

# Analysis and Modeling of Demographic Cohorts in the Population of PhD Recipients in Science and Engineering



Antonio Sanfilippo – PNNL  
Katy Börner, James P. Crutchfield,  
Dowman P. Varn et al. – IU, UC Davis  
Marlene Lee – PRB

*Progress report, SWAM Program, January 26, 2012*

# Overview

- ▶ Research Aims
- ▶ Progress report (by team)
  - Data Access and Analysis
  - Modeling
  - Future work

# Research Aims

- ▶ Collect and analyze evidence relevant to demographic dynamics of the scientific workforce and related activities
- ▶ Use evidence collected to develop, calibrate and evaluate evidence-based dynamic models of population change in the scientific workforce
- ▶ Answer questions such as
  - What is the impact of the increasing proportion of foreign graduate students and postdoctoral scholars?
  - What are the proportions of female and minority students?
  - What are the effects of geolocation on a successful career?
  - How much should the rate of participation be increased to get to population proportions?

# PNNL Progress Report

- » Data Access and Analysis
  - Initial data analysis insights
  - Initial modeling efforts
  - Future work

PRB



$\Psi$  INDIANA UNIVERSITY

# Data Access and Analysis

- ▶ Received NSF approval for licensing restricted SDR, NSRCG and integrated SESTAT Data 1993–2008
- ▶ Developed knowledge discovery wiki for SWAM with 2003 and 2006 SDR public data
  - Use faceted search to create data subset and export these as MS Excel CVS files suitable for analysis and modeling

