

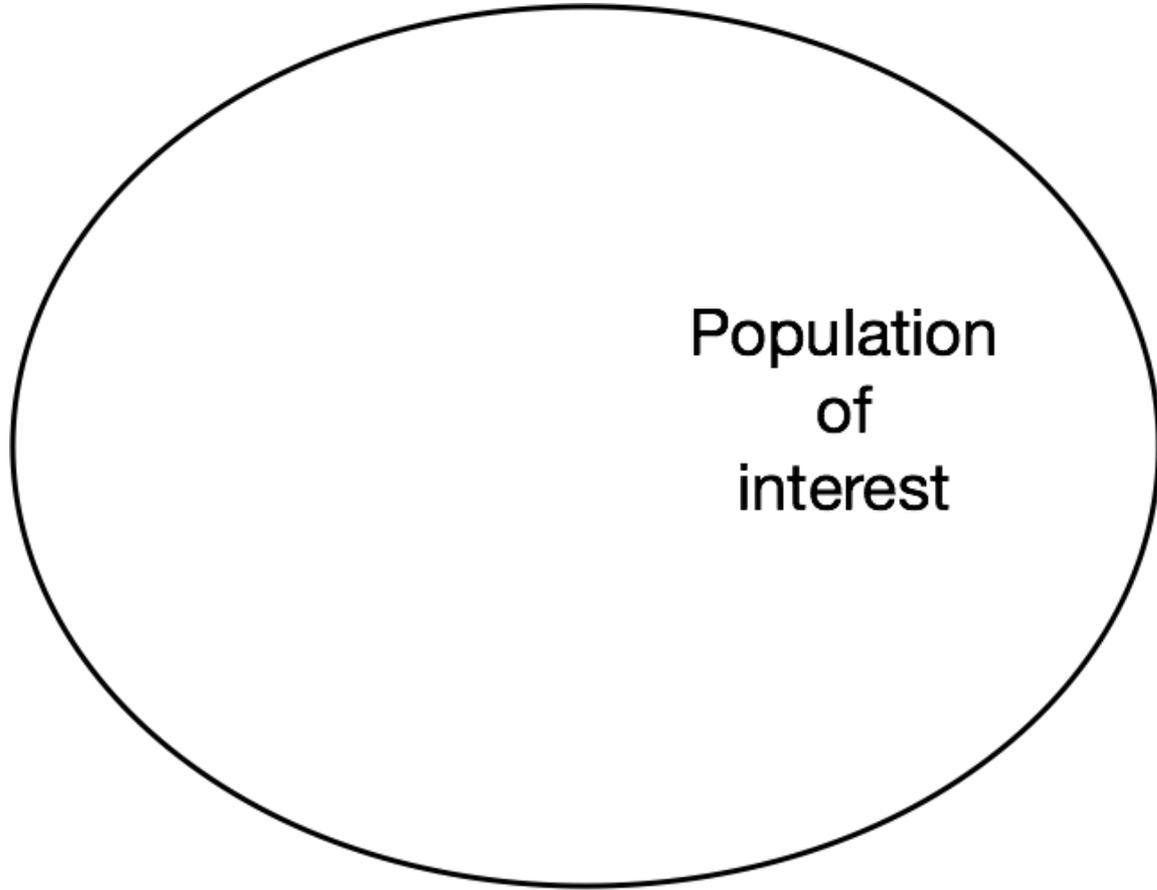
# **Estimating internet adoption around the world using a sample of Facebook users**

**PAA 2018**

Dennis M. Feehan  
UC Berkeley

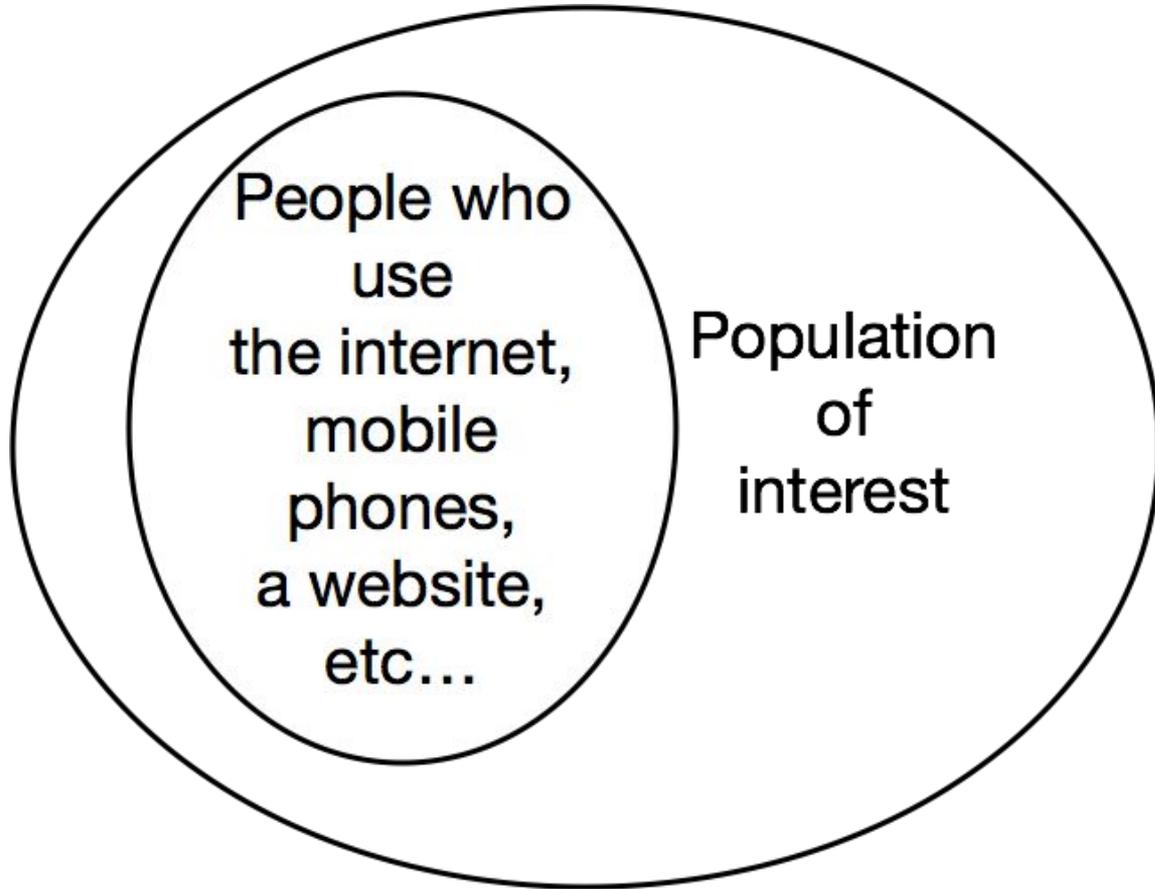
Curtiss Cobb  
Facebook

The problem



Population  
of  
interest

# The problem



# The problem

## Digital inequalities and why they matter

Laura Robinson<sup>a</sup>, Shelia R. Cotten<sup>b</sup>, Hiroshi Ono<sup>c</sup>, Anabel Wenhong Chen<sup>f</sup>, Jeremy Schulz<sup>g\*</sup>, Timothy M. Hale<sup>h</sup> and

<sup>a</sup>Department of Sociology, Santa Clara University, Santa Clara, CA, USA; <sup>b</sup>Department of Information, Michigan State University, East Lansing, MI, USA; <sup>c</sup>Hitotsubashi University School of International Corporate Strategy, Hitotsubashi University, Tokyo, Japan; <sup>d</sup>Department of Anthropology, University of Western Ontario, London, ON, Canada; <sup>e</sup>Department of Anthropology, University of Texas at Austin, Austin, TX, USA; <sup>f</sup>Institute for Connected Health, Part of the Center for Excellence in Survey Research, University of Michigan, Ann Arbor, MI, USA; <sup>g</sup>Department of Sociology, University of Michigan, Ann Arbor, MI, USA; <sup>h</sup>Department of Sociology, University of Michigan, Ann Arbor, MI, USA

### The Arrival of Fast Internet and Skilled Job Creation in Africa\*

Jonas Hjort  
Columbia University  
& BREAD & NBER

Jonas Poulsen  
Harvard University

September 10, 2016

WWW.ECONSTOR.EU

ECONSTOR

Der Open-Access-Publikationsserver der ZBW – Leibniz-Informationzentrum Wirtschaft  
The Open Access Publication Server of the ZBW – Leibniz Information Centre for Economics

Chinn, Menzie D.; Fairlie, Robert W.

Working Paper

The Determinants of the Global Digital Divide : A Cross-Country Analysis of Computer and Internet Penetration

IZA Discussion paper series, No. 1305

Provided in Cooperation with:  
Institute for the Study of Labor (IZA)

Suggested Citation: Chinn, Menzie D.; Fairlie, Robert W. (2004) : The Determinants of the Global Digital Divide : A Cross-Country Analysis of Computer and Internet Penetration, IZA Discussion paper series, No. 1305

ing its current frontier in devel-  
opment, the network of submarine  
cables that connects those cities  
covering 14 countries show large  
fact by a bigger increase in jobs  
between more and less educated  
s to investigate how higher aver-  
age job creation. We find an increase  
d the productivity of workers in  
ry in South Africa in the sectors  
comes available, and (iii) work-  
ers in developing countries. Finally, we show  
mes. Our findings shed light on  
b creation, structural change, job

January 2015)

continues to expand in  
other forms of inequality  
in and outside the field  
stantive problem and as  
on multiple aspects  
age, skills, and self-perce  
s makes the case that di  
of inequality in the t  
ould not be only the pr  
cial scientists concerned  
e trajectories. As we a  
ad range of individual-  
e, and class, as well as

munication; digital divi

A World Bank Group Flagship Report

102725

world development report

2016

Public Disclosure Authorized

Public Disclosure Authorized

Public Disclosure Authorized



# DIGITAL DIVIDENDS

# The plan

- Methods: network reporting with an online sample
- Study design: estimating internet adoption in 5 countries
- Results: estimates and sensitivity
- Next steps

Methods: network reporting

# Methods: network reporting

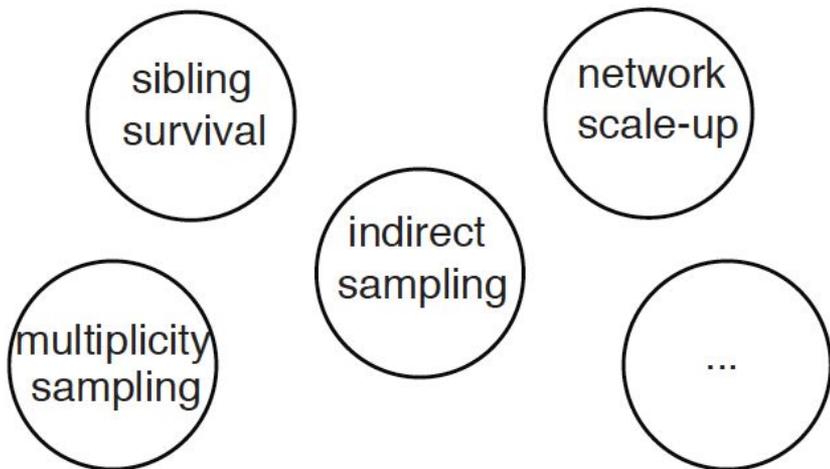
The idea: survey respondents are connected to other people through many different kinds of personal networks

We can ask respondents questions about their personal network and learn about more than just the respondent.

# Network reporting

Approaches like this have been used in lots of different situations

- Deaths
- Epidemiologically important groups (drug injectors, sex workers)
- Migrants
- ... and many others



## network reporting

sibling  
survival

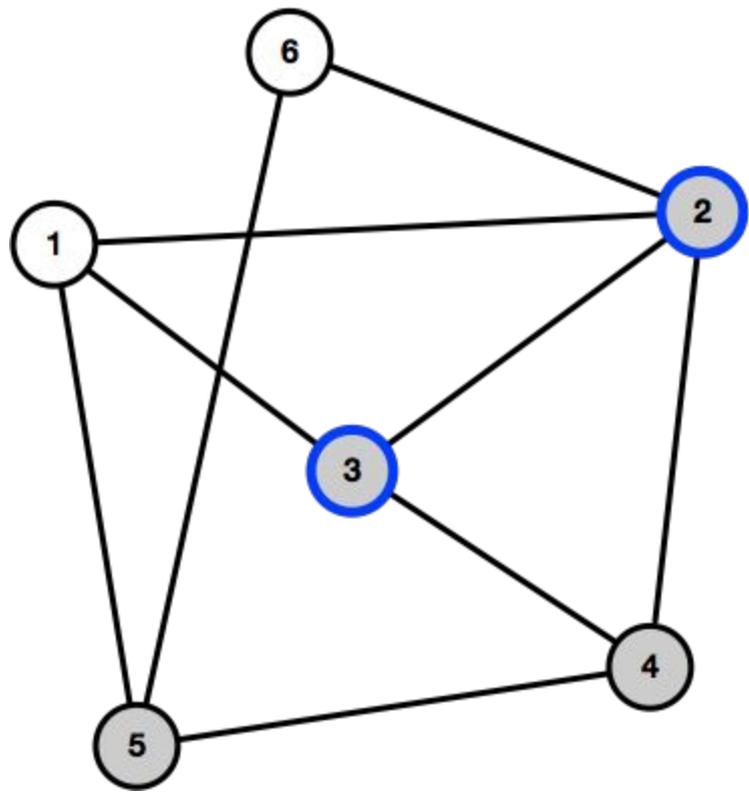
network  
scale-up

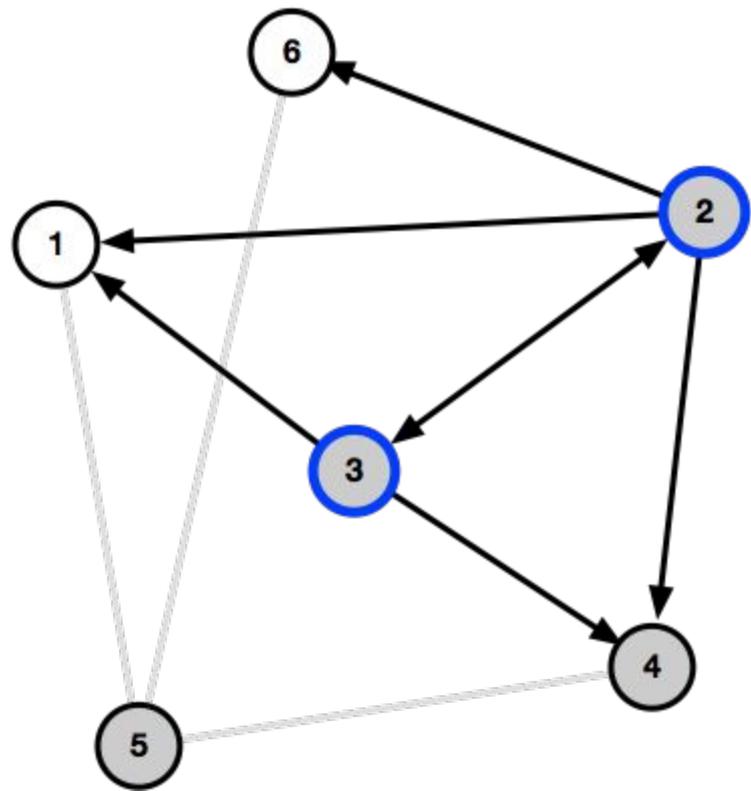
indirect  
sampling

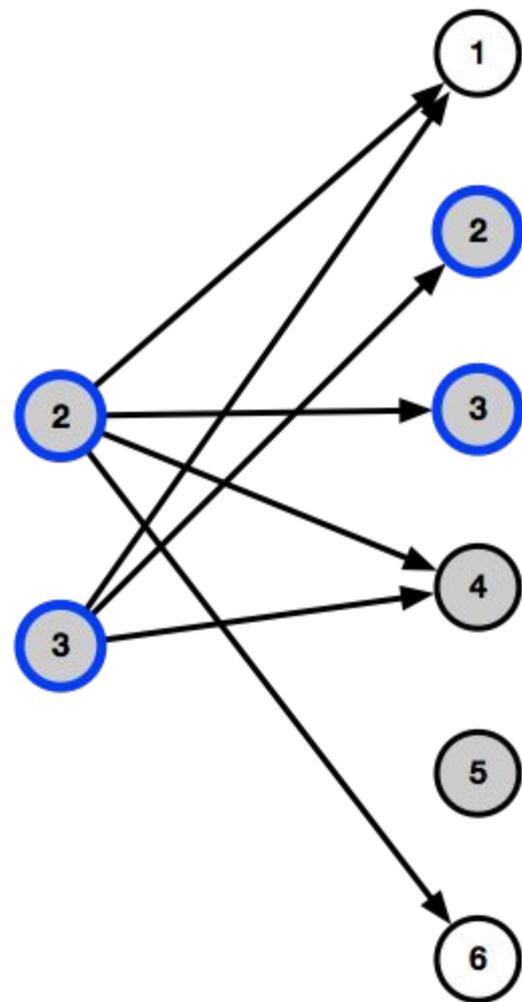
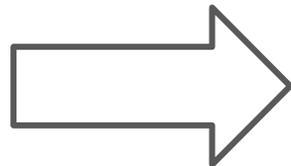
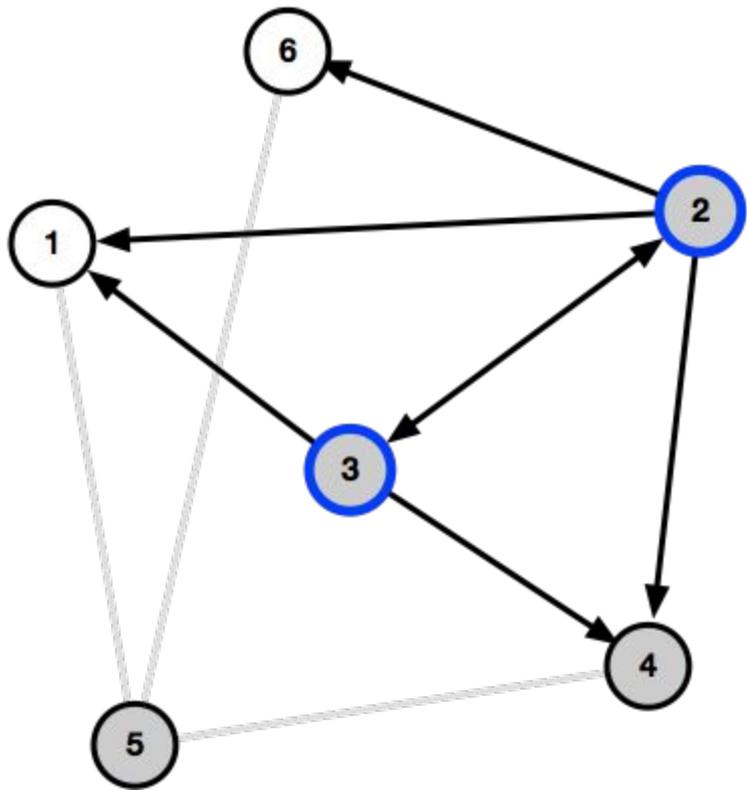
multiplicity  
sampling

...

