

# Automatically assembling a census of an academic field

Allison Morgan  
University of Colorado, Boulder

with Samuel Way and Aaron Clauset

# About Me

Second year PhD Student  
in CS at CU Boulder

Collaborators and I study the  
“sociology of science”

Interested in computational  
methods to study under-  
representation in academia

## RESEARCH ARTICLE

### NETWORK SCIENCES

#### **Systematic inequality and hierarchy in faculty hiring networks**

**Aaron Clauset,<sup>1,2,3\*</sup> Samuel Arbesman,<sup>4</sup> Daniel B. Larremore<sup>5,6</sup>**

*Science Advances* 1(1), e1400005, 2015.

#### **Gender, Productivity, and Prestige in Computer Science Faculty Hiring Networks**

Samuel F. Way,<sup>1,\*</sup> Daniel B. Larremore,<sup>2,†</sup> and Aaron Clauset<sup>1,3,2,†</sup>

<sup>1</sup>*Department of Computer Science, University of Colorado, Boulder CO, 80309 USA*

<sup>2</sup>*Santa Fe Institute, Santa Fe NM, 87501 USA*

<sup>3</sup>*BioFrontiers Institute, University of Colorado, Boulder CO, 80303 USA*

*Proc. 25th Int'l World Wide Web Conf. (WWW), (2016)*

#### **The misleading narrative of the canonical faculty productivity trajectory**

**Samuel F. Way<sup>a,1</sup>, Allison C. Morgan<sup>a</sup>, Aaron Clauset<sup>a,b,c,2</sup>, and Daniel B. Larremore<sup>a,b,c,1,2</sup>**

<sup>a</sup>Department of Computer Science, University of Colorado, Boulder, CO 80309; <sup>b</sup>BioFrontiers Institute, University of Colorado, Boulder, CO 80303; and <sup>c</sup>Santa Fe Institute, Santa Fe, NM 87501

*Proceedings of the National Academy of Sciences* Oct 2017, 201702121

# Motivation

*Nobel Prize winners*



*Chemists*



*and those who leave academia*



Much of the sociology of science studies small samples of the academic workforce at a single point in time.

Can we build a tool to efficiently collect the employment information of **all faculty** across institutions, **across time**?

# Challenge



Jane  
Professor  
jane@example.edu



Mark  
Associate Professor  
mark@example.edu



Susan  
Assistant Professor  
susan@example.edu

Every department contains a public directory of its faculty

With the same information:  
names, titles, email addresses,  
and webpages

But, information is distributed  
and not well structured



# Our Approach

Department Homepage

Courses | Faculty ...

Identify the directory's  
HTML structure & extract  
faculty information



Jane  
Professor  
jane@example.edu

**faculty\_name:** Jane  
**title:** Professor  
**website:** ...  
**email:** ...



Mark  
Associate Professor  
mark@example.edu

Filter non-tenure-track  
faculty for further analyses



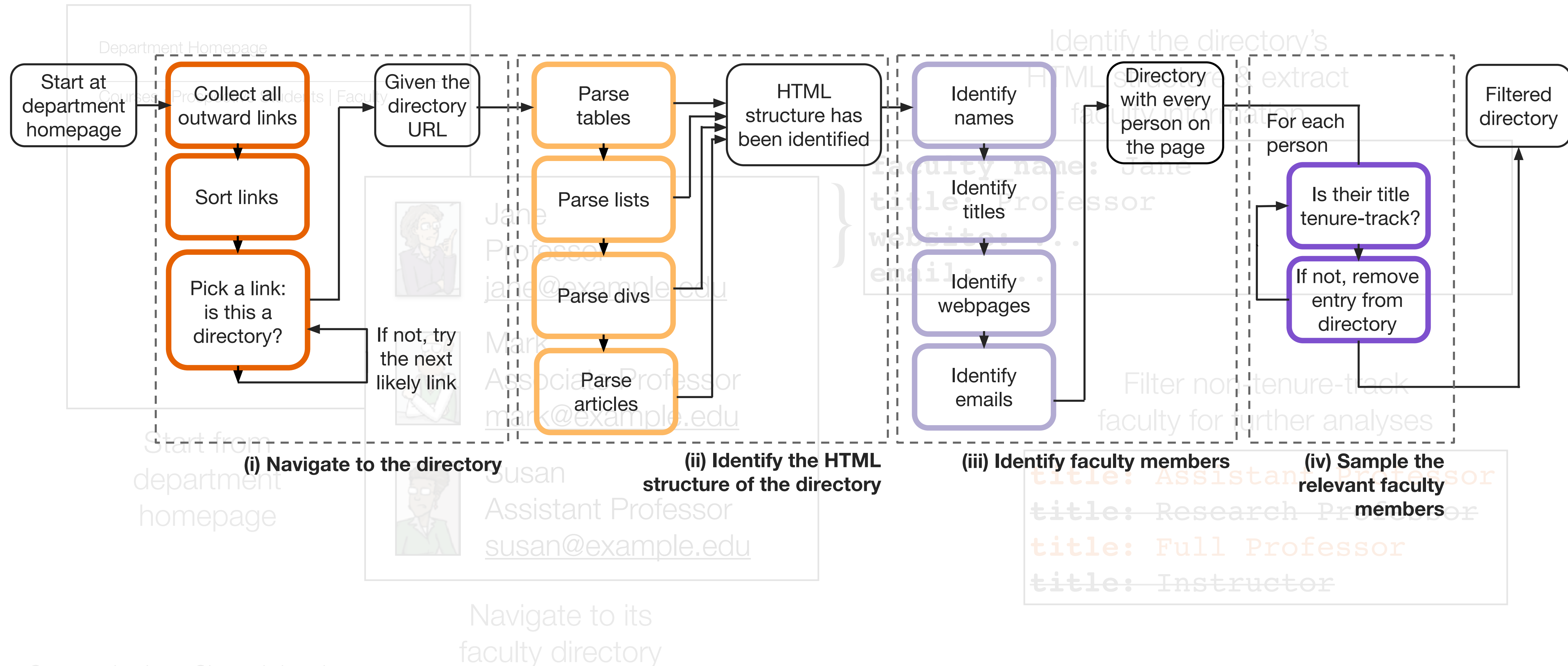
Susan  
Assistant Professor  
susan@example.edu