Swam Data 2006

Export SWAM 2006 Data Set

102 items

sorted by labels, then by... grouped as sorted

NIH 2006 11805

Gender Male

Race Asian, non-Hispanic ONLY

Employment Status Employed

Highest Degree Field Biological sciences

Survey Mode World Wide Web

NIH 2006 16525

Gender Male

Race Under-represented Minorities

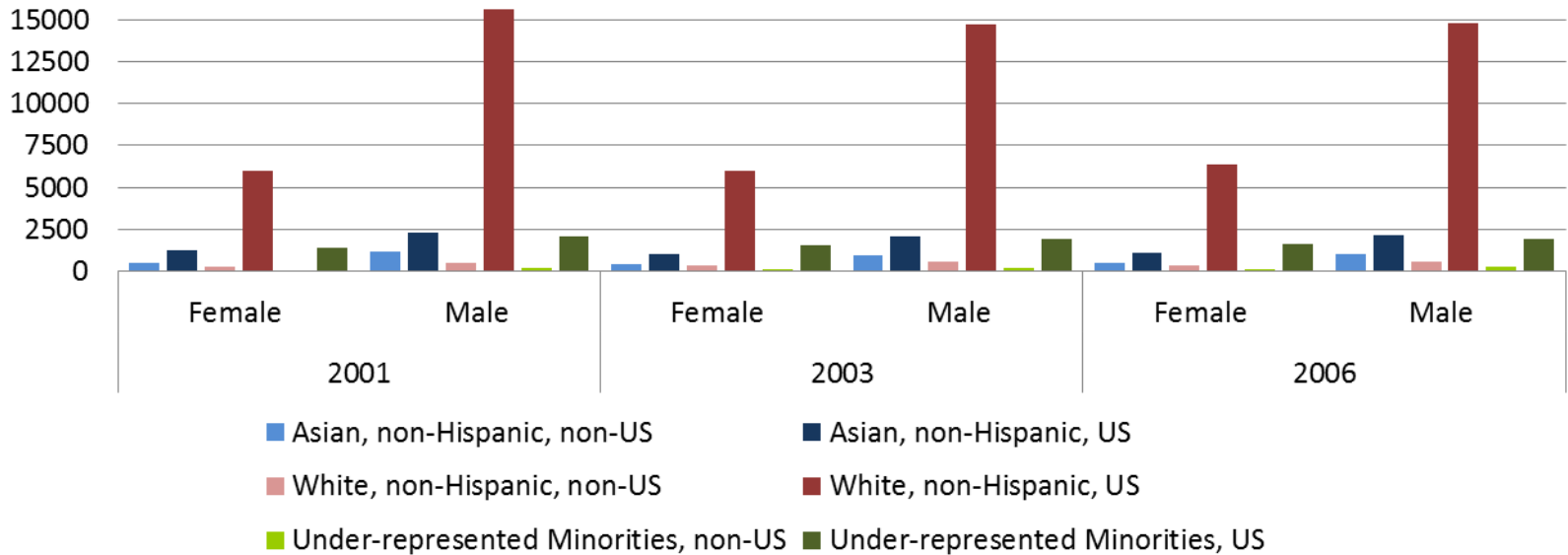
Employment Status Employed



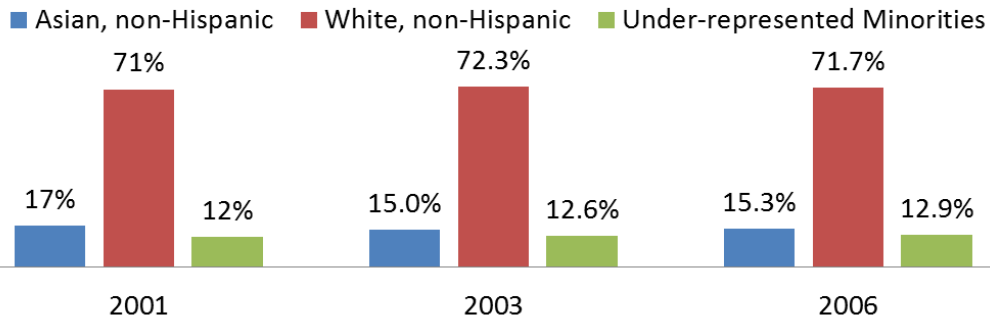
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	ACADADIF	ACADADMN	ACADANA	ACADOTHF	ACADPODC	ACADRCH	ACADTCHF	ACTCAP	ACTDED	ACTMGT	ACTRD	ACTRDT	ACTRES	ACTTCH	AGEP
2	No	No	No	No	No	No	Yes	No	No	No	No	Yes	No	Yes	53
3	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	49
4	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	50
5	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	55
6	No	No	Yes	No	No	No	No	No	No	No	No	Yes	No	Yes	58
7	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	53
8	No	No	No	No	No	No	Yes	No	No	No	Yes	Yes	Yes	Yes	51
9	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	51
10	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	No	No	No	No	No	53
11	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	50
12	Yes	No	No	No	No	No	No	No	No	No	No	Yes	No	Yes	63
13	No	No	No	No	No	Yes	Yes	No	No	No	Yes	Yes	Yes	Yes	55
14	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	Yes	Yes	Yes	Yes	No	No	55
15	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	No	No	No	No	No	67
16	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	Yes	No	Yes	Yes	No	No	56
17	No	Yes	No	No	No	No	No	No	No	Yes	No	No	No	No	52
18	No	No	No	No	No	No	Yes	No	No	No	No	Yes	No	Yes	56
19	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	55
20	Yes	No	No	No	No	No	No	Yes	No	No	No	No	No	Yes	52
21	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	No	Yes	Yes	Yes	No	51
22	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	55
23	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	No	Yes	Yes	Yes	No	53
24	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Yes	No	Yes	No	No	No	No	56
25	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Yes	Yes	No	Yes	Yes	Yes	Yes	54
26	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	38
27	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	No	Yes	Yes	Yes	No	41
28	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	Yes	Yes	Yes	No	36
29	No	No	No	No	No	Yes	Yes	No	No	No	Yes	Yes	Yes	Yes	38
30	No	No	No	No	No	Yes	No	No	No	No	No	Yes	No	Yes	37
31	No	No	No	No	No	Yes	No	No	No	No	Yes	Yes	Yes	Yes	36
32	No	Yes	No	No	No	No	Yes	No	No	Yes	Yes	Yes	Yes	No	51
33	No	Yes	No	No	No	No	Yes	No	No	Yes	Yes	Yes	Yes	No	59
34	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	57
35	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	No	Yes	No	No	No	No	57
36	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Yes	No	No	Yes	Yes	Yes	No	50
37	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	Yes	No	Yes	Yes	Yes	No	61
38	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	Logical Skip	No	Yes	Yes	Yes	Yes	No	No	54

# Initial data analysis insights

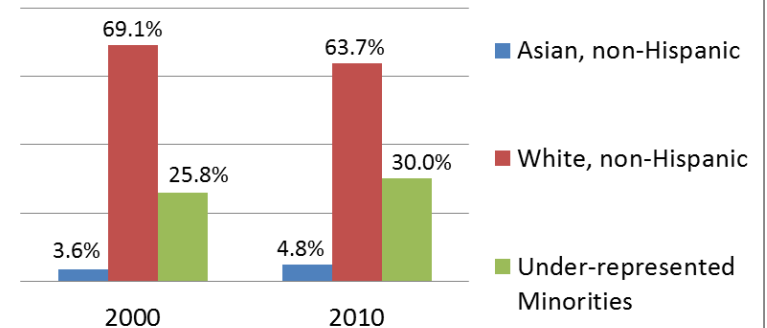
Number of PhDs by race, gender and citizenship status (Source: SDR 2001-2006)