How it works



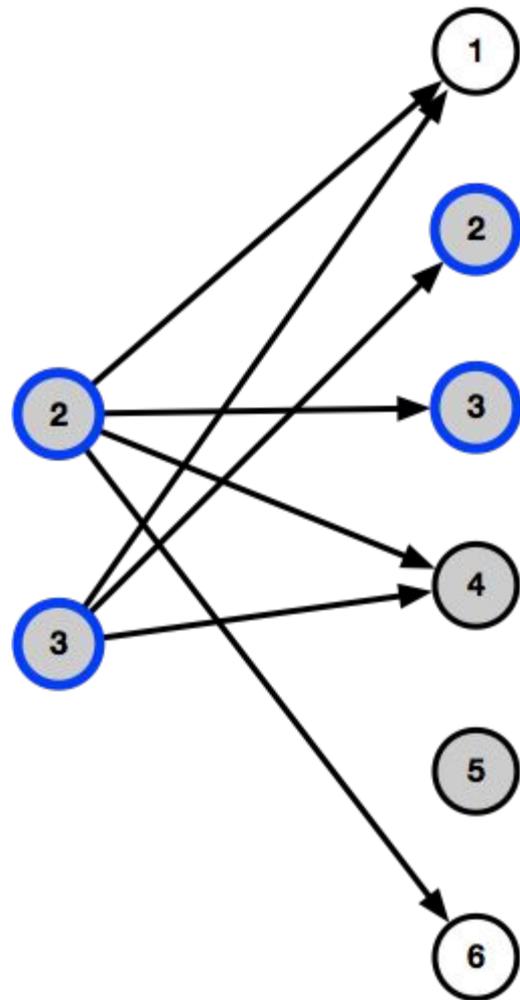




total out-reports = total in-reports

$\Leftrightarrow$  total out-reports = (number of internet users  $\times$   
average in-reports per internet user)

$\Leftrightarrow$  number of internet users =  $\frac{\text{total out-reports}}{\text{average in-reports per internet user}}$



$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

# Study design

$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

We can ask respondents questions like “how many people are in your network?”

And then, “which of these people uses the internet?”

We can ask respondents questions like “how many people are in your network?”

... but what does it mean to ‘know’ someone?

=> we need to choose a **tie definition**

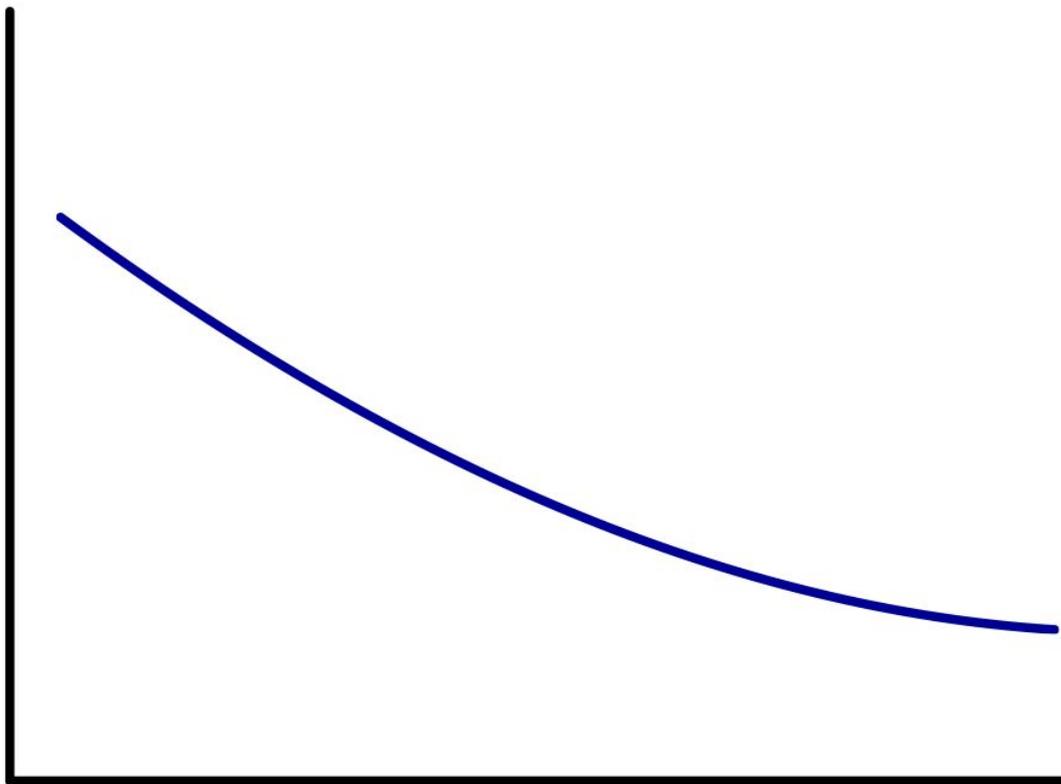
A graph with a vertical y-axis and a horizontal x-axis. The y-axis is labeled "error in estimate" and the x-axis is labeled "stronger tie" on the left and "weaker tie" on the right. The graph area is currently blank.

error in  
estimate

stronger  
tie

weaker  
tie

error in  
estimate

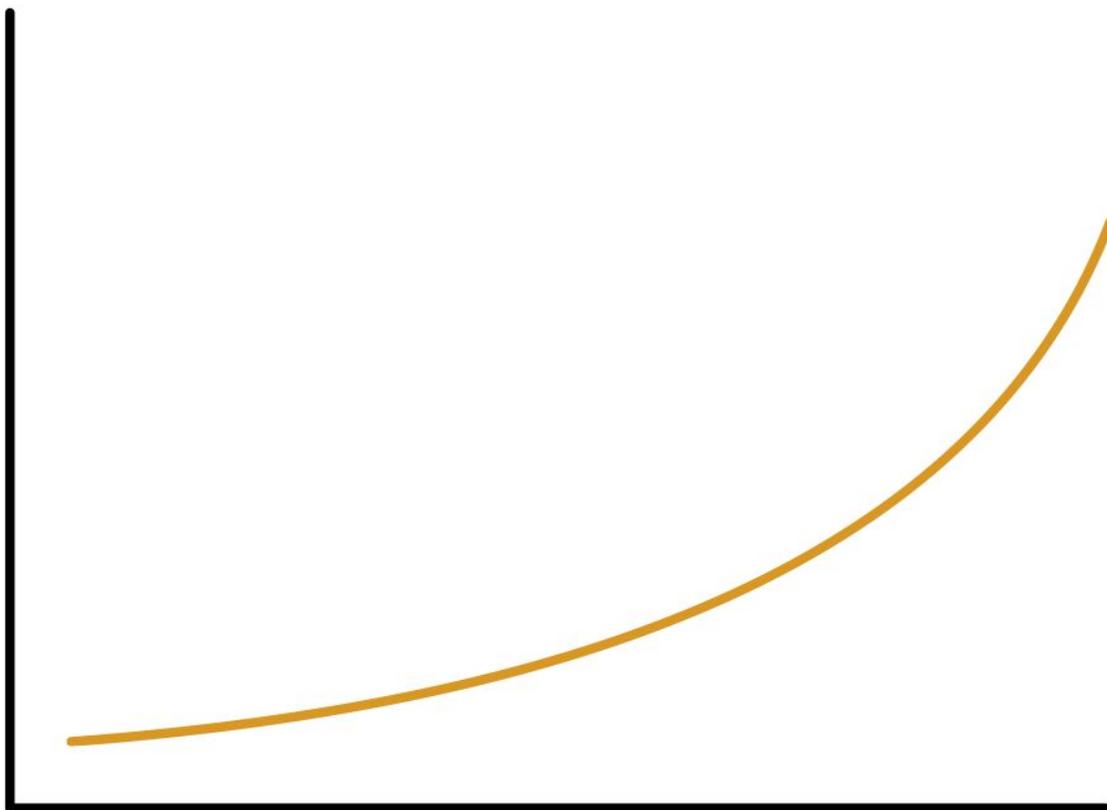


stronger  
tie

weaker  
tie

sampling  
error

error in  
estimate

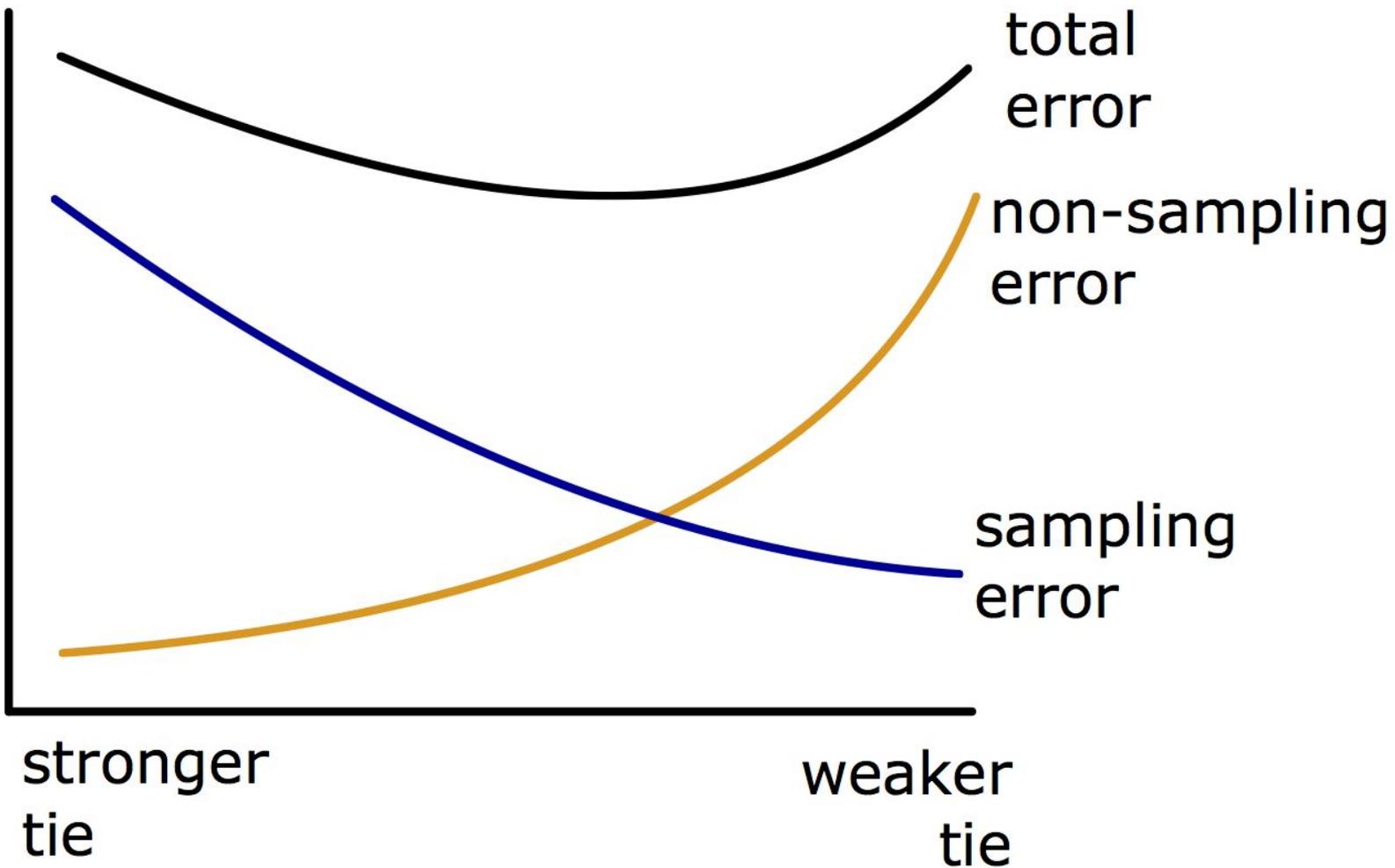


non-sampling  
error

stronger  
tie

weaker  
tie

error in estimate



total error

non-sampling error

sampling error

stronger tie

weaker tie

# Tie definition: survey experiment

- Previous research has found some evidence of a tie strength / accuracy tradeoff
- We designed an experiment to further test this question in our setting

# Tie definition: survey experiment

- Previous research has found some evidence of a tie strength / accuracy tradeoff
- We designed an experiment to further test this question in our setting

## **Conversational Contact Network**

- How many people did you have conversational contact with yesterday? By conversational contact, we mean anyone you spoke with face to face for at least three words.

# Tie definition: survey experiment

- Previous research has found some evidence of a tie strength / accuracy tradeoff
- We designed an experiment to further test this question in our setting

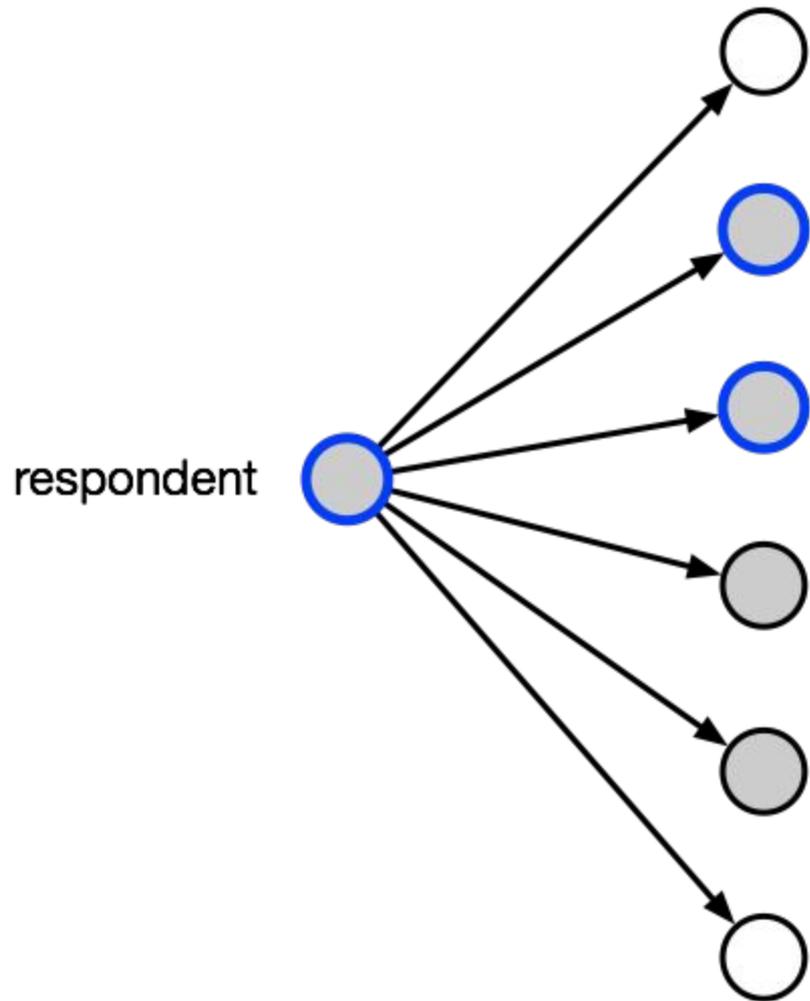
## **Conversational Contact Network**

- How many people did you have conversational contact with yesterday? By conversational contact, we mean anyone you spoke with face to face for at least three words.