**title:** Assistant Professor  
~~**title:** Research Professor~~  
**title:** Full Professor  
~~**title:** Instructor~~

Navigate to its  
faculty directory

Start from  
department  
homepage

# Our Approach



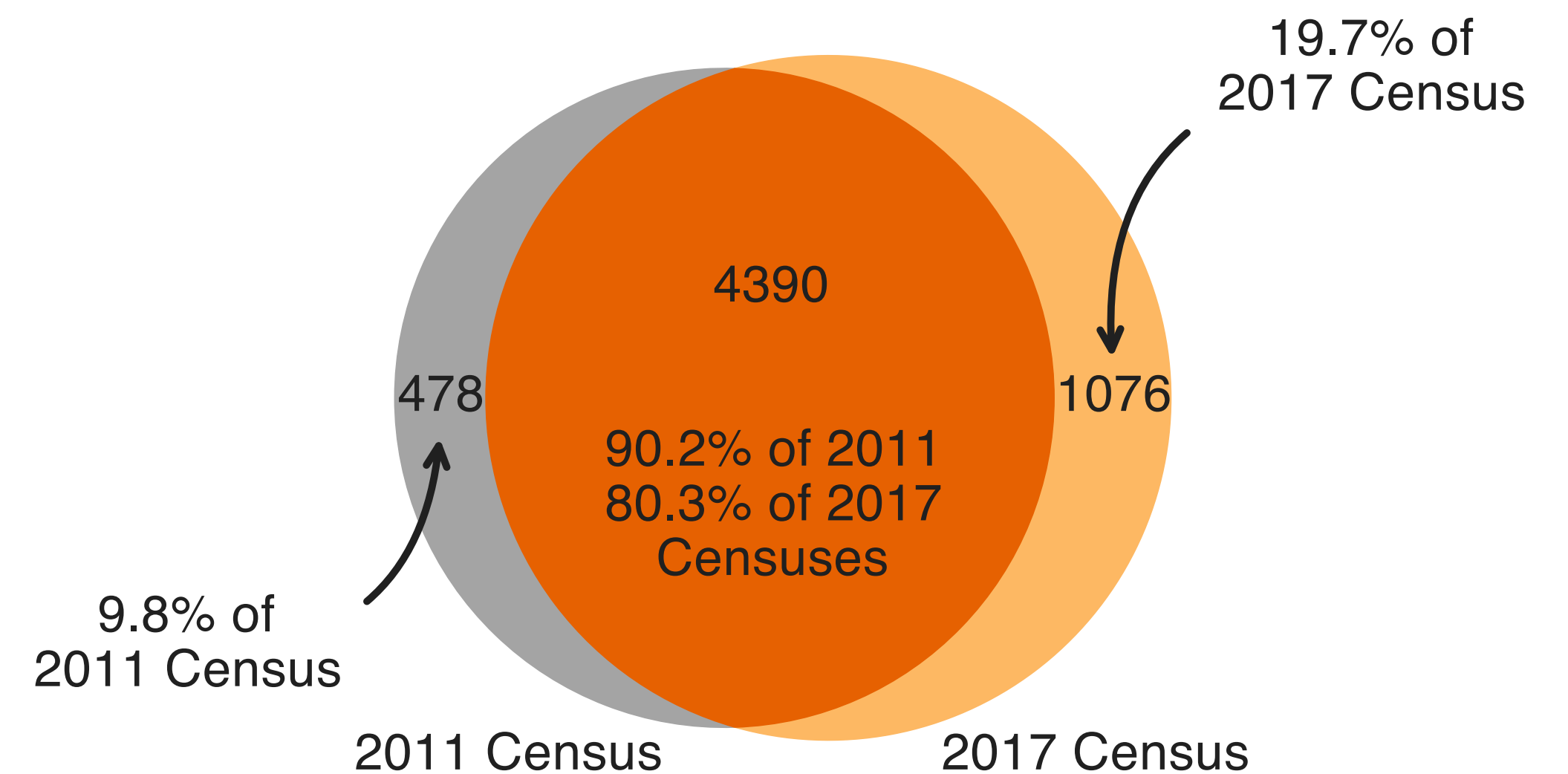
# Engineering Results

**Fast:** average < 1 minute vs ~8 hours to produce a single department's faculty directory

**Accurate:** 99% recall (nearly all tenure-track faculty are retrieved) and precision (few non-tenure-track faculty are retrieved)