Race percentages of total US PhDs  
Source: SDR 2001-2006



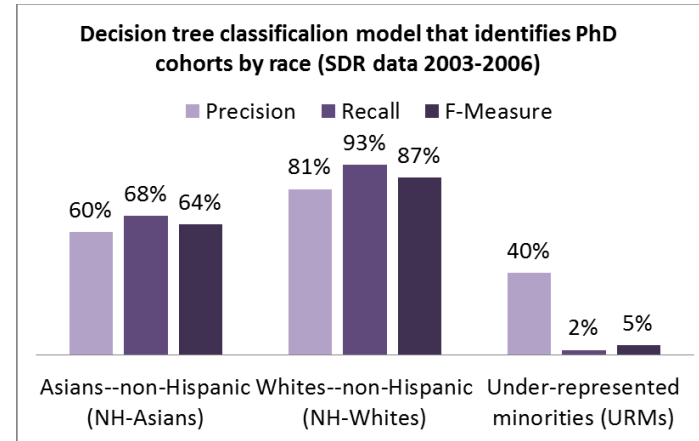
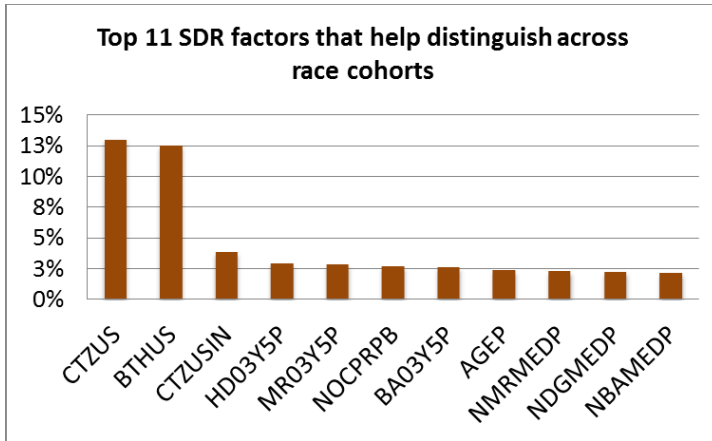
Percentage of total US population  
Source: US Census Bureau



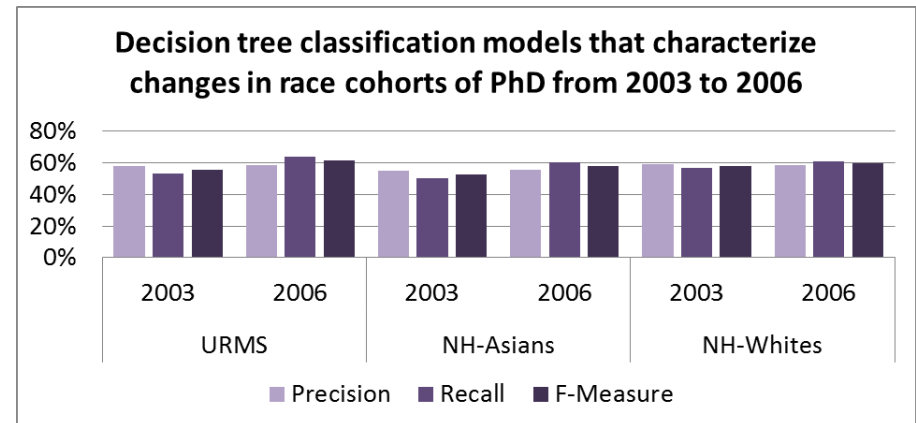
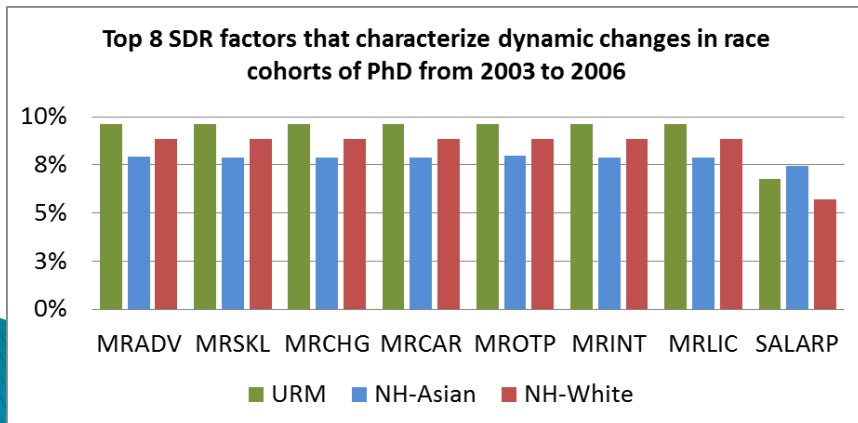