## **Meal Network**

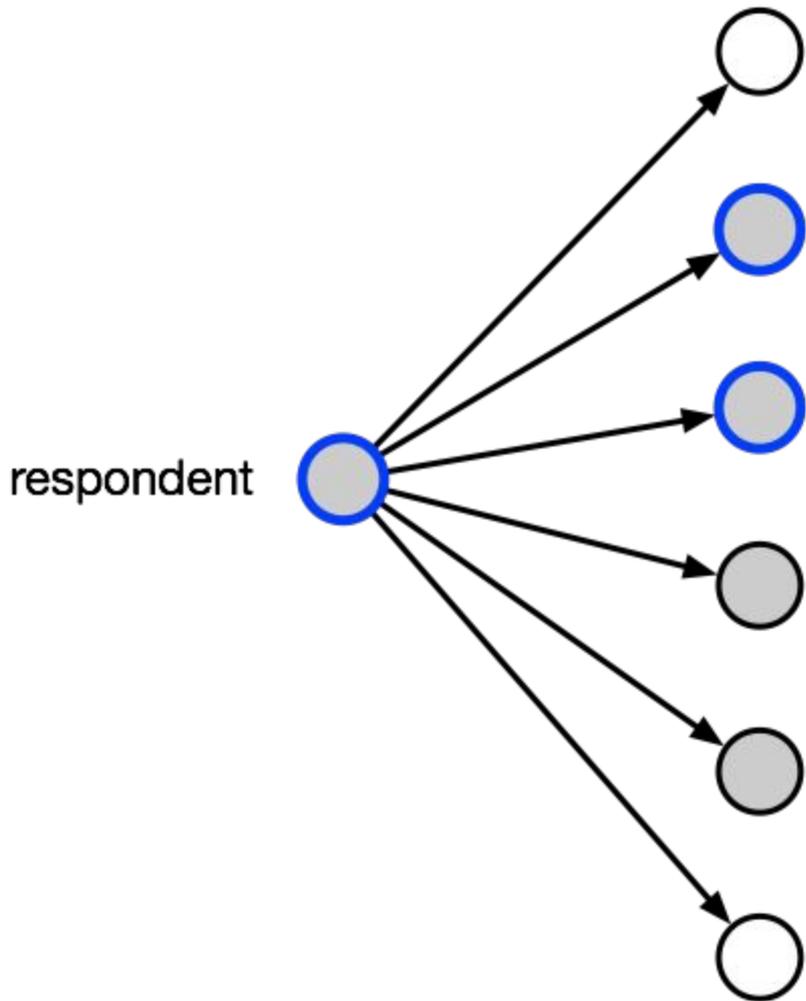
- How many people did you share food or drink with yesterday? These people could be family members, friends, co-workers, neighbors, or other people. Please include all food and drink taken at any location, including at home, at work, at a cafe, or in a restaurant.

How many people did you share food or drink with yesterday?



How many people did you share food or drink with yesterday?

=> response tells us about network size

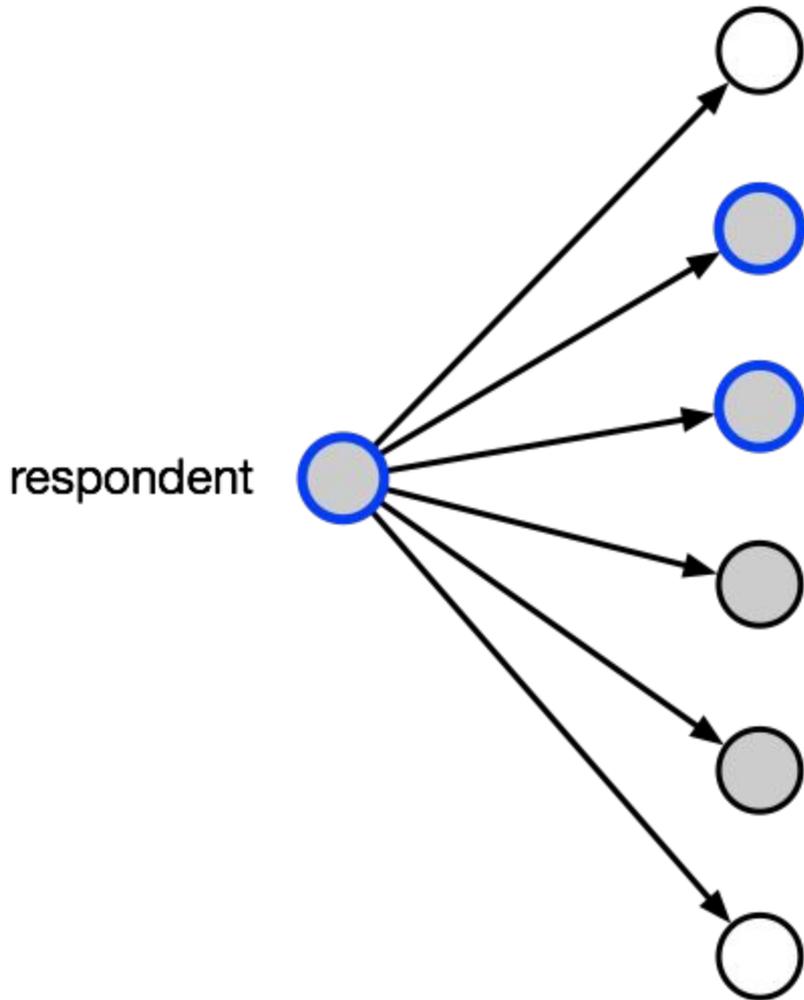


How many people did you share food or drink with yesterday?

=> response tells us about network size

Next, we want to know what proportion of respondent's network uses the internet.

Ideally: ask respondent about each person in her network, one after another



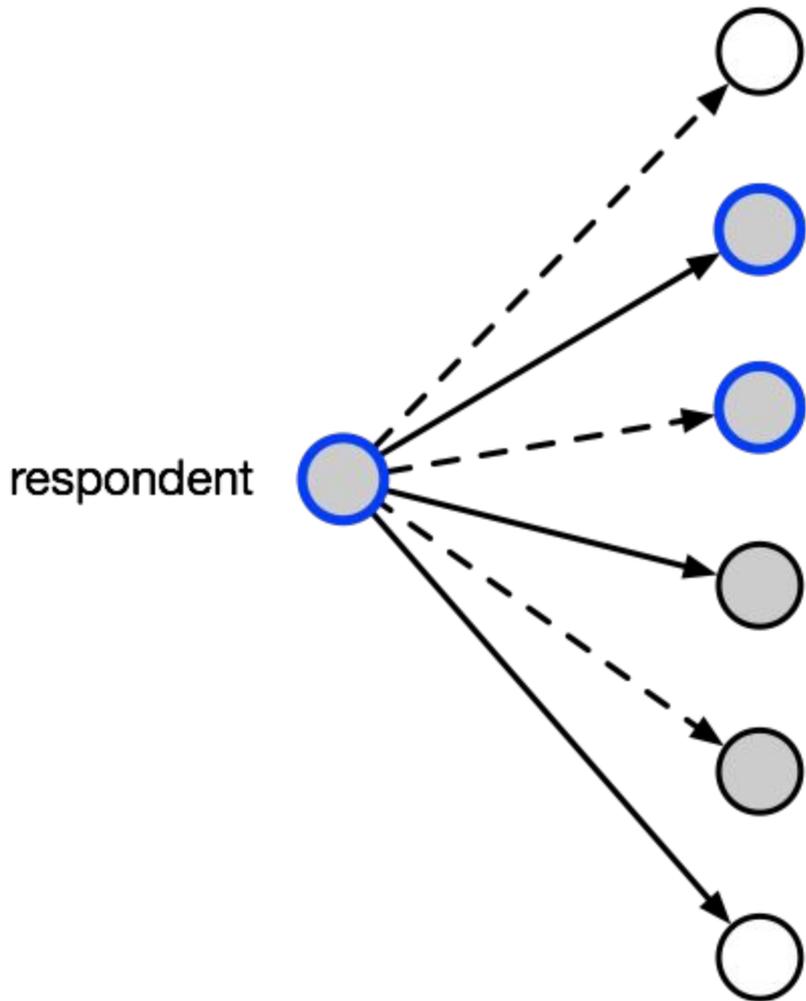
How many people did you share food or drink with yesterday?

=> response tells us about network size

Next, we want to know what proportion of respondent's network uses the internet.

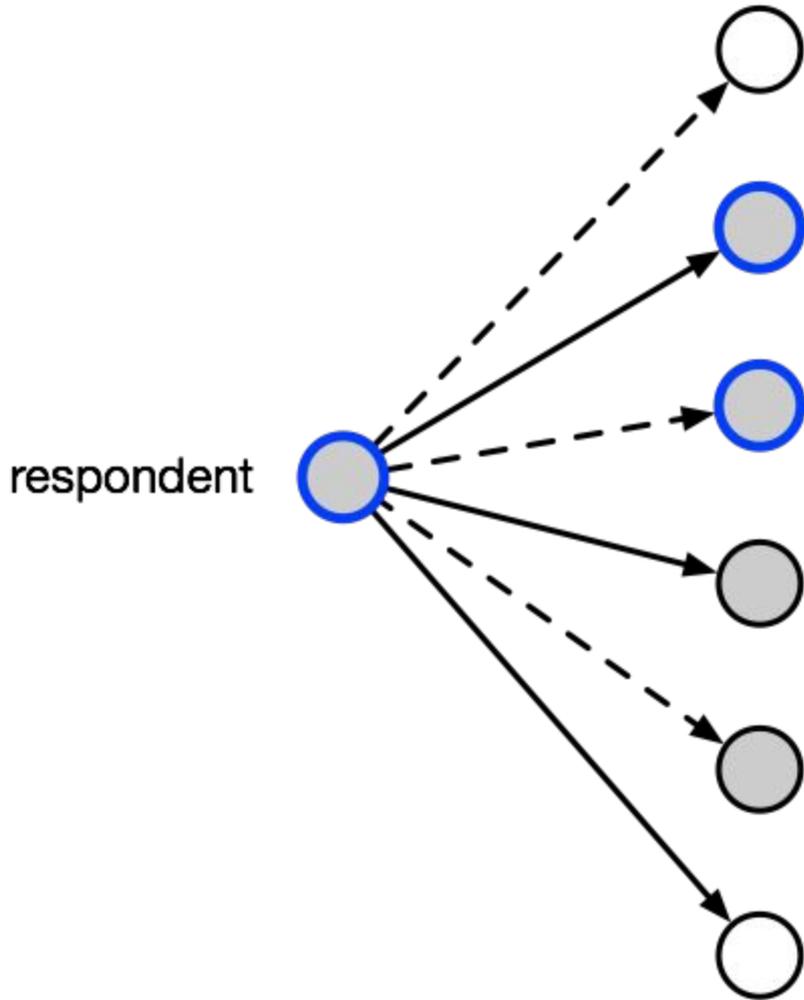
Ideally: ask respondent about each person in her network, one after another

Problem: this would likely cause unacceptable levels of respondent fatigue



Instead, we ask respondents about a subset of their network contacts; we call this subset the **detailed alters**

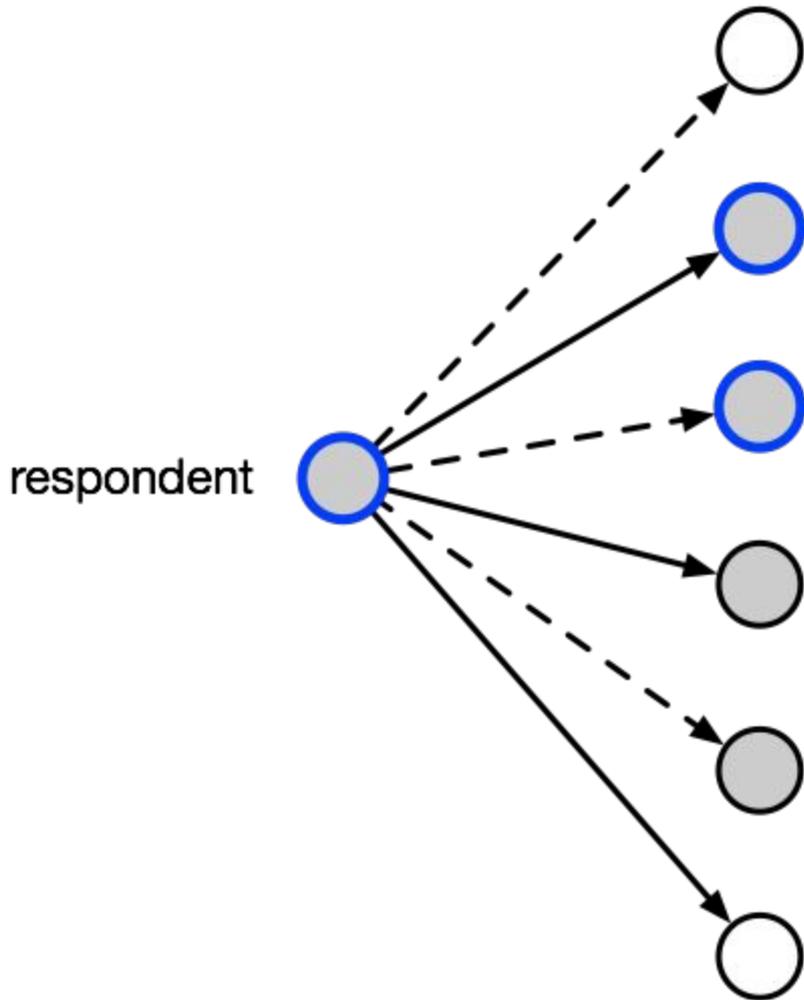
We ask for information about the three network members who 'come to mind' first



Instead, we ask respondents about a subset of their network contacts; we call this subset the **detailed alters**

We ask for information about the three network members who 'come to mind' first

We treat these three detailed alters as if they were a simple random sample of the respondent's network members



Instead, we ask respondents about a subset of their network contacts; we call this subset the **detailed alters**

We ask for information about the three network members who 'come to mind' first

We treat these three detailed alters as if they were a simple random sample of the respondent's network members

-> in reality, some alters are probably more likely to come to mind than others

-> paper mathematically describes how estimates are sensitive to this condition

-> and we can check this empirically

$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

# Estimating visibility

Can imagine many different approaches to this

# Estimating visibility

Can imagine many different approaches to this

We chose something very simple: assume that people do not pay attention to whether or not they are on Facebook when they share meals with one another

# Estimating visibility

Can imagine many different approaches to this

We chose something very simple: assume that people do not pay attention to whether or not they are on Facebook when they share meals with one another

Our approach works if two quantities are equal:

- The rate at which people on the internet share meals with someone on FB
- The rate at which people on FB share meals with someone else on FB

# Estimating visibility

Can imagine many different approaches to this

We chose something very simple: assume that people do not pay attention to whether or not they are on Facebook when they share meals with one another

Our approach works if two quantities are equal:

- The rate at which people on the internet share meals with someone on FB
- The rate at which people on FB share meals with someone else on FB

We can estimate the second quantity from our survey responses

# Putting it all together

$$\# \text{ of internet users} = \frac{\text{total reported connections to internet users}}{\text{average in-reports per internet user}}$$

# Recap: 3 key conditions

- Accurate reporting
- Detailed alters picked at random
- Meals shared between people without paying attention to whether they are on Facebook or not

Our paper has sensitivity framework for understanding what impact violating these conditions would have on estimates

Framework also shows how these conditions can be relaxed or eliminated if additional data can be collected

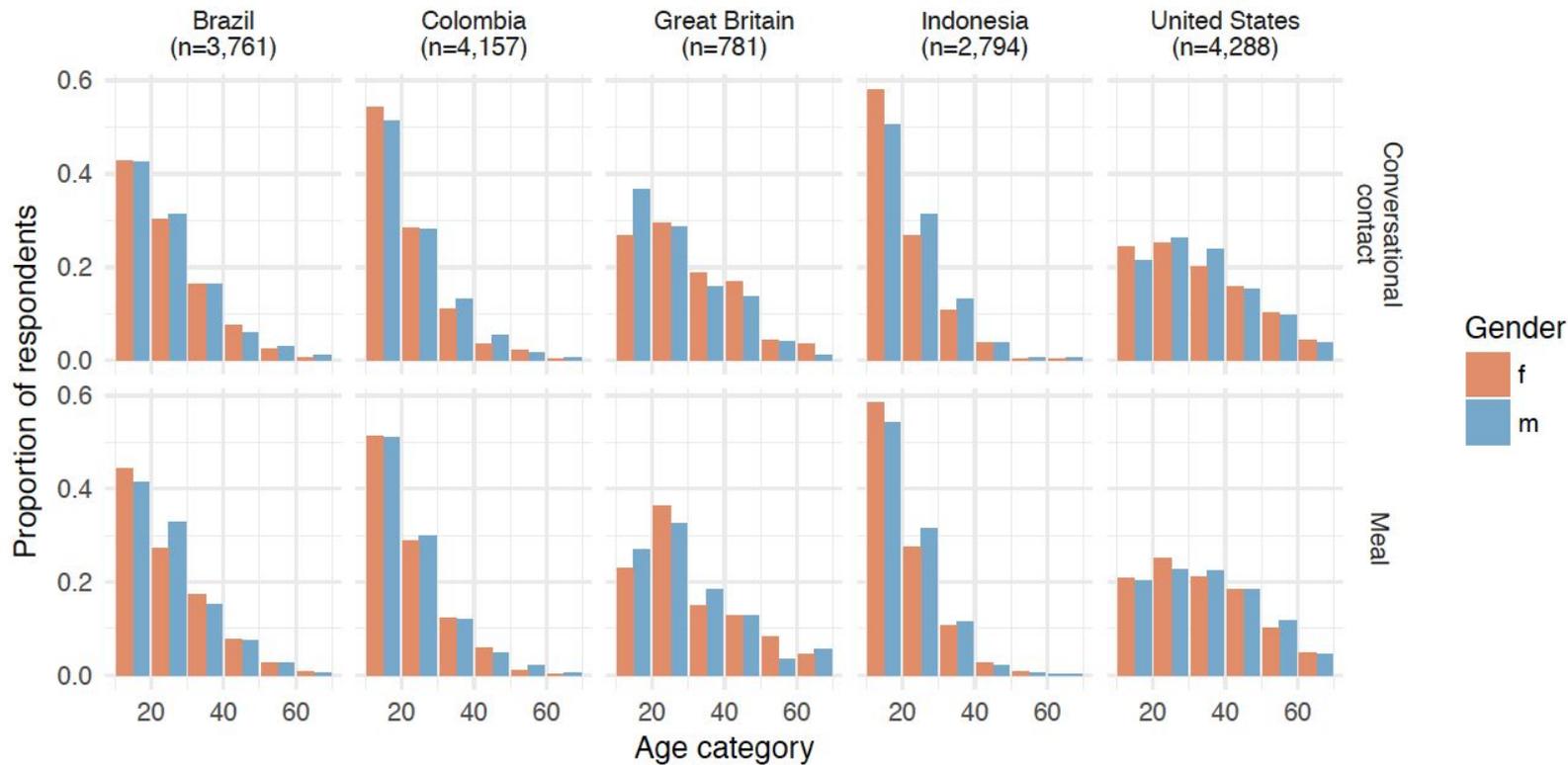
We'll see that the first two conditions can be checked empirically