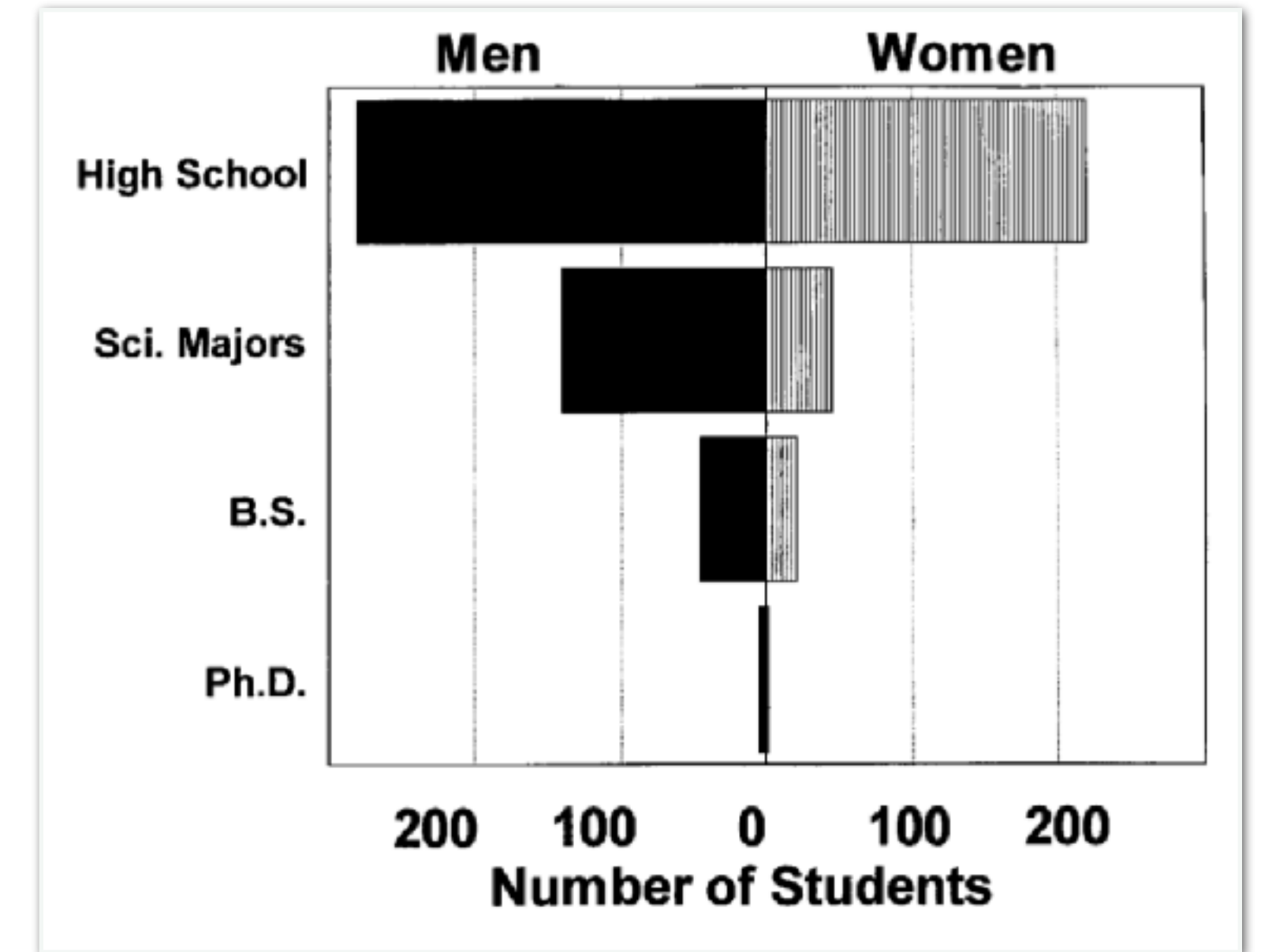
**Reproduces findings of major survey**

**organization:** 12% vs 11%  
net growth in the number of  
faculty from the CRA

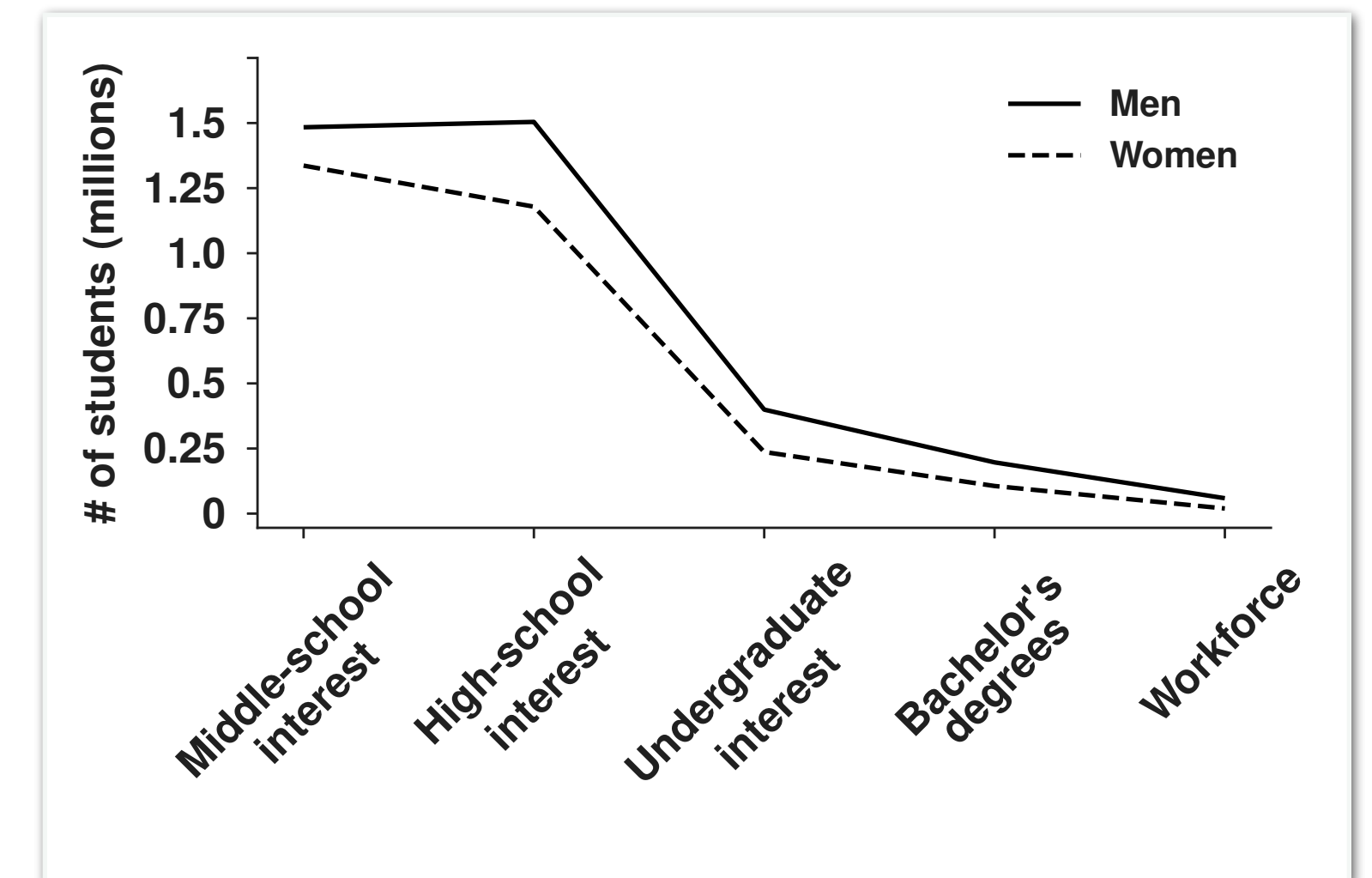


# So what can we do with this tool?

