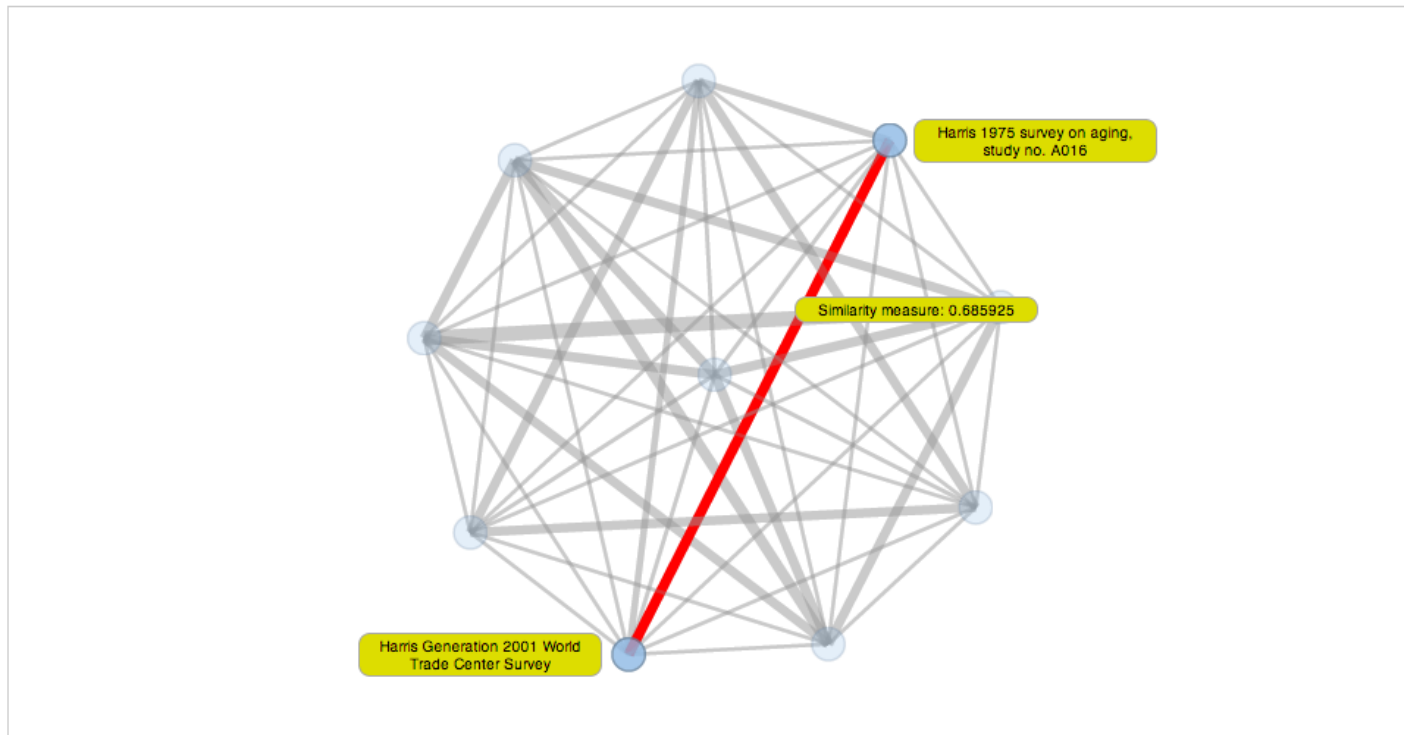
# DataBridge Implementation

# DataBridge Progress to Date: JavaScript based network visualization tool

# DataBridge Team

- **PI: Arcot Rajasekar RENCI and SILS, UNC-Chapel Hill**
- **Collaborators:**
  - **Odum Institute, UNC-Chapel Hill**
  - **Population Informatics Research Group, UNC-Chapel Hill, Texas A & M University**
  - **iLab, North Carolina A&T University**
  - **The Institute for Quantitative Social Science, Harvard University**
- **Funded by: NSF Office of Cyberinfrastructure Awards OCI-1247562, OCI-1247602 and OCI-1247663**