# Results

# Sample

- Random sample of Facebook users, taken using FB's survey infrastructure
- Short survey, taken over web or mobile
- Looked at lots of calibration and post-stratification approaches, found that these mattered very little
- All analyses use rescaled bootstrap to estimate sampling uncertainty

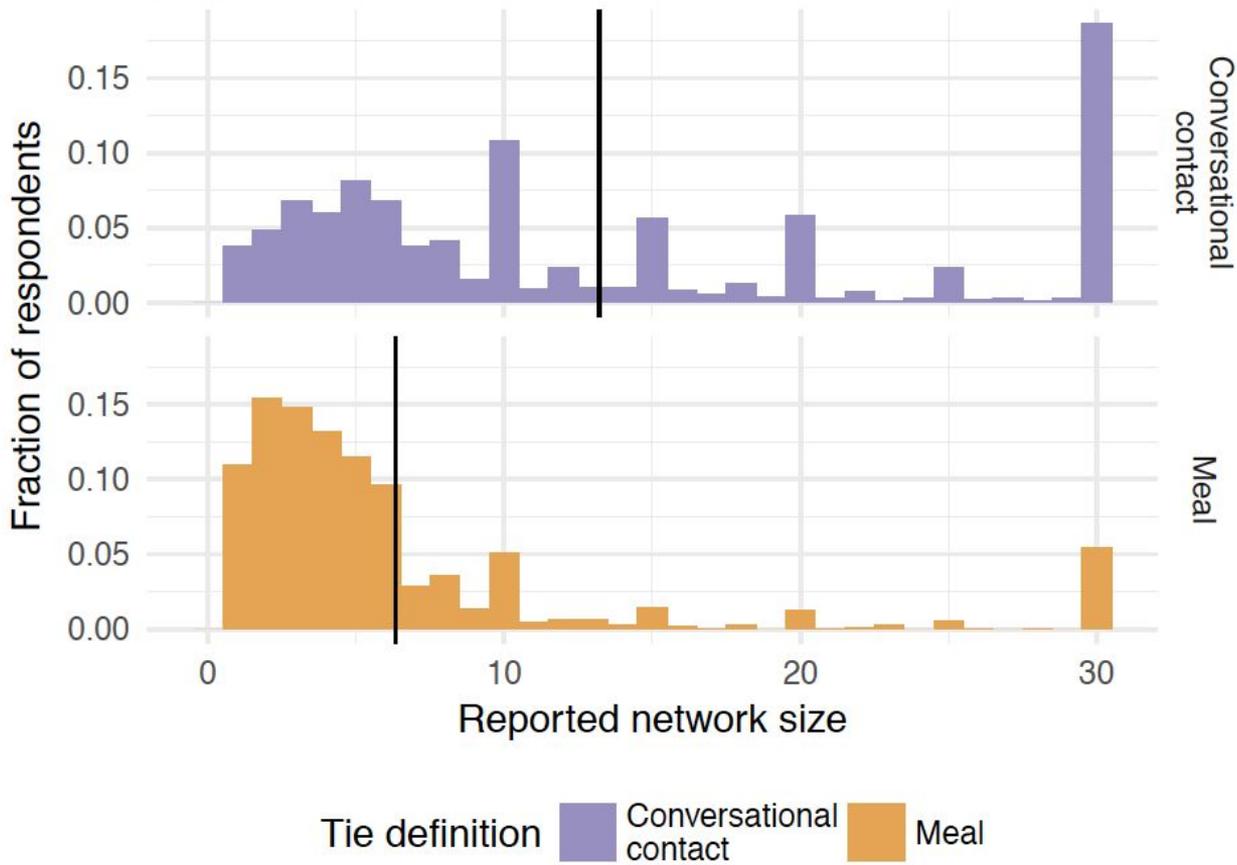
# Sample: 5 countries



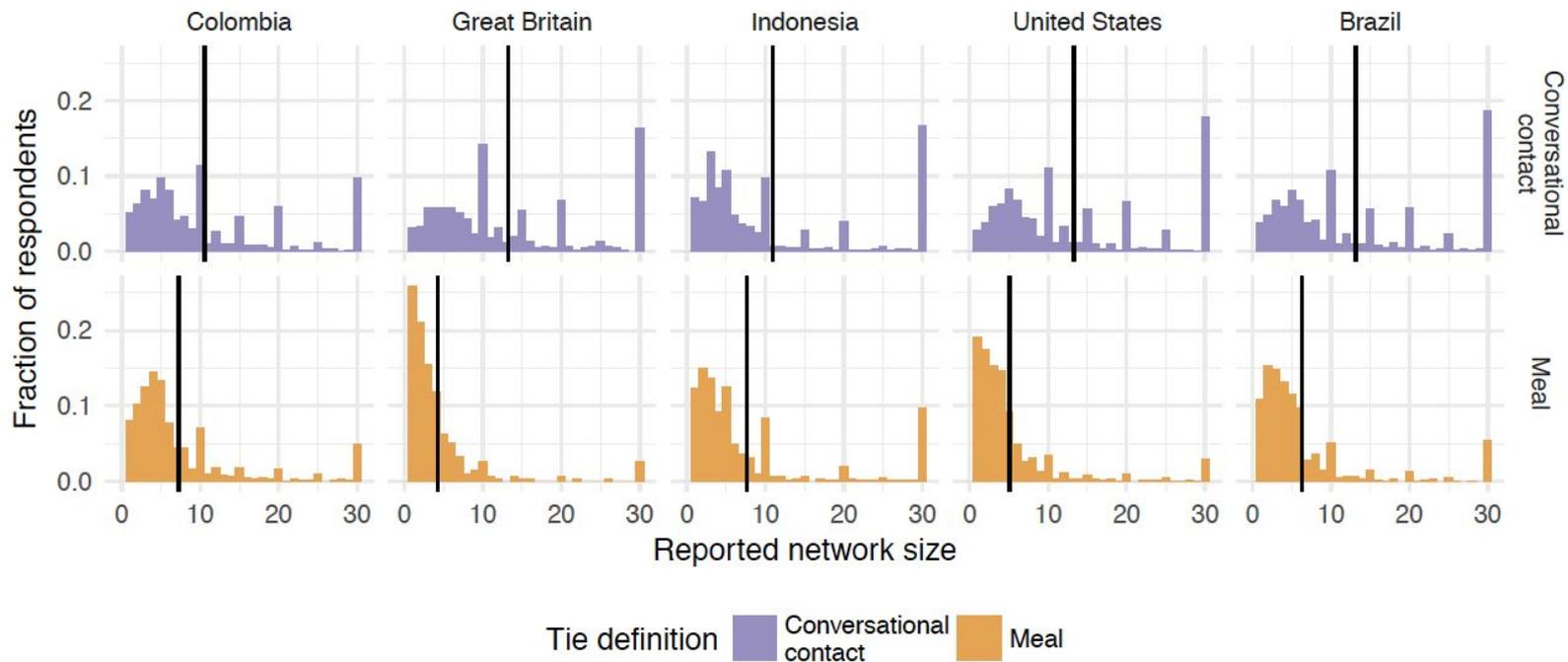
Degrees

# Degrees

Distribution of reported network sizes: Brazil  
(topcoded at 30)



## Distribution of reported network sizes (topcoded at 30)



# Internal consistency checks

# Internal consistency checks

Idea: come up with **two independent ways** of estimating the **same quantity** from network reports

Compare these independent estimates to one another

When all of the technical conditions are satisfied, estimates will agree (up to sampling noise)

Some reporting errors or other violations of conditions can be detected with IC checks

# Internal consistency checks

$$\# \text{ connections from men to women} = \# \text{ connections from women to men}$$

# Normalized difference

$$\Delta_{\alpha} = \frac{1}{N_F} (\hat{d}_{F_{-\alpha}, F_{\alpha}} - \hat{d}_{F_{\alpha}, F_{-\alpha}})$$

Example: reported connections to women made by men

Example: reported connections to men made by women

# Normalized difference

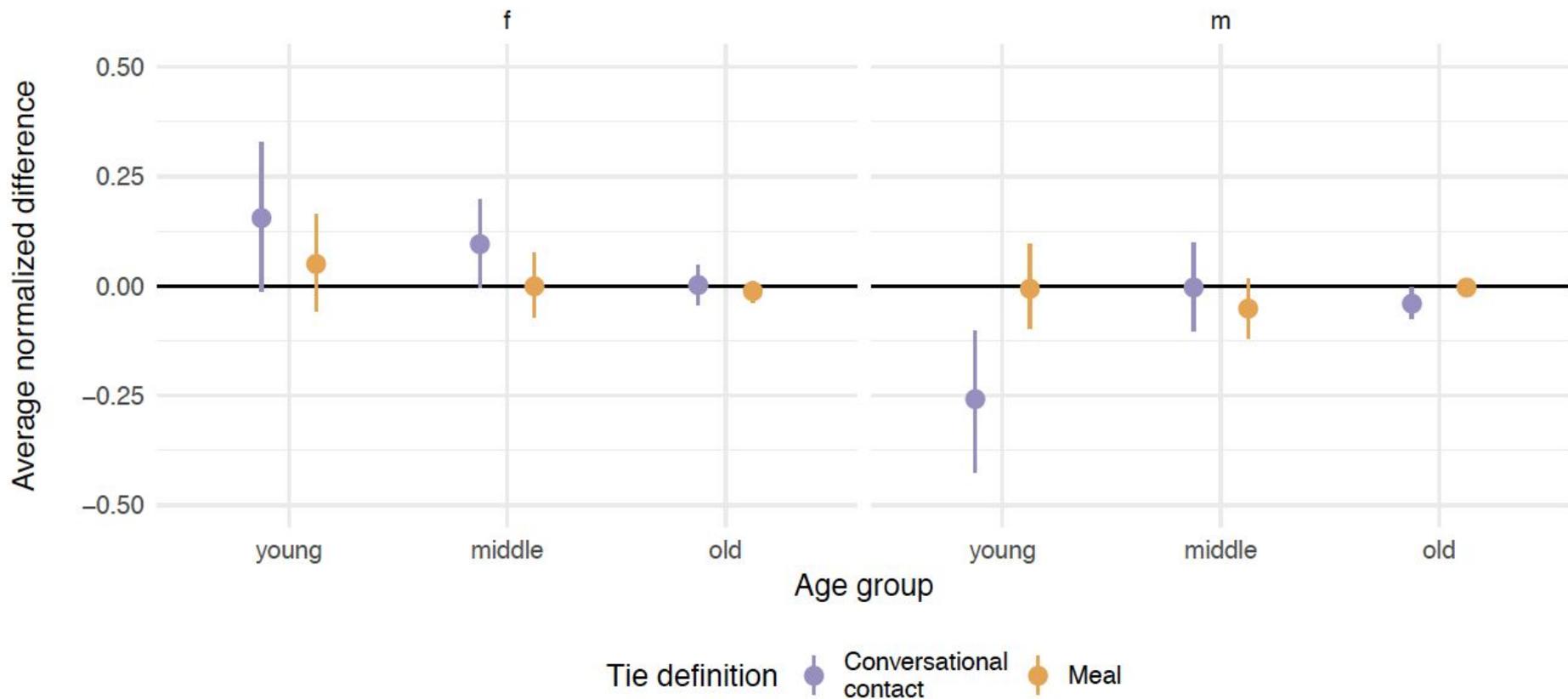
$$\Delta_{\alpha} = \frac{1}{N_F} \left( \hat{d}_{F_{-\alpha}, F_{\alpha}} - \hat{d}_{F_{\alpha}, F_{-\alpha}} \right)$$

Example: reported connections to women made by men

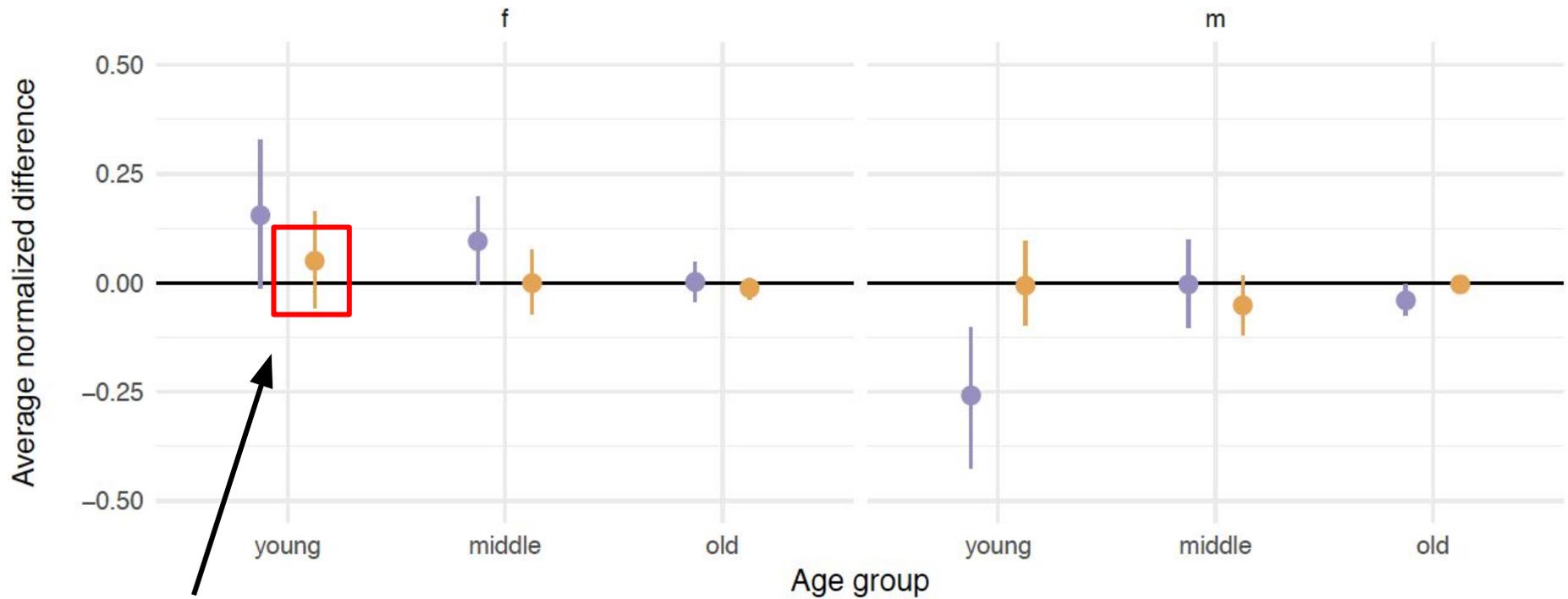
Example: reported connections to men made by women