We investigate the “**leaky pipeline**”: women leave STEM at various career stages, resulting in their under-representation at the faculty level



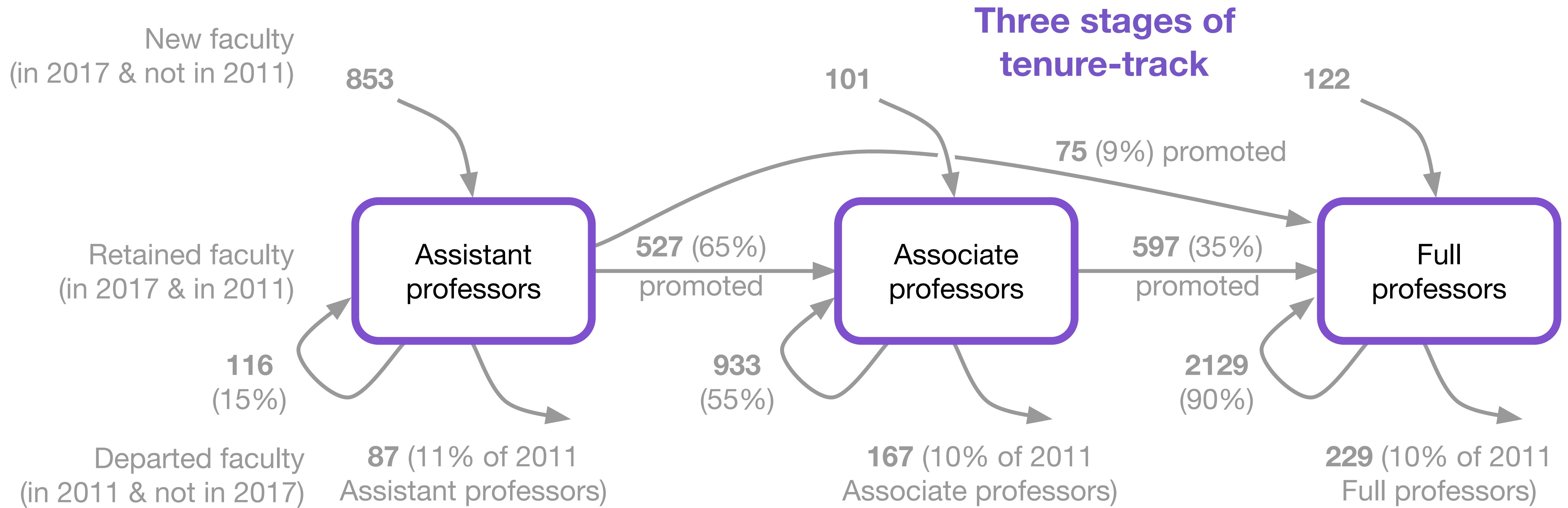
*Journal of Animal Science*, 74(11), 2843-2848, 1996