# Initial modeling efforts

- ▶ Developed and evaluated classification models that
  - Identify URM, NH-White and NH-Asian PhD cohorts



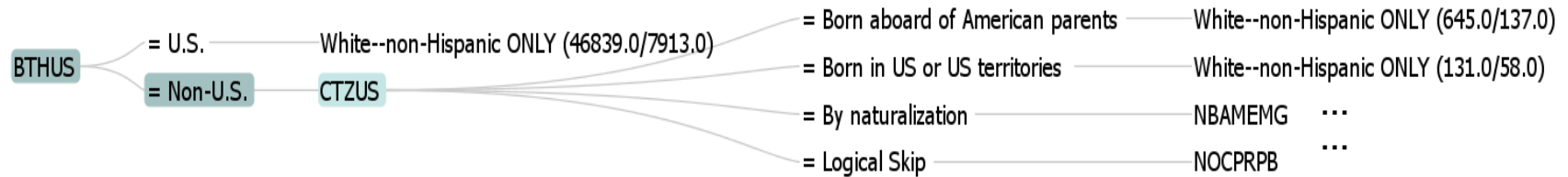
- Characterize dynamic changes across each race cohort



# Next Steps

## ▶ Modeling

- Improve performance, visualization and analysis of models



- Develop gender-based models
- Normalize analyses and models using census data (w/PRB)
- Start bringing science policy factors to bear on modeling task
- Develop dynamic models of population change in the scientific workforce
- Start working with restricted data

## ▶ SWAM Wiki

- Make available to SWAM members at [swam-us.org](http://swam-us.org)
- Load remaining public SESTAT data
- Improve usability and add visual analytic functionality
- Create standalone SWAM Wiki for restricted data





# Indiana University Progress Report

» Data Acquisition & Analysis  
Theory & Definitions  
Modeling Scholarly Dynamics  
Future work



Dowman P Varn,<sup>1,2</sup> Katy Börner,<sup>1</sup> Robert P Light,<sup>1</sup>  
Scott B Weingart,<sup>1</sup> & James P Crutchfield<sup>2</sup>

<sup>1</sup>*Cyberinfrastructure for Network Science Center  
School of Library & Information Science  
Indiana University*

<sup>2</sup>*Department of Physics & Complexity Sciences Center  
University of California, Davis*

# Data Acquisition & Analysis

- The National Science Foundation Survey of Doctoral Recipients (NSF–SDR) data contains information about recent PhDs in the sciences for the years 1993, 1995, 1997, 1999, 2001, 2003, 2006, & 2008. This longitudinal data give a rich, detailed picture of the scholars' career evolution. We are in the process of acquiring the complete surveys from NSF.
- In collaboration with Vincent Larivière of Université du Québec à Montréal, we now have bibliometric data from the Web of Science (WoS) for the 13,513 most moved physicists between 1980 and 1987, each of whom has published at least six papers. This sample contains 258,021 publications with a total of 4,120,342 citations and is expected to be rich in postdoctoral researchers. We can track their geolocation vs. time as well as any topical changes in their research interests.