These can be estimated **independently**  
but they are the **same quantity**

# Internal consistency checks: Brazil



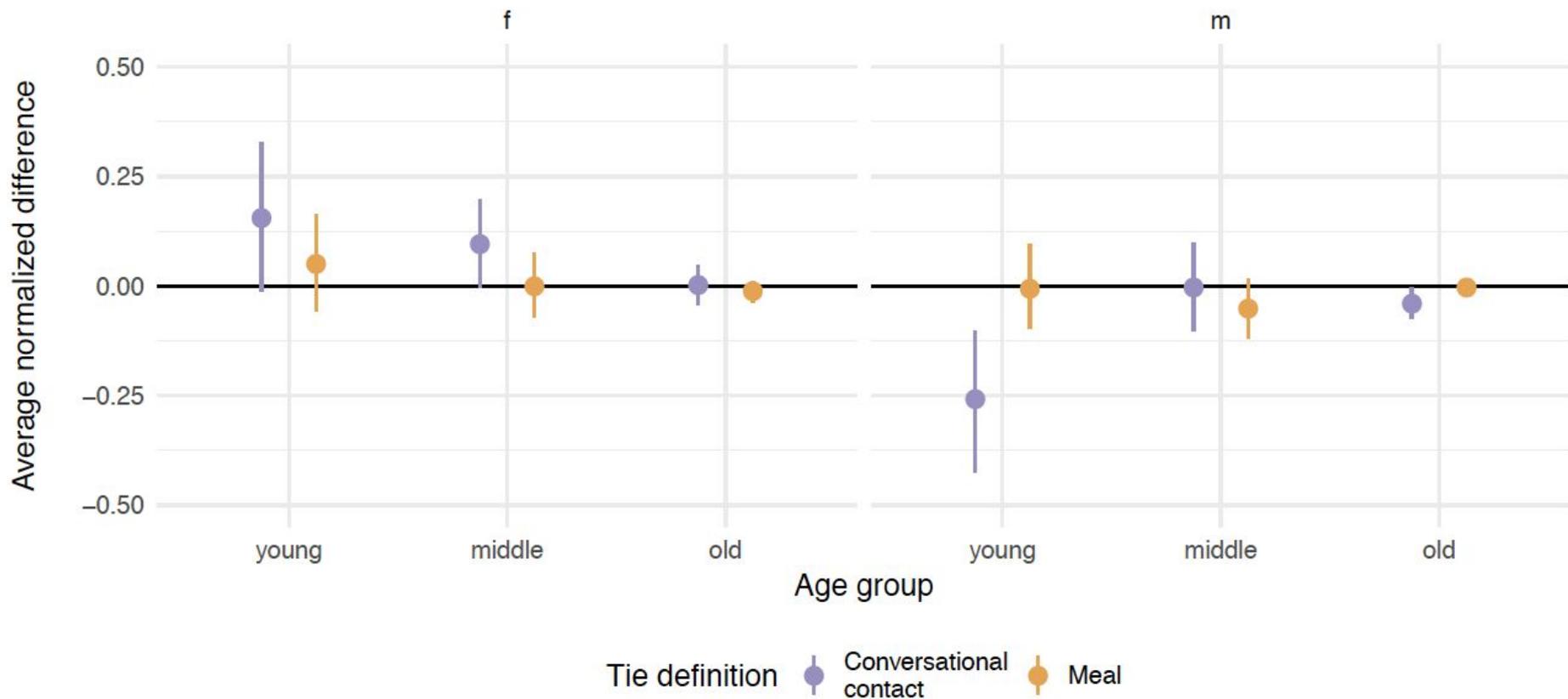
# Internal consistency checks: Brazil



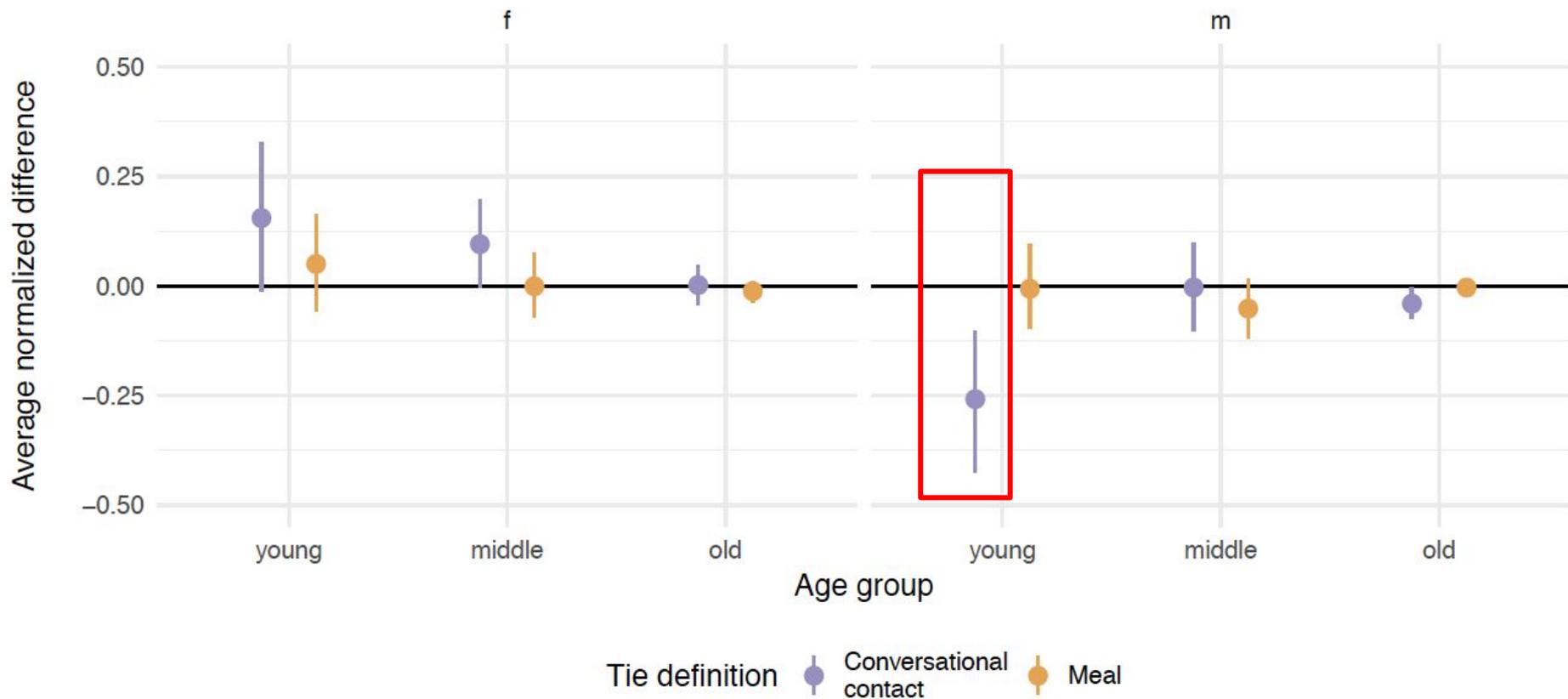
Meals everyone else reports sharing with young women MINUS meals young women report sharing with everyone else

Tie definition    ● Conversational contact    ● Meal

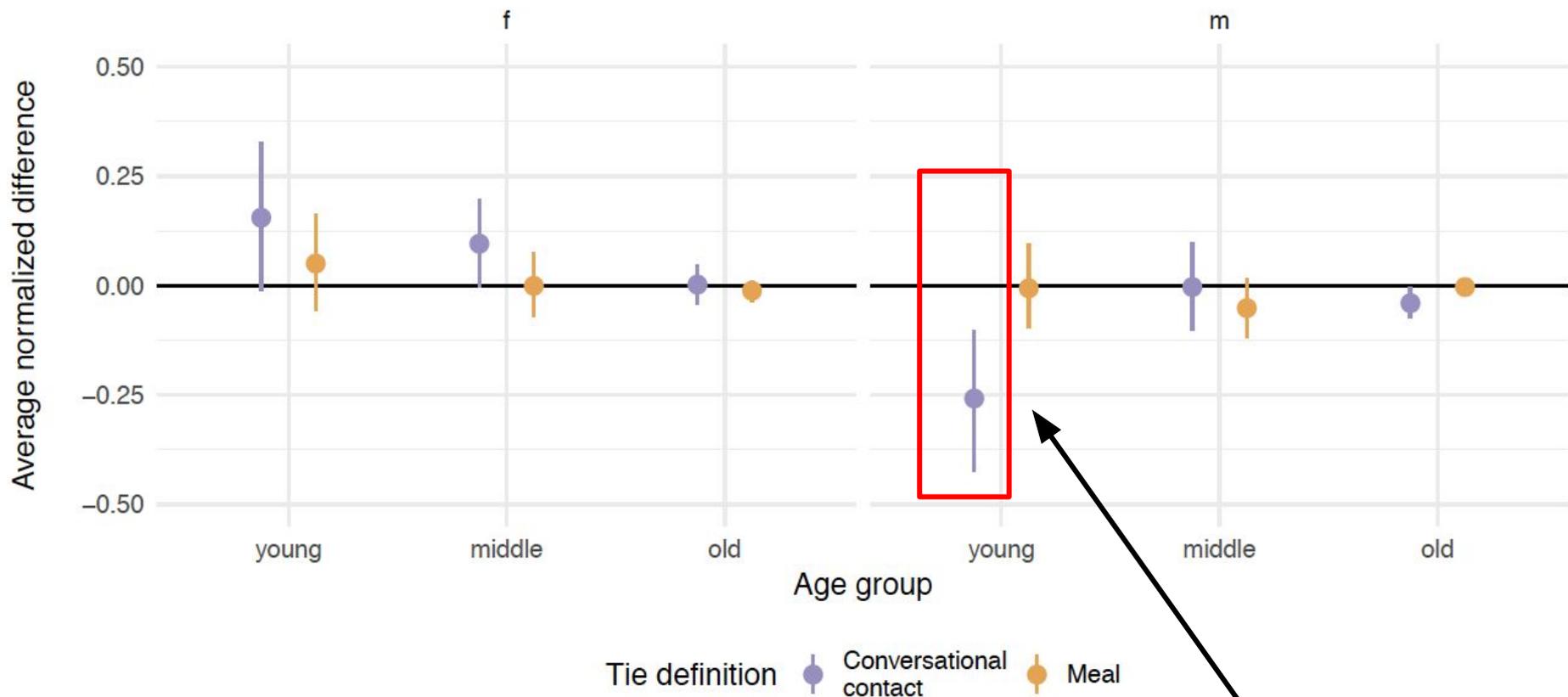
# Internal consistency checks: Brazil



# Internal consistency checks: Brazil

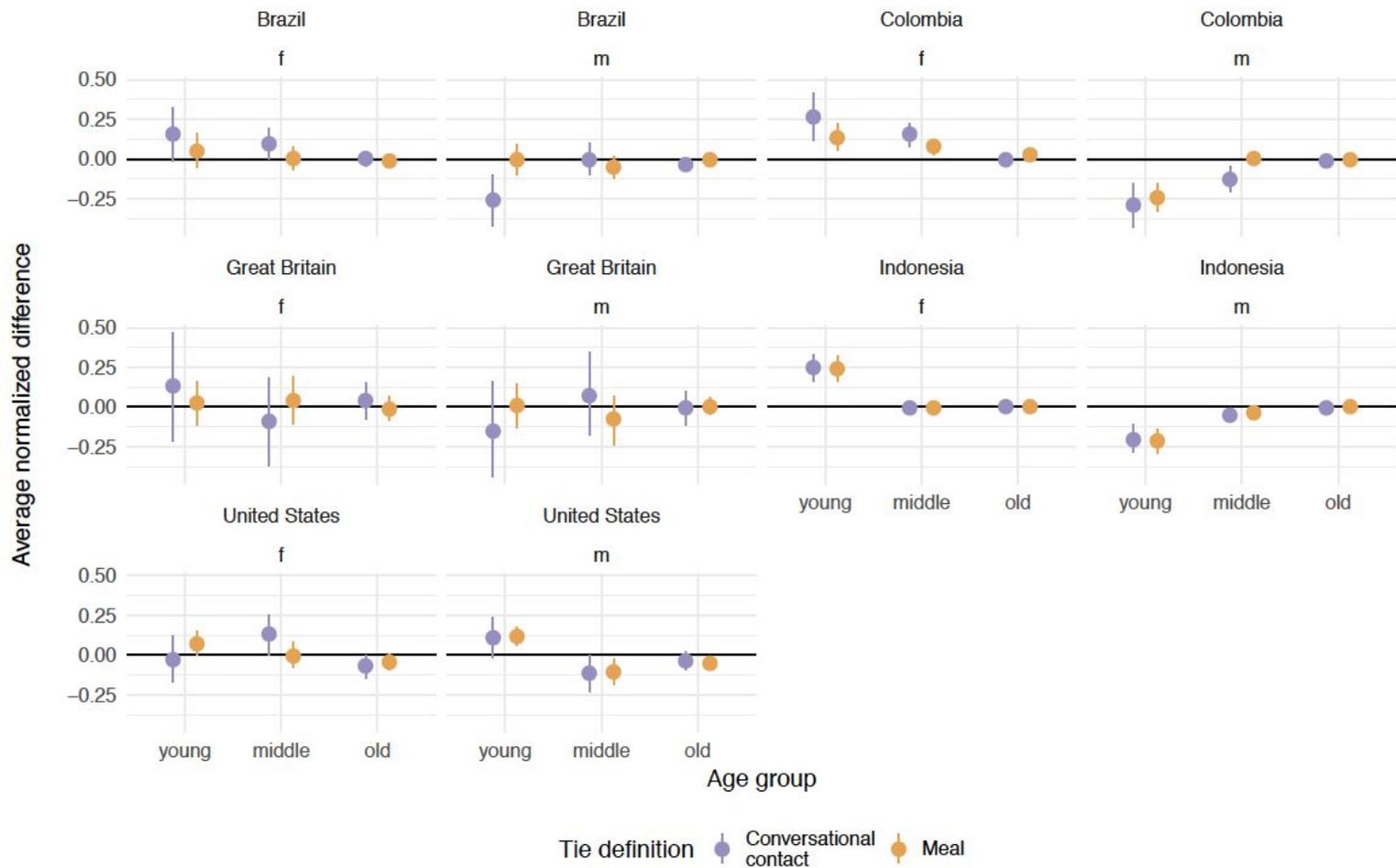


# Internal consistency checks: Brazil

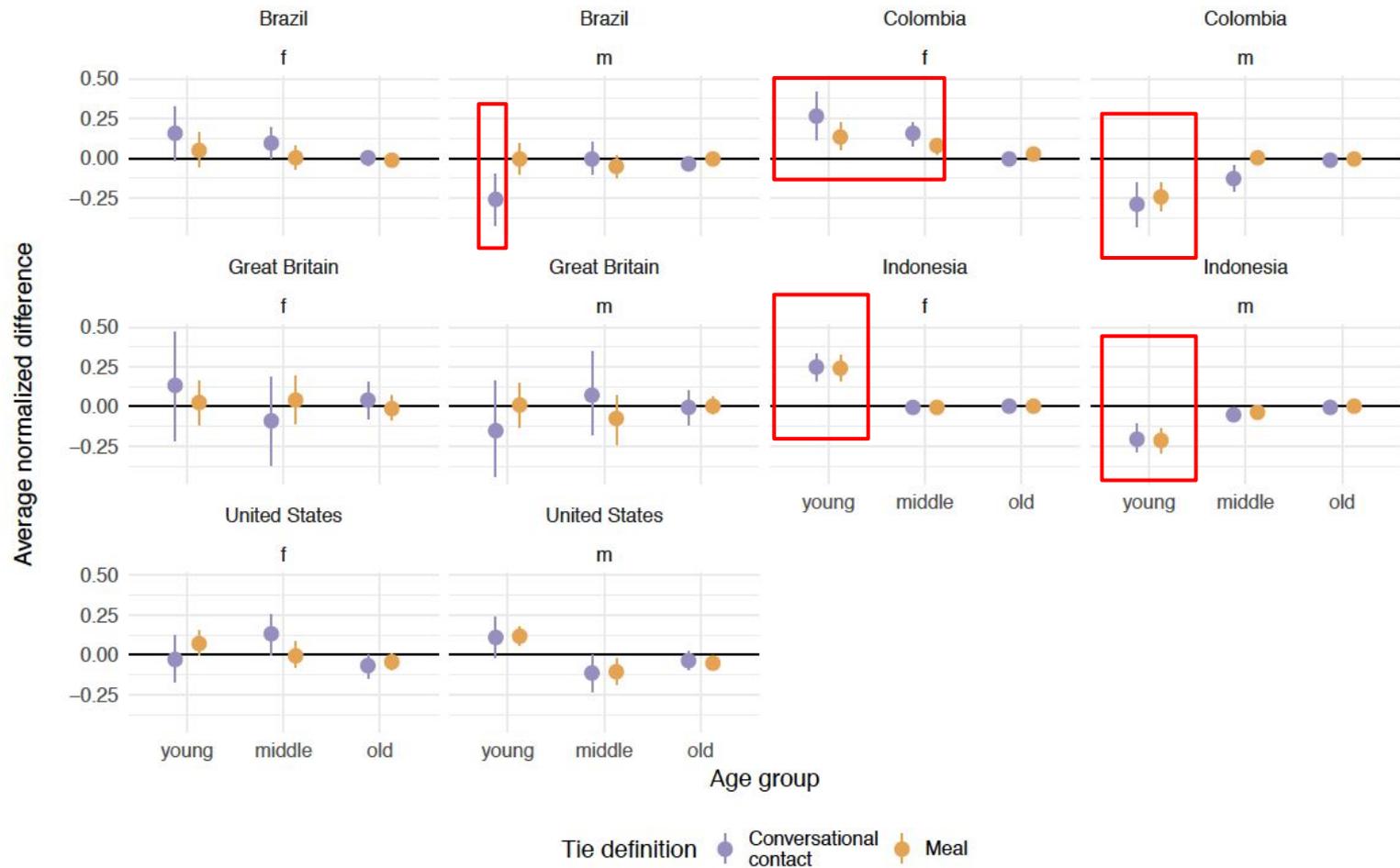


Young men report more connections to everyone else than everyone else reports to young men