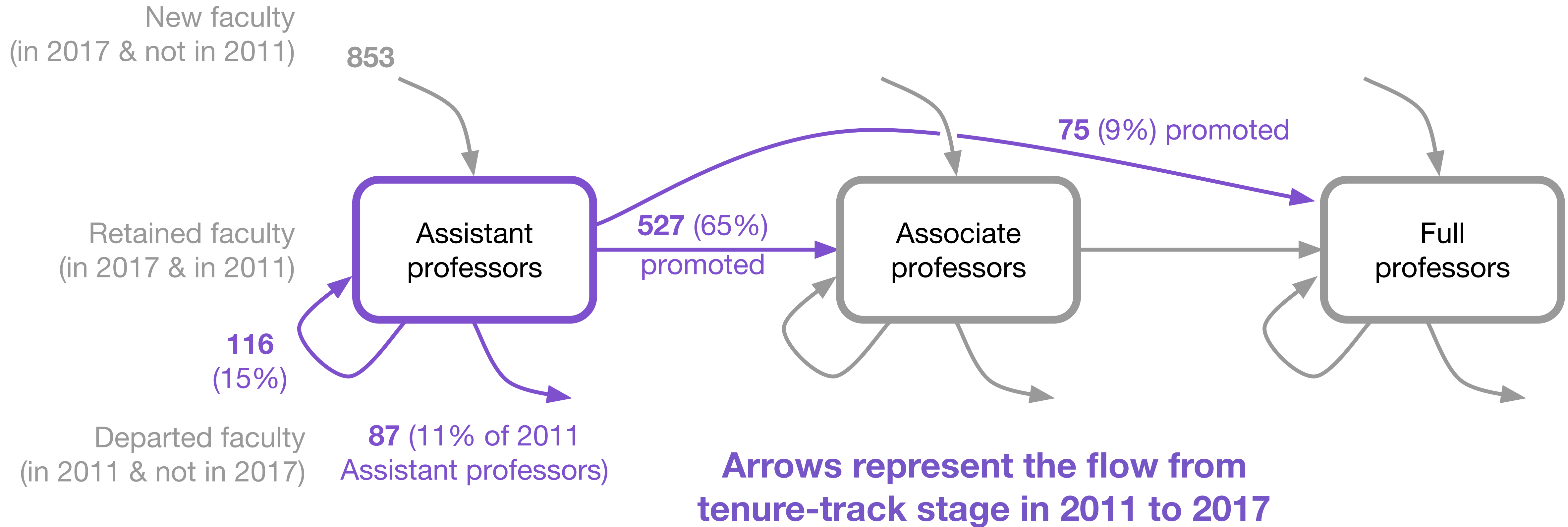
*PloS ONE*, 11(7), e0157447, 2016



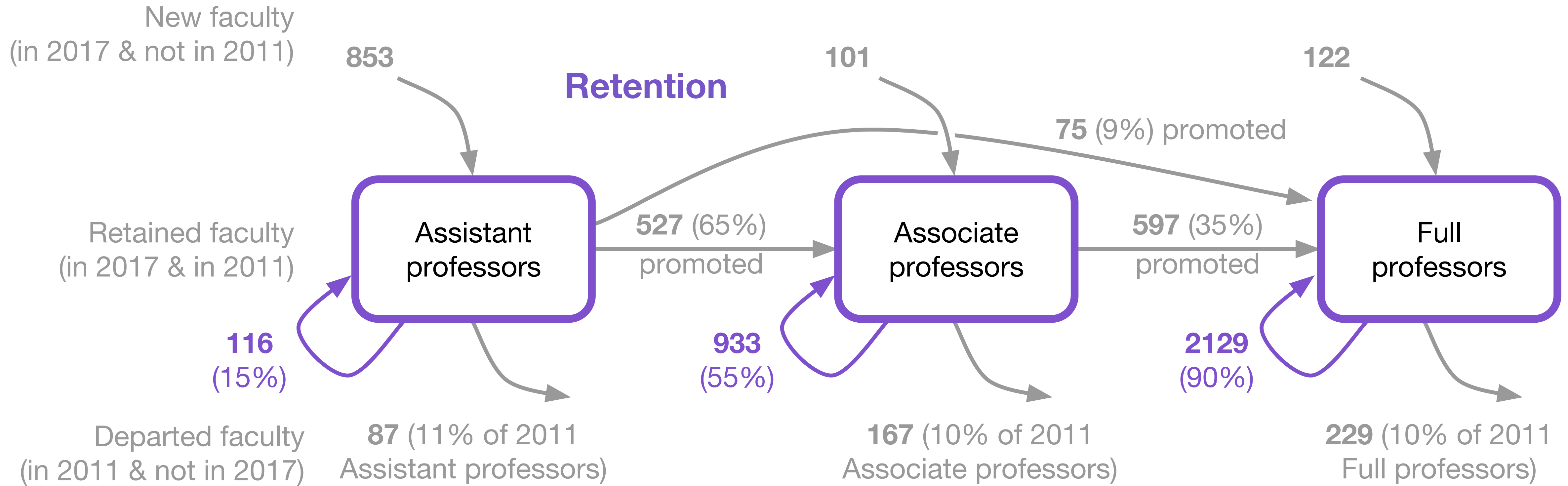
# Leaky Pipeline



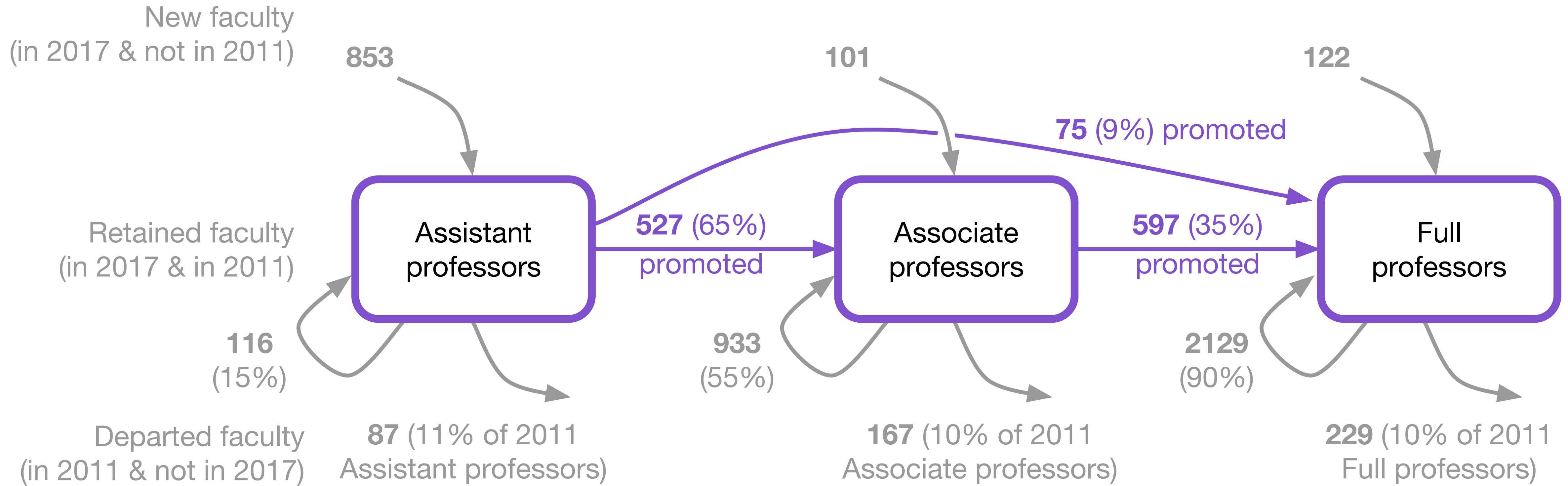
# Leaky Pipeline



# Leaky Pipeline



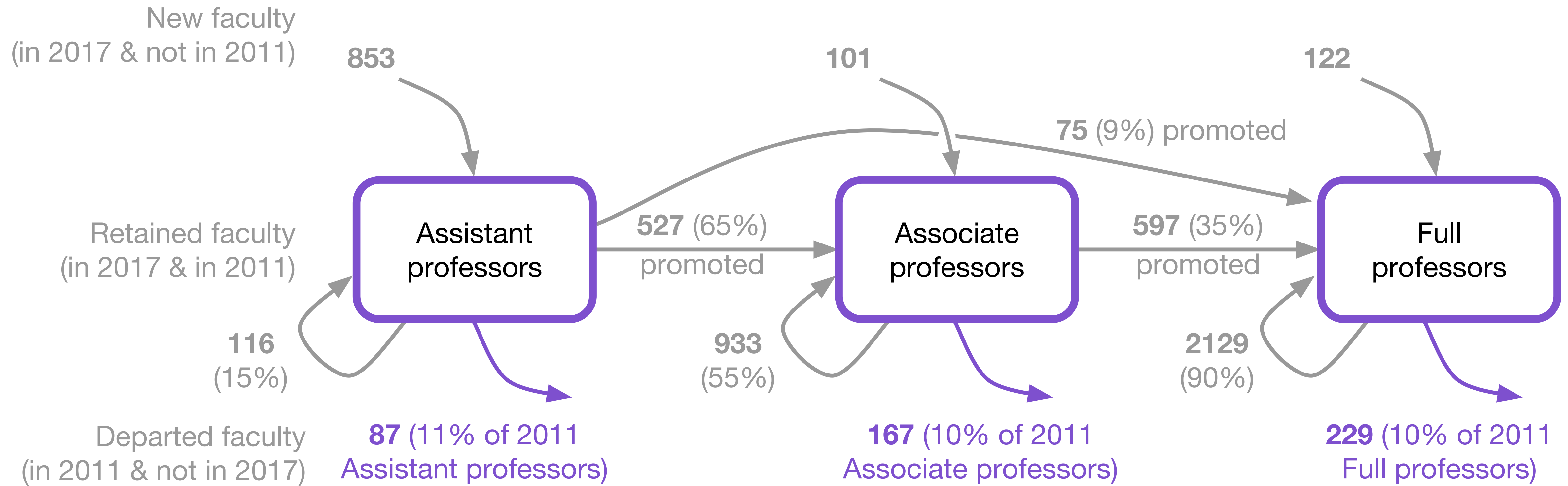
# Leaky Pipeline



**Promotion**

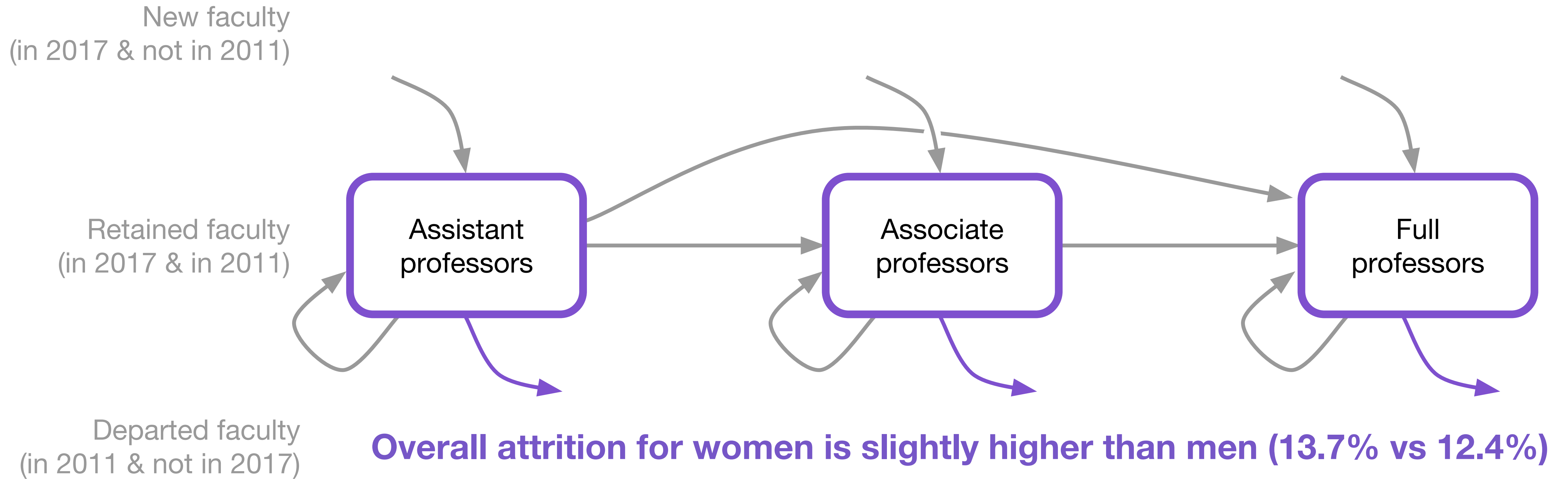