# Theory & Definitions

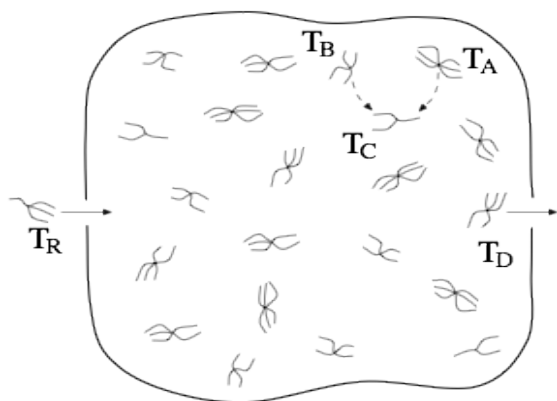
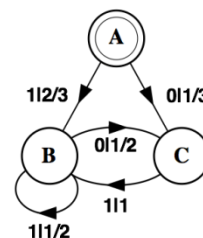
- ▶ Under what theoretical framework do we couch our questions? A necessary first task is to introduce mathematical definitions that parameterize our problem.
- ▶ **Scholar:** Let a scholar,  $\alpha^{(i,j,k)}$ , be an autonomous agent that can perform certain actions,  $i$ , (e.g., publish papers, work at institutions), has certain attributes,  $j$ , (e.g., gender, ethnicity), and can store and process information,  $k$ .
- ▶ **Event:** Let an event,  $\omega_m$ , be something that happens to an agent or its environment, or that the agent does to itself or its environment. Examples are the publication of a paper or movement to another location.
- ▶ **Career:** Let the career of a scholar,  $\Psi_\alpha$ , be the time ordered sequence of events that spans the active life of the scholar, i.e.,  $\Psi_\alpha = (\omega_1, \omega_2, \dots, \omega_M)$ .

# Modeling Scholarly Dynamics

- ▶ What parts of the dynamics of scholarly communities and of individual careers are endogenous and what parts are exogenous? That is, what is due to internal, spontaneously formed patterns inherent in the social population dynamics, and what is due to external funding, scholarly fashion, or other environmental factors? More pointedly, what in the observed dynamics of scholarly activity is amenable to change from external influences? What can policy makers really do?

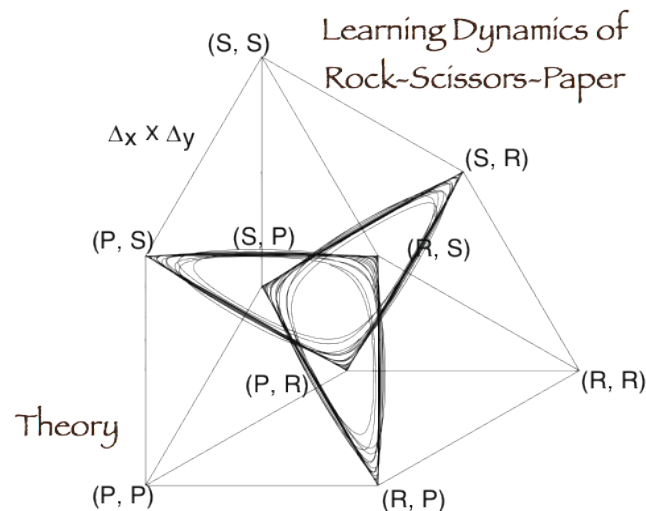
# Modeling Scholarly Dynamics II

- Recent work in game-theoretical dynamical systems (Sato & Crutchfield, 2002; Sato et. al. 2004) and population dynamics (Gornerup & Crutchfield, 2008) shows that a rich variety of endogenous behaviors, including collective adaption and spontaneous hierarchical organization, can form spontaneously in multiagent systems.
- We are extending these models to scholarly dynamics.



To the left, an interacting “soup” of scholars.

To the right, the complex dynamics of the rock-scissors-paper game.



# Next Steps

- ▶ Complete the theoretical (conceptual) framework necessary to parameterize the problem.
- ▶ Continue our analysis of the WoS data, in particular looking for patterns of changing geolocation and research interests and how these relate to productivity and success.
- ▶ Continue our preparations for the NSF–SDR survey data, and once acquired, begin extracting information about gender, ethnicity, movement, and topical interests of the participants.
- ▶ Develop models of active, interacting scholars. How well do these models predict publication patterns (both over topic space and geolocation space) and career trajectories? Do the parameters of this model differ for different genders and ethnic groups? What predictions do these models give when we introduce exogenous shocks to the population?



# PRB Progress Report

»» Data Access  
Analysis and Modeling  
Future work



# Status Update – Data Access

- ▶ Application for secure data in progress
- ▶ Computer and software acquisitions completed
- ▶ Testing new version of SPACE software for life tables with simulations

# Status Update – Analysis & Modeling

- ▶ Definitions of population groups
- ▶ Definition of employment states
- ▶ Functional form of hazard models for transitions

# Status Update – Next Steps

- ▶ Transition Matrices
- ▶ Life Table Simulations
- ▶ Involvement of collaborators (MPI)
- ▶ Research Assistant