# Internal consistency checks



# Internal consistency checks

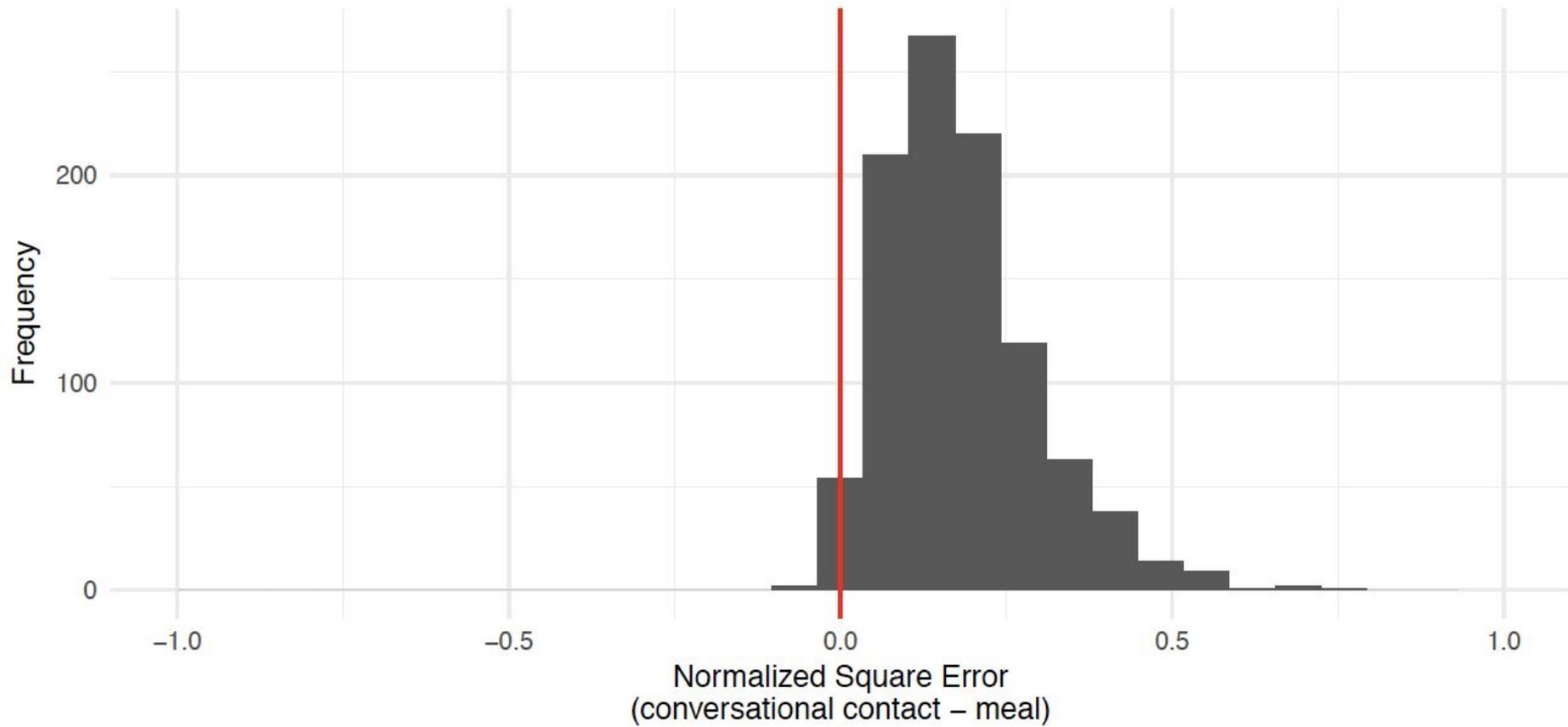


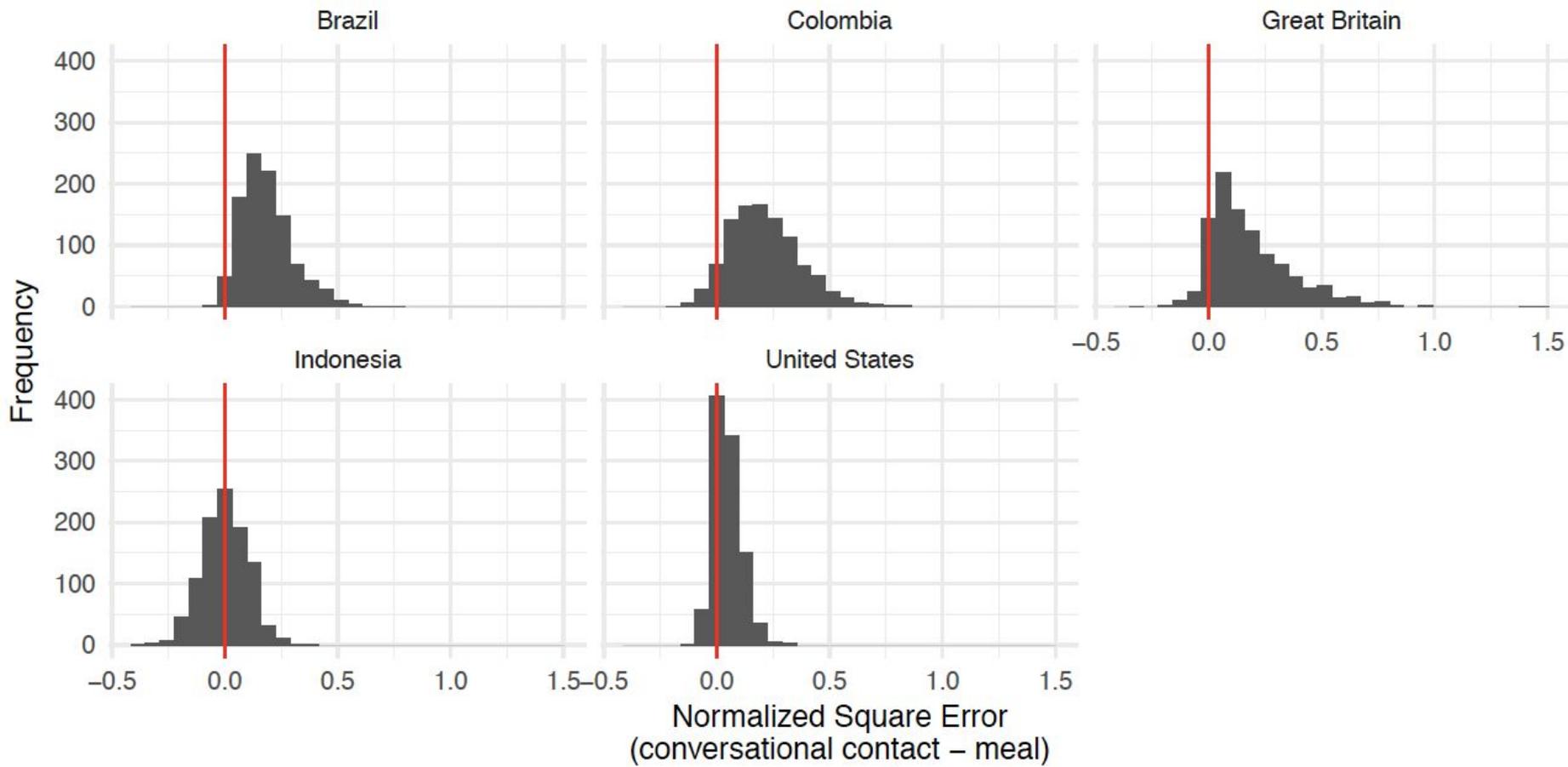
# Internal consistency checks

So the IC checks give us a way to detect when conditions are not exactly met

We can also use the IC checks to compare the two different tie definitions to better understand which one is more accurate

IC check accuracy difference (cc – meal): Brazil



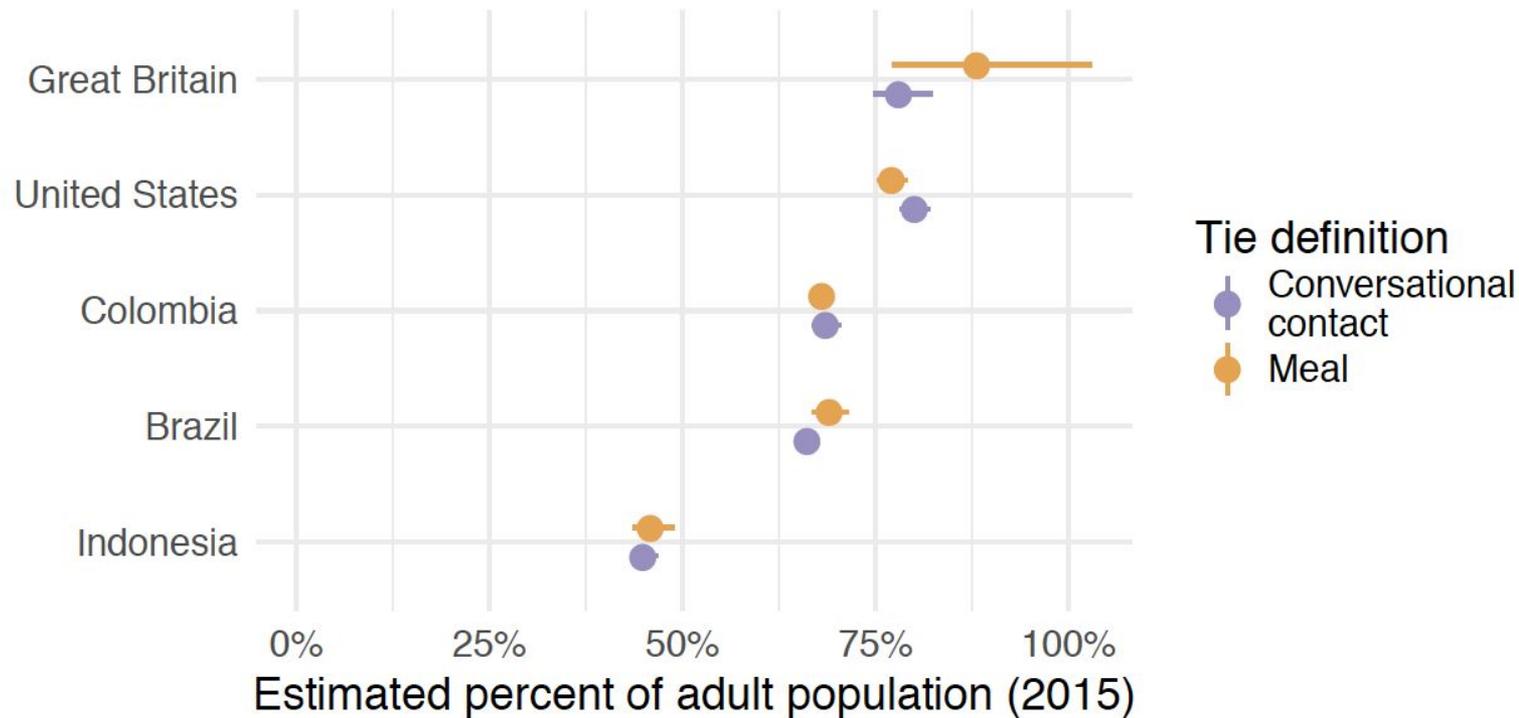


# Internal consistency checks: summary

- Built-in way to assess quality of reporting
- This is very useful for building up a picture of what kind of networks people can accurately report about
- Some evidence of reporting error (especially in Indonesia and Colombia); also suggestive of differential social visibility
- They can also form the basis for model-based approaches to improving estimates from a given network
- Results from these five countries and two networks show that meal network reports tend to be more accurate

