# Potentials and Pitfalls of Social Media Data

*Indira Sen & Katrin Weller*
GESIS workshop - December 2022

# Today's Schedule

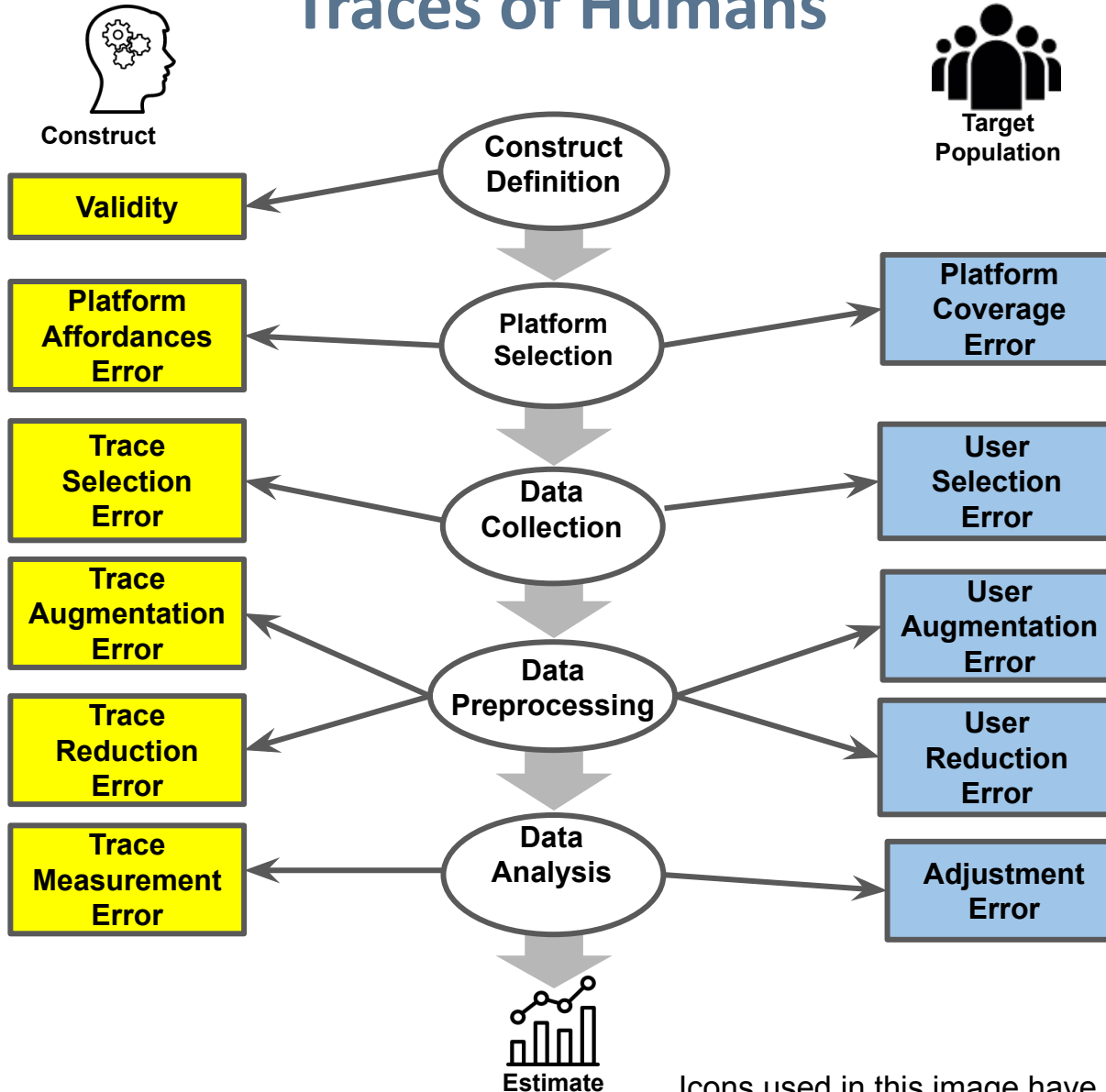| Monday, 12.12. | |
|---|---|
| 9:30-11:00 | Recap of hands-on examples from last week and exemplary research scenarios |
| 11:00-11:15 | *Break* |
| 11:15-12:00 | Potential error sources in social media research |
| 12:00-12:30 | *Break* |
| 12:30-14:00 | Hands-on error identification + Case Study: Mitigating potential error sources in social media research |