# Leaky Pipeline

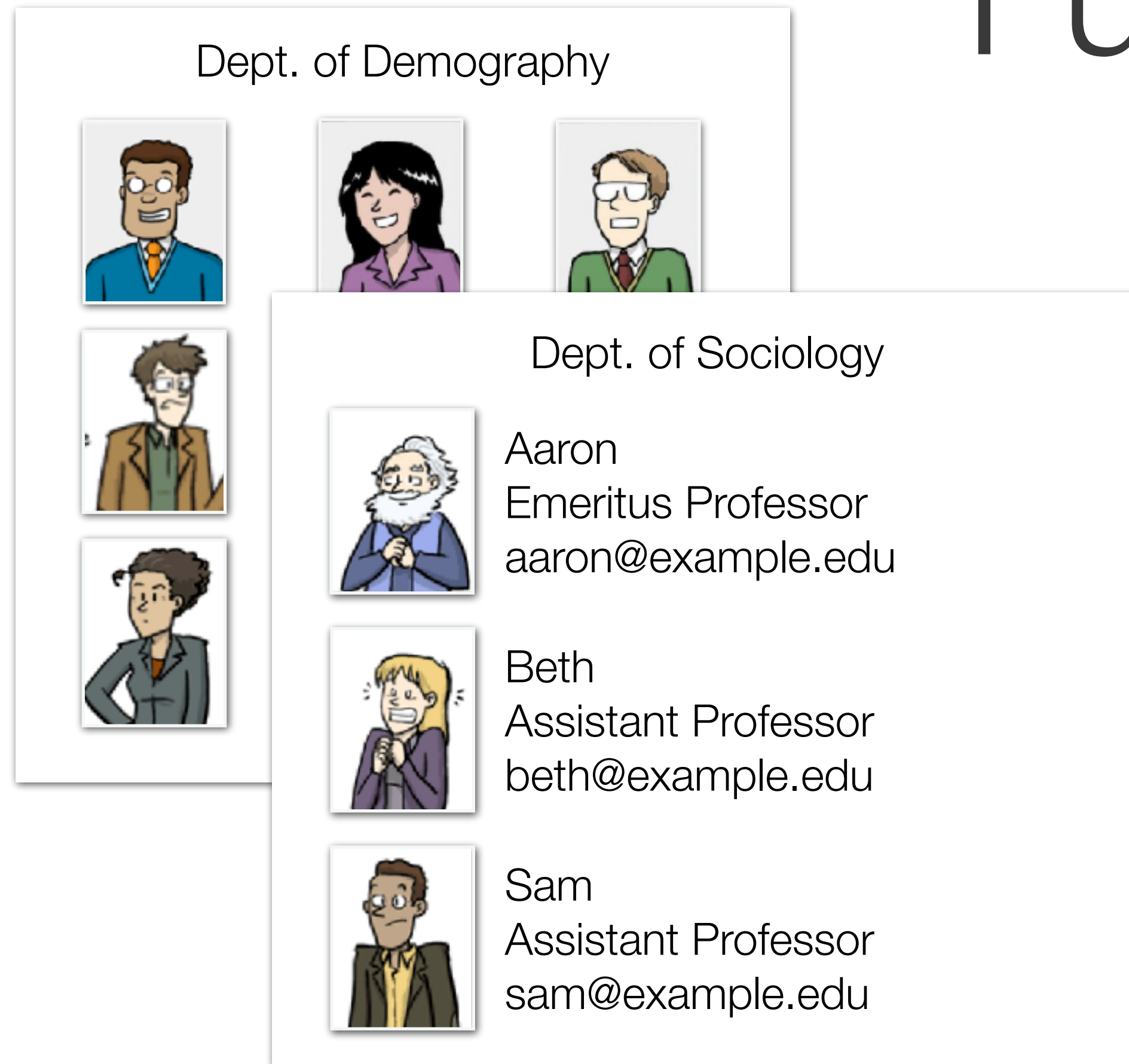


**Attrition**

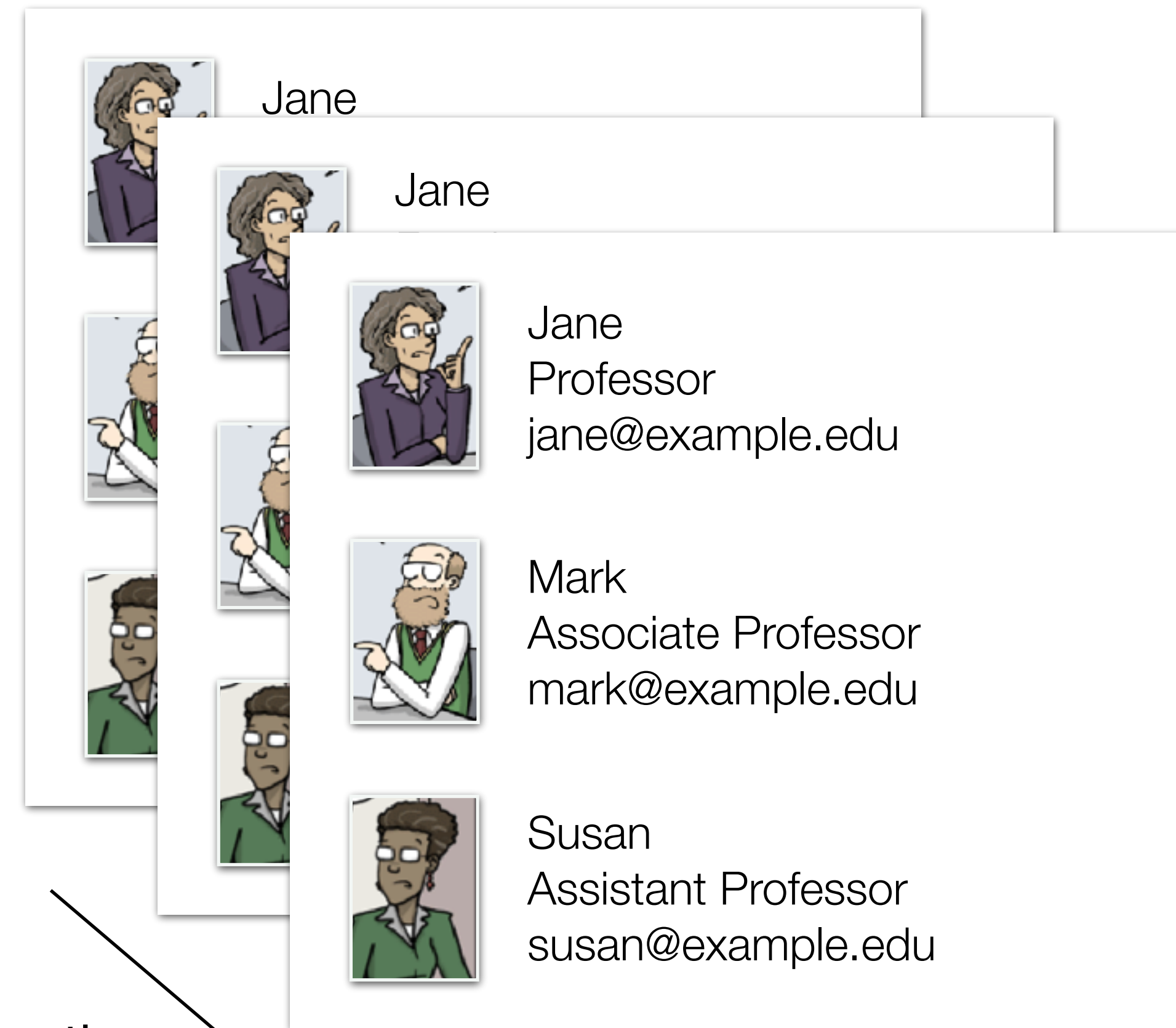
# Leaky Pipeline



# Future Work



Expand support to other  
academic fields



Use the InternetArchive to  
collect the historical data

# Thanks!

## Automatically assembling a full census of an academic field

Allison C. Morgan,<sup>1,\*</sup> Samuel F. Way,<sup>1,†</sup> and Aaron Clauset<sup>1,2,3,‡</sup>

<sup>1</sup>*Department of Computer Science, University of Colorado, Boulder, CO, USA*

<sup>2</sup>*BioFrontiers Institute, University of Colorado, Boulder, CO, USA*

<sup>3</sup>*Santa Fe Institute, Santa Fe, NM, USA*

<https://arxiv.org/abs/1804.02760>



Prof. Aaron Clauset  
PhD Computer Science  
[aaron.clauset@colorado.edu](mailto:aaron.clauset@colorado.edu)



Dr. Sam Way  
PhD Computer Science  
[samuel.way@colorado.edu](mailto:samuel.way@colorado.edu)



University of Colorado **Boulder**