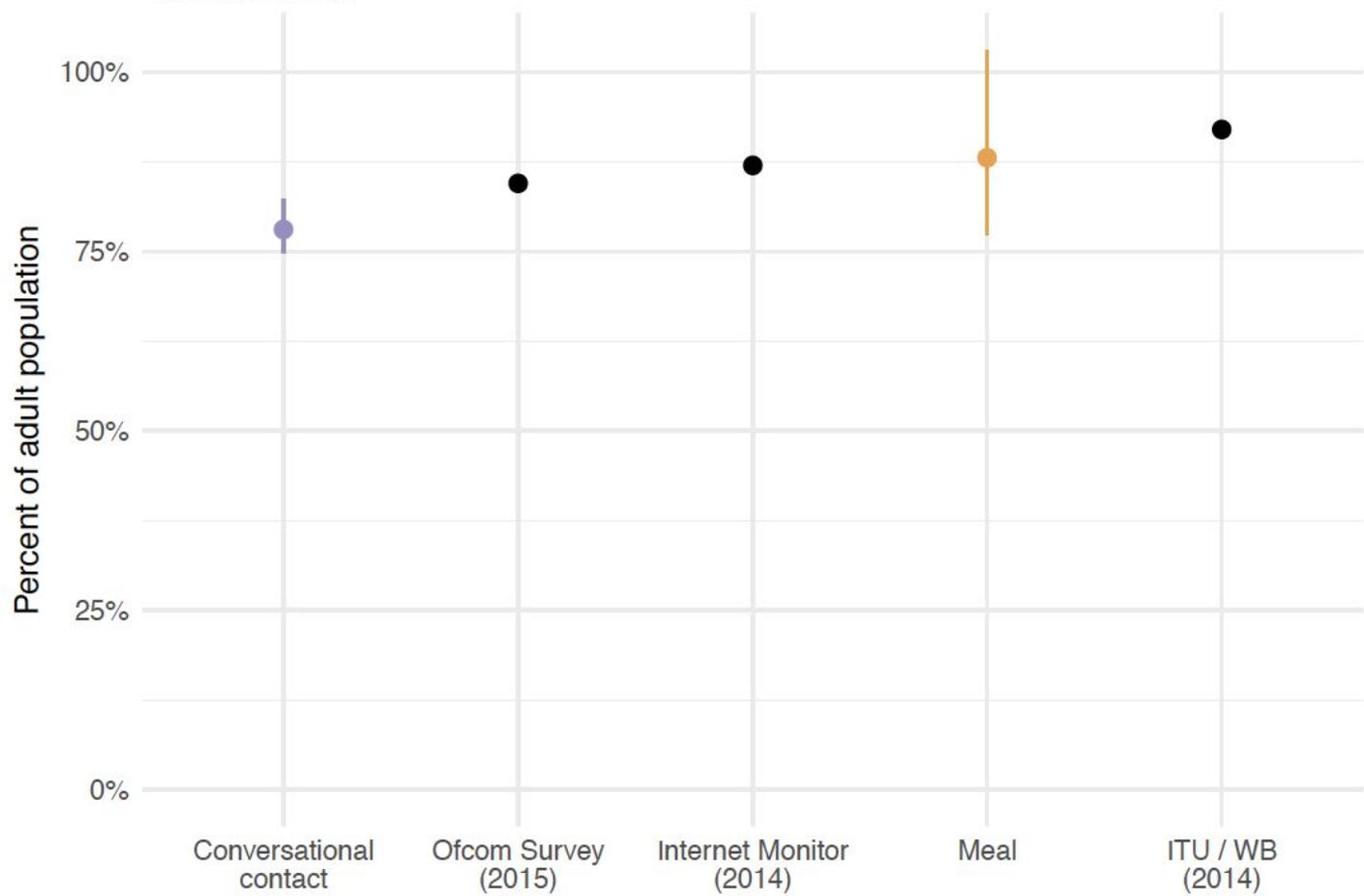
# Estimates

# Estimates

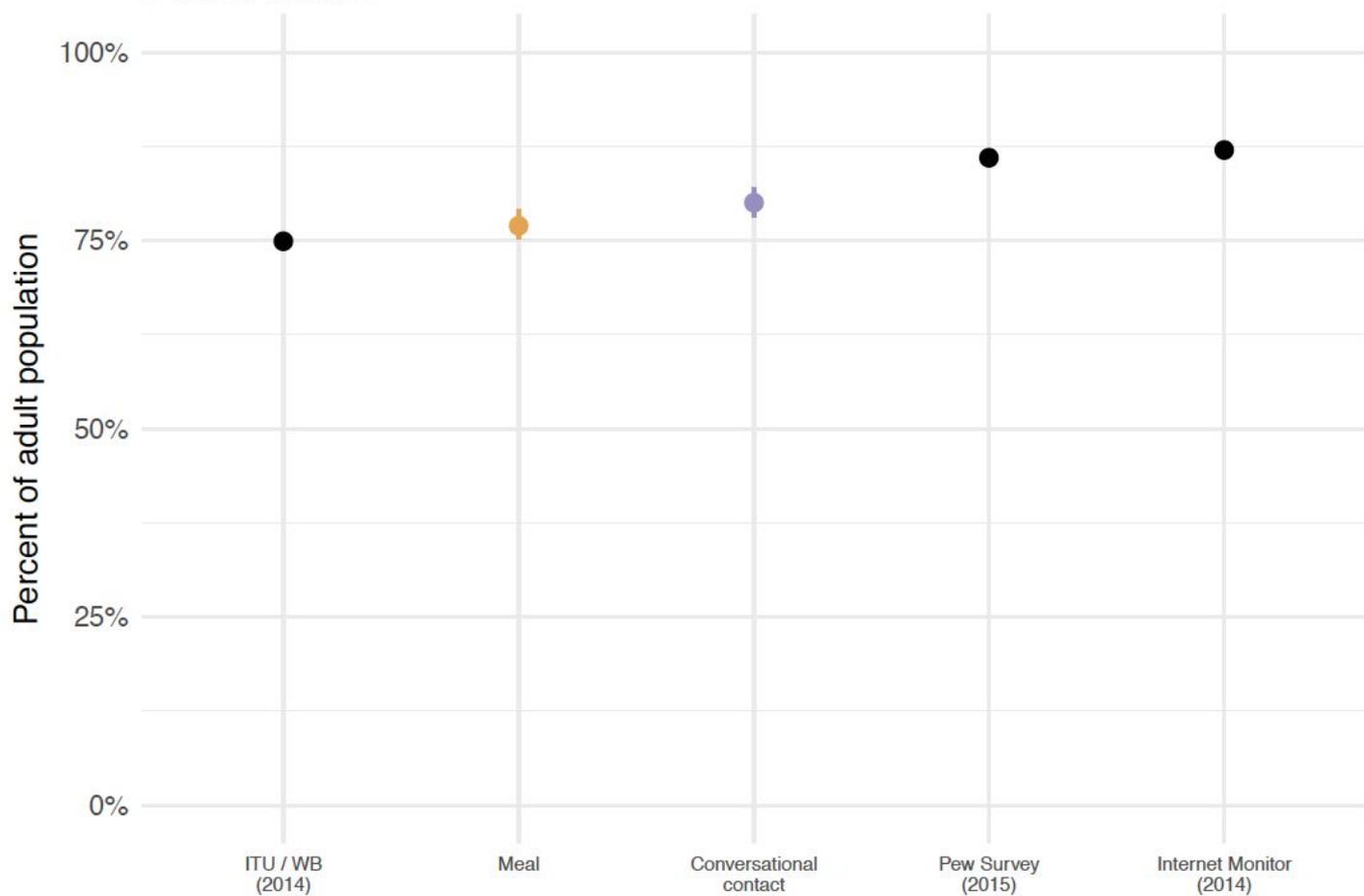


Estimates: comparisons

# Great Britain



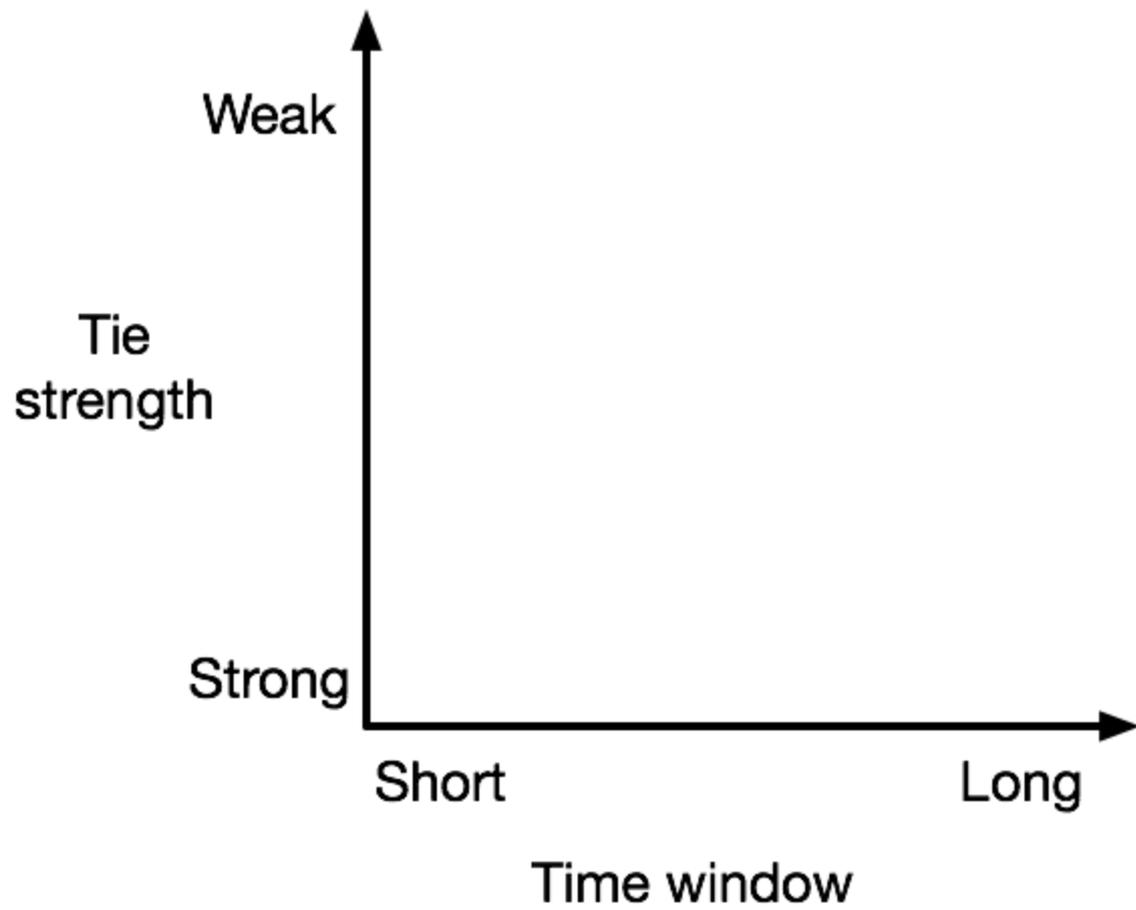
# United States

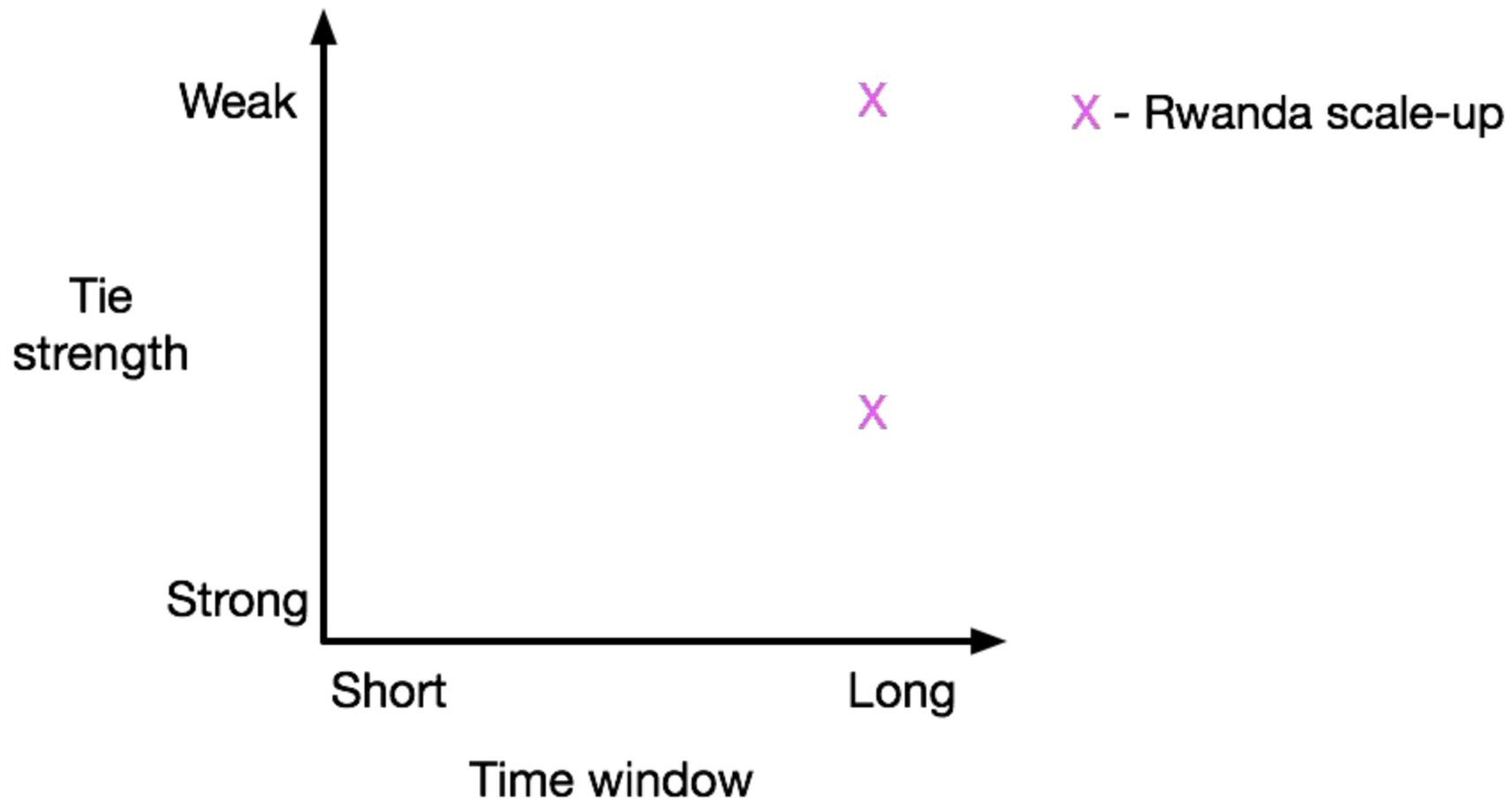


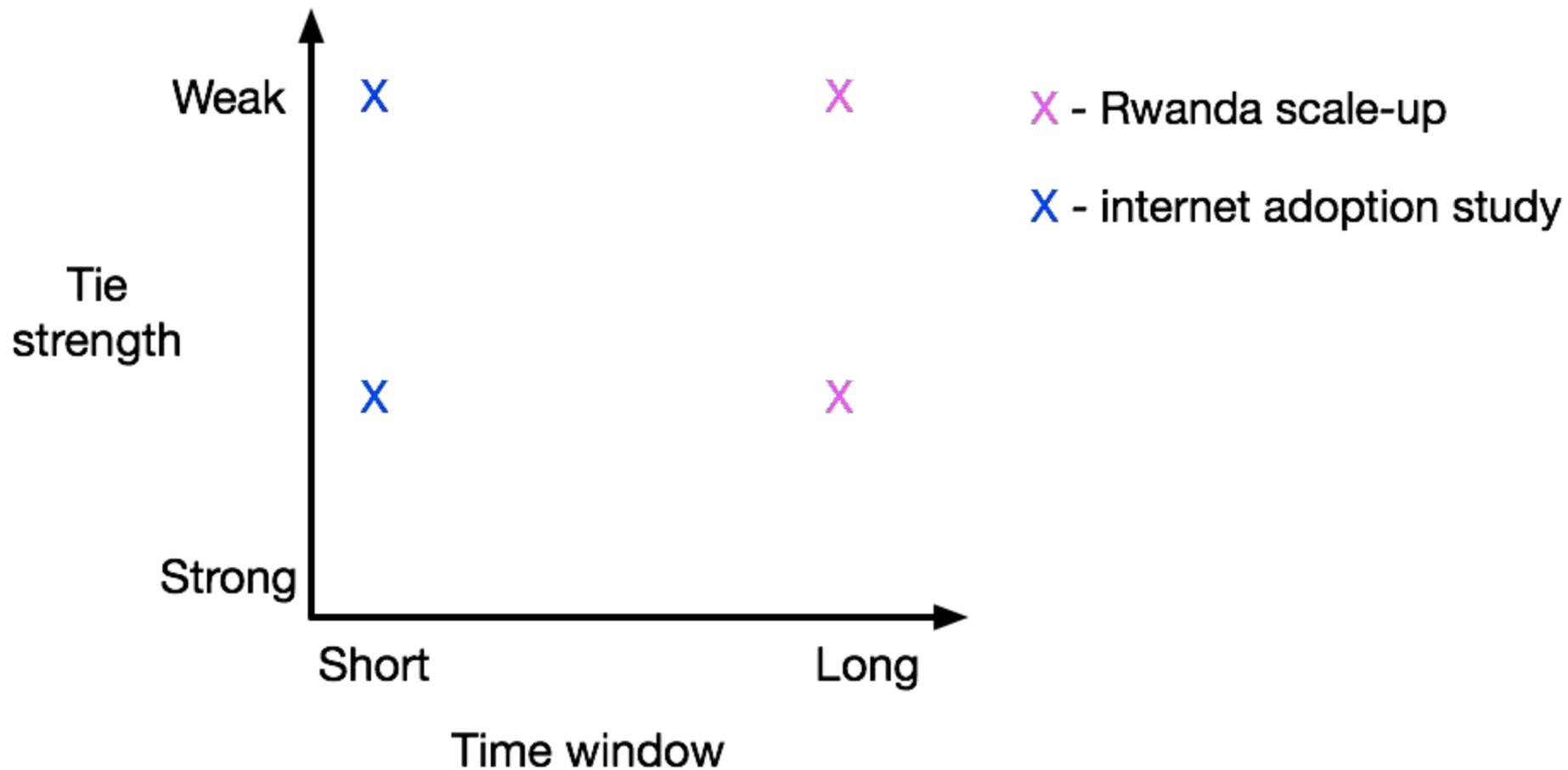
# Estimates: summary

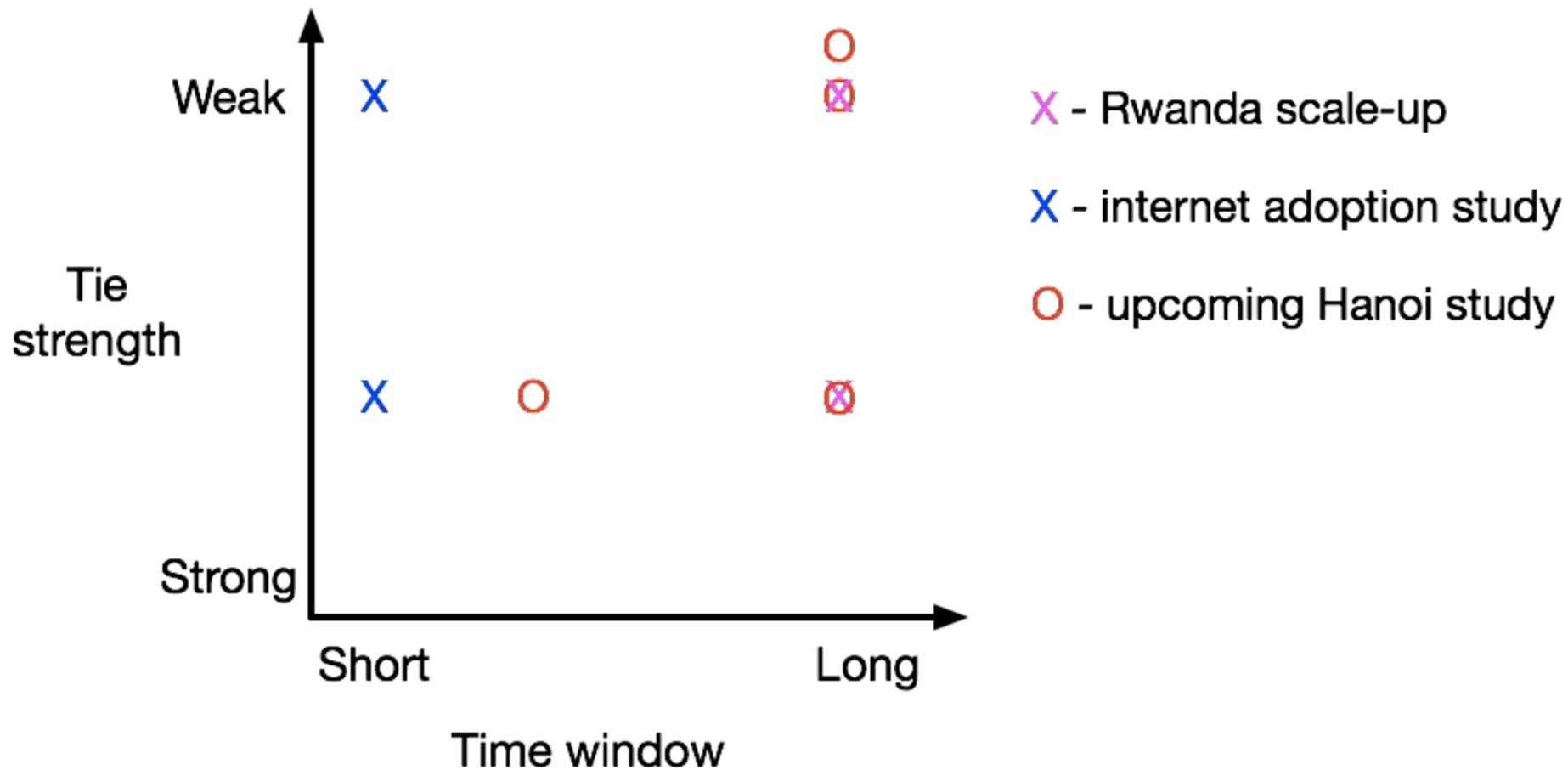
- No gold standard data to compare against, so we can't assess estimates directly
- Comparisons to other estimates in US and GB suggest our estimates are similar to other approaches, maybe slightly low
- Internal consistency checks show some evidence of reporting error (and modeling may help with this)
- Paper has sensitivity framework that can be used to formally understand what impact violating different conditions would have on estimates