# Agenda

- Session 1: Introduction to Research with Social Media Data (SMD)
- Session 2: SM Data Collection
- Session 3: SMD Preprocessing and Analysis
- Session 4: Potential Pitfalls of SMD
- Session 5: Identifying Pitfalls with help from surveys
- **Session 6: Identifying Pitfalls in SMD**
- Session 7: Mitigating Pitfalls
- Session 8: Documenting Pitfalls
- Session 9: Recap and Conclusions

# TED-On: A Total Error Framework for Digital Traces of Humans



MEASUREMENT

Construct

REPRESENTATION

Target Population

Construct Definition

Validity

Platform Selection

Platform Affordances Error

Platform Coverage Error

Data Collection

Trace Selection Error

User Selection Error

Data Preprocessing

Trace Augmentation Error

User Augmentation Error

Trace Reduction Error

User Reduction Error

Data Analysis

Trace Measurement Error

Adjustment Error

Estimate

Icons used in this image have been designed by Becris, EliasBikbulatov and Pixel perfect from www.flaticon.com

# Example study

- How would a researcher study influenza prevalence in a national population using digital traces?

**Construct Definition**

↓

Platform Selection

↓

Data Collection

↓

Data Preprocessing

↓

Data Analysis

MEASUREMENT

Construct Definition

Platform Selection

Preprocessing

Data Analysis

TRACES

TRACES

TRACES

search queries related to flu

flu related information usage
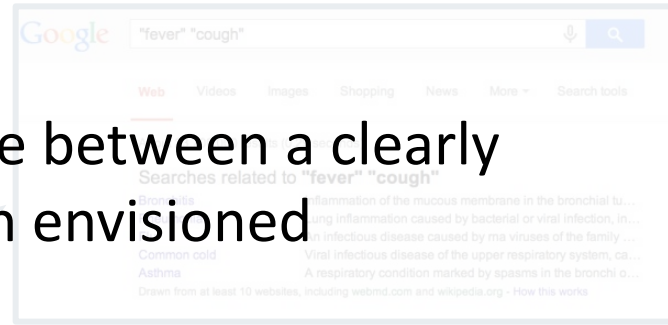
posts about having the flu

7

**Validity**

**Construct Definition**

**searching for flu related information or tweeting about flu may misrepresent incidences of flu**
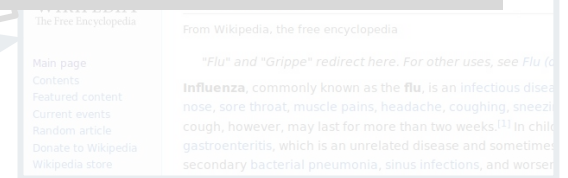
= The insufficient linkage between a clearly defined construct and an envisioned measurement

**TSE**

~ Construct definition + questionnaire design

Selection

SIGNALS

Preprocessing

SIGNALS

SIGNALS

posts about having the flu

I have the flu. Put that in your flu model (and also thanks for your patience with all the emails I'm too tired to answer 😢)

1:42 PM - 4 Jan 2019

Data Analysis

Google "fever" "cough"

**MEASUREMENT**

**Validity**

searching for flu related information or tweeting about flu may misrepresent incidences of flu

**Construct Definition**

**Platform Selection**

Data Collection

Data Preprocessing

Data Analysis

**MEASUREMENT**

**Validity**

searching for flu related information or tweeting about flu may misrepresent incidences of flu

**Construct Definition**