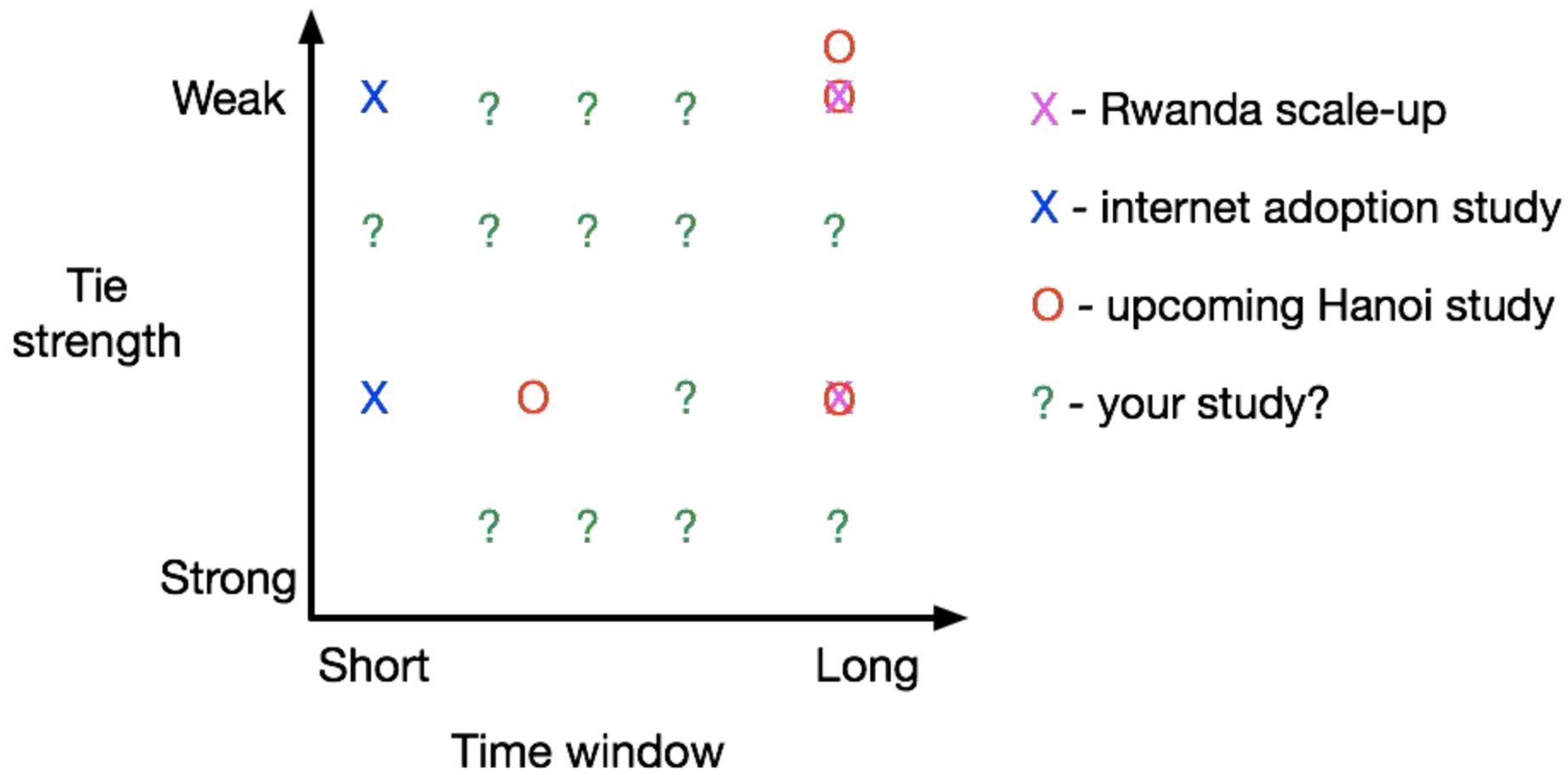
Future directions



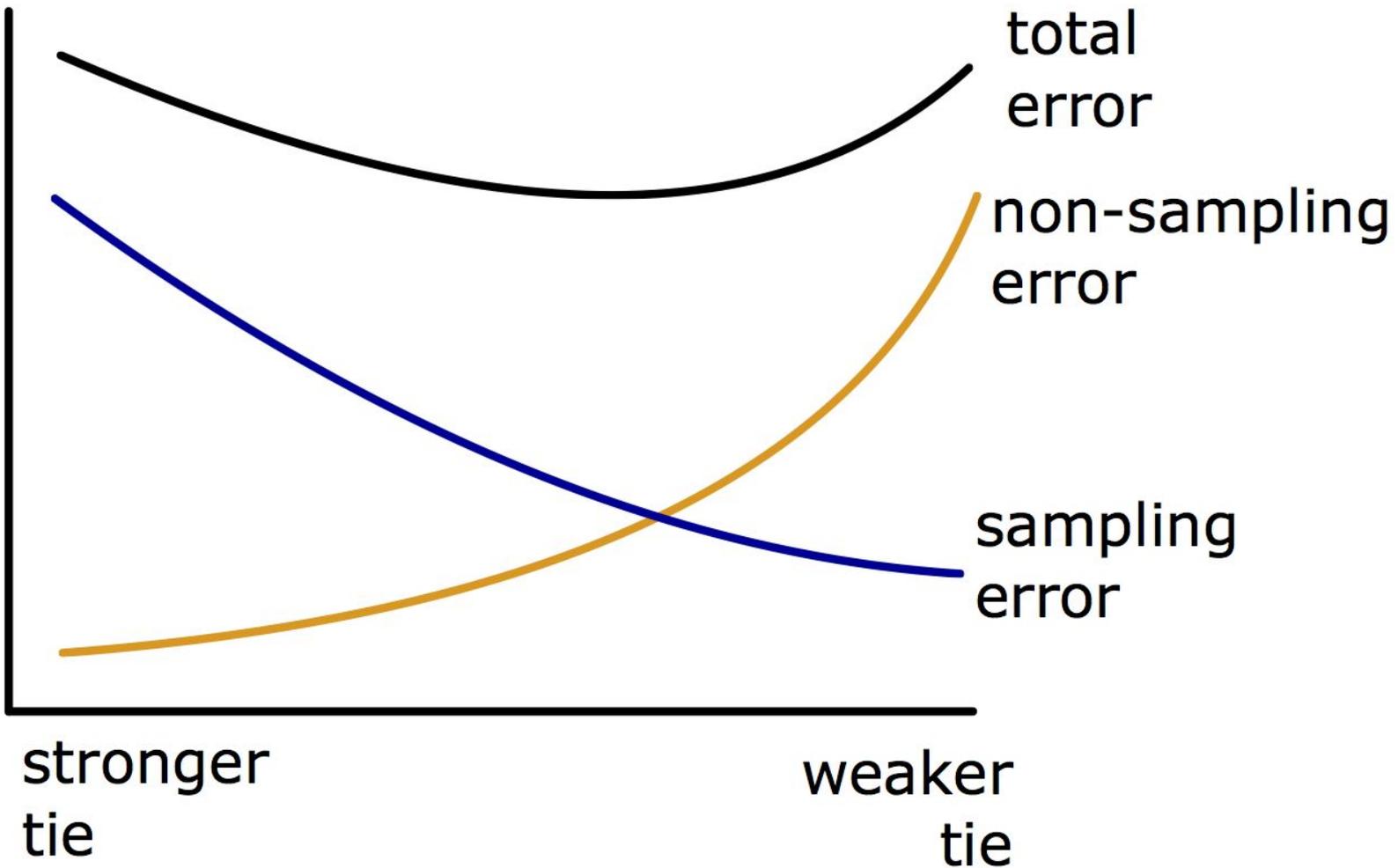








error in estimate



total error

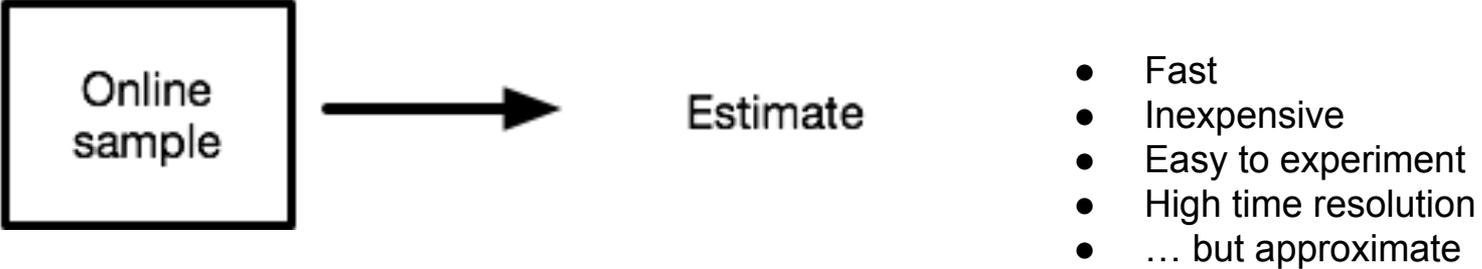
non-sampling error

sampling error

stronger tie

weaker tie

Online  
sample



Estimate

- Fast
- Inexpensive
- Easy to experiment
- High time resolution
- ... but approximate

Conventional  
probability  
sample



Estimate

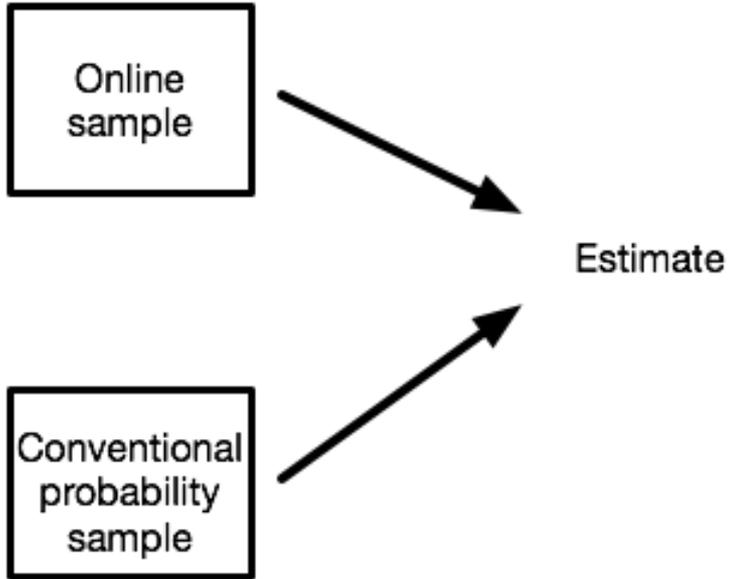
- Slow
- Expensive
- High quality estimate

Conventional  
probability  
sample

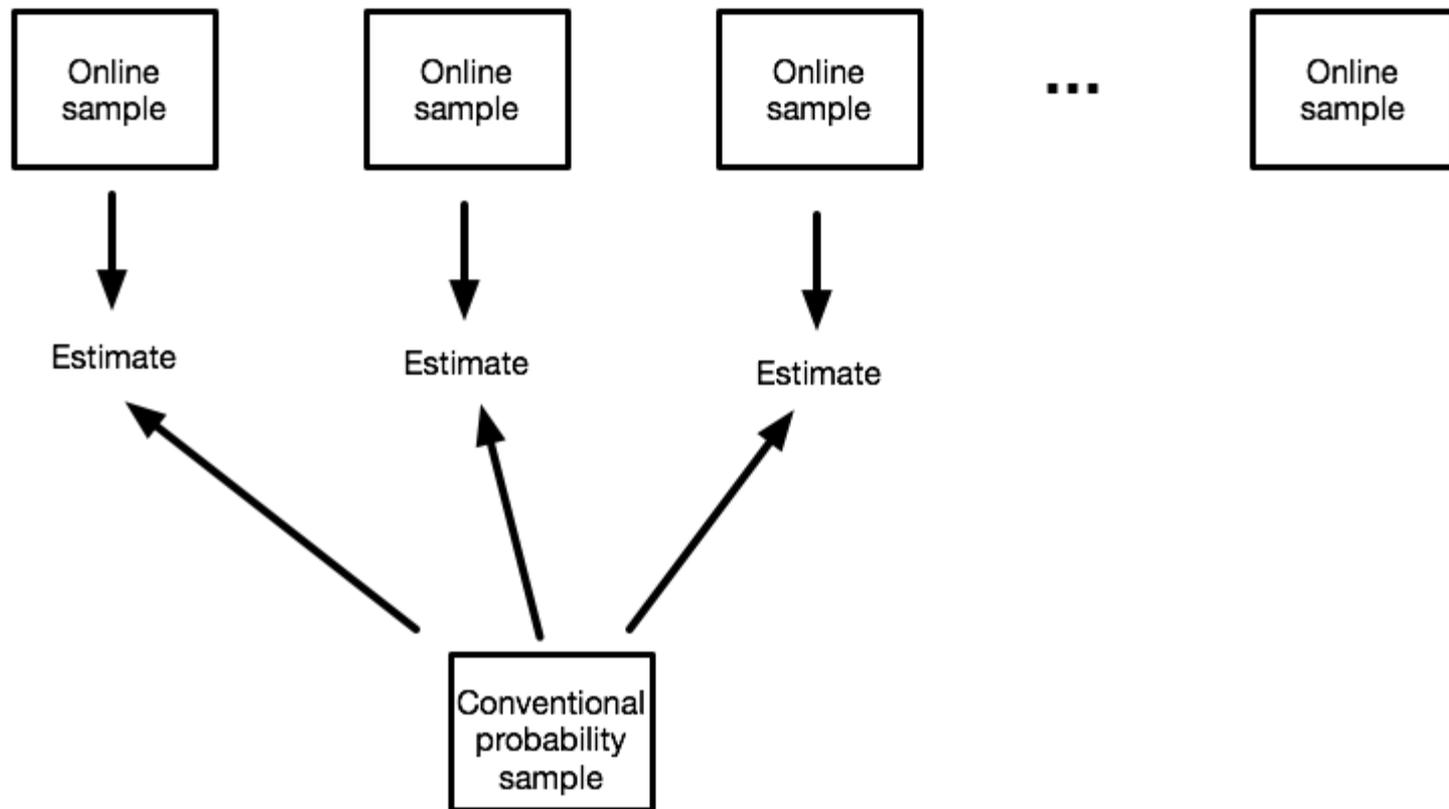


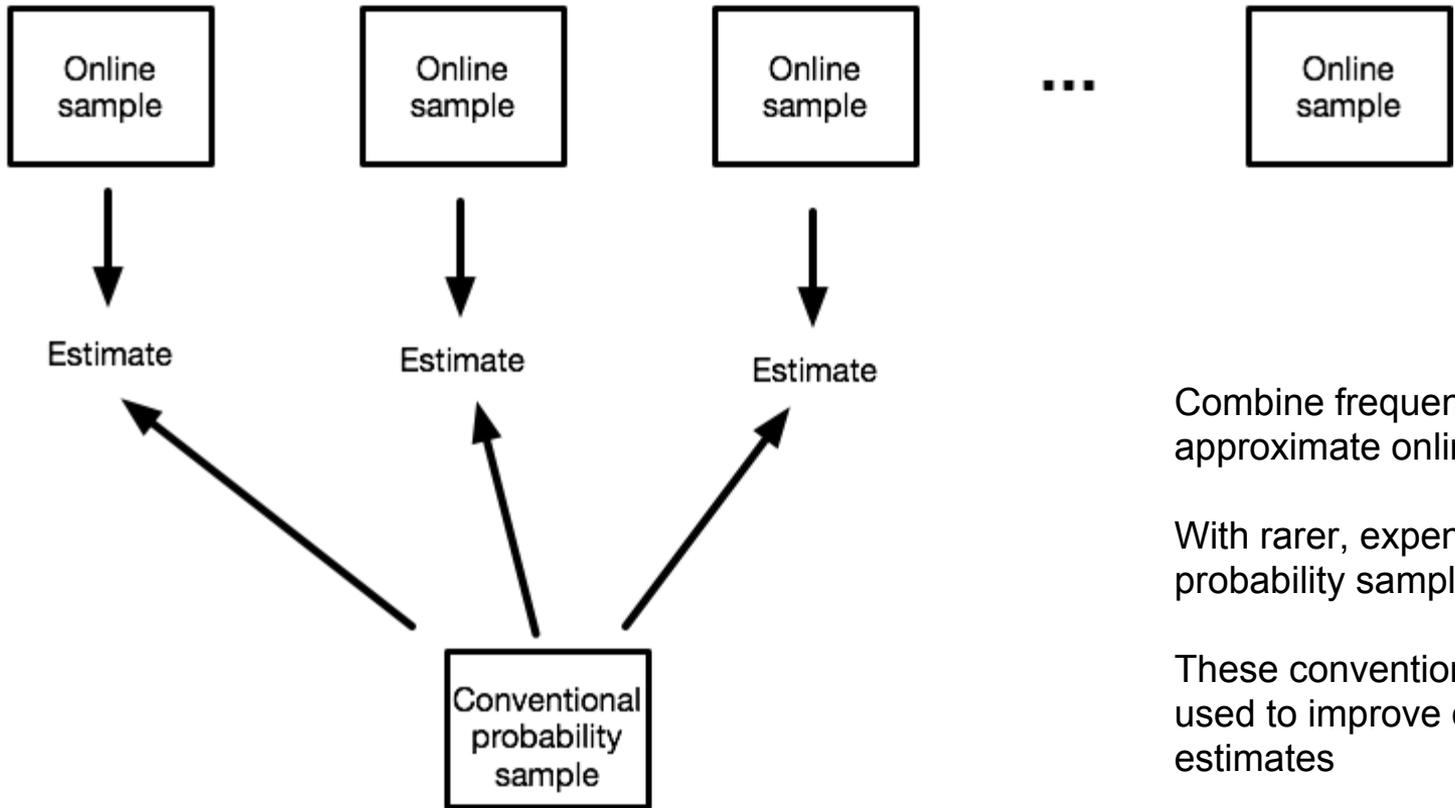
Estimate

- Slow
- Expensive
- High quality estimate
- Can collect information to help online estimates



- Network reporting framework can be used to understand how to measure things in conventional sample to improve online estimates





Combine frequent, inexpensive, approximate online-based estimates

With rarer, expensive conventional probability samples

These conventional samples can be used to improve online-based estimates

# Coming next...

- Internet adoption
  - Full sensitivity framework
  - Explore models to adjust for IC checks
  - Can also calculate estimated adoption by age and gender
  - And it's possible to do some reporting adjustments from data we collected

## Coming next...

- Internet adoption
  - Full sensitivity framework
  - Explore models to adjust for IC checks
  - Can also calculate estimated adoption by age and gender
  - And it's possible to do some reporting adjustments from data we collected
- Sibling histories (PAA 2018 session 68-4, Thurs)
- Brazil: probability sample of 25,000 respondents
  - Validate network survival methods for adult mortality
  - Test estimating out-migration using network reports
- Hanoi network scale-up for key populations at risk of HIV
- Guidance on sampling and study design

# Thanks!

- Collaborator, Curtiss Cobb
- My R packages *networkreporting* and *surveybootstrap* are available on CRAN
- Rwanda data are downloadable from the DHS website
- Feehan, Umubyeyi, Mahy, Hladik, and Salganik (2016) “Quantity vs quality: a survey experiment to improve the network scale-up method”, *American Journal of Epidemiology*
- Feehan and Salganik “Generalizing the network scaleup method”, *Sociological Methodology*.
- Feehan, Mahy, and Salganik “The network survival estimator for adult mortality: evidence from Rwanda”, *Demography*

See <http://www.dennisfeehan.org> for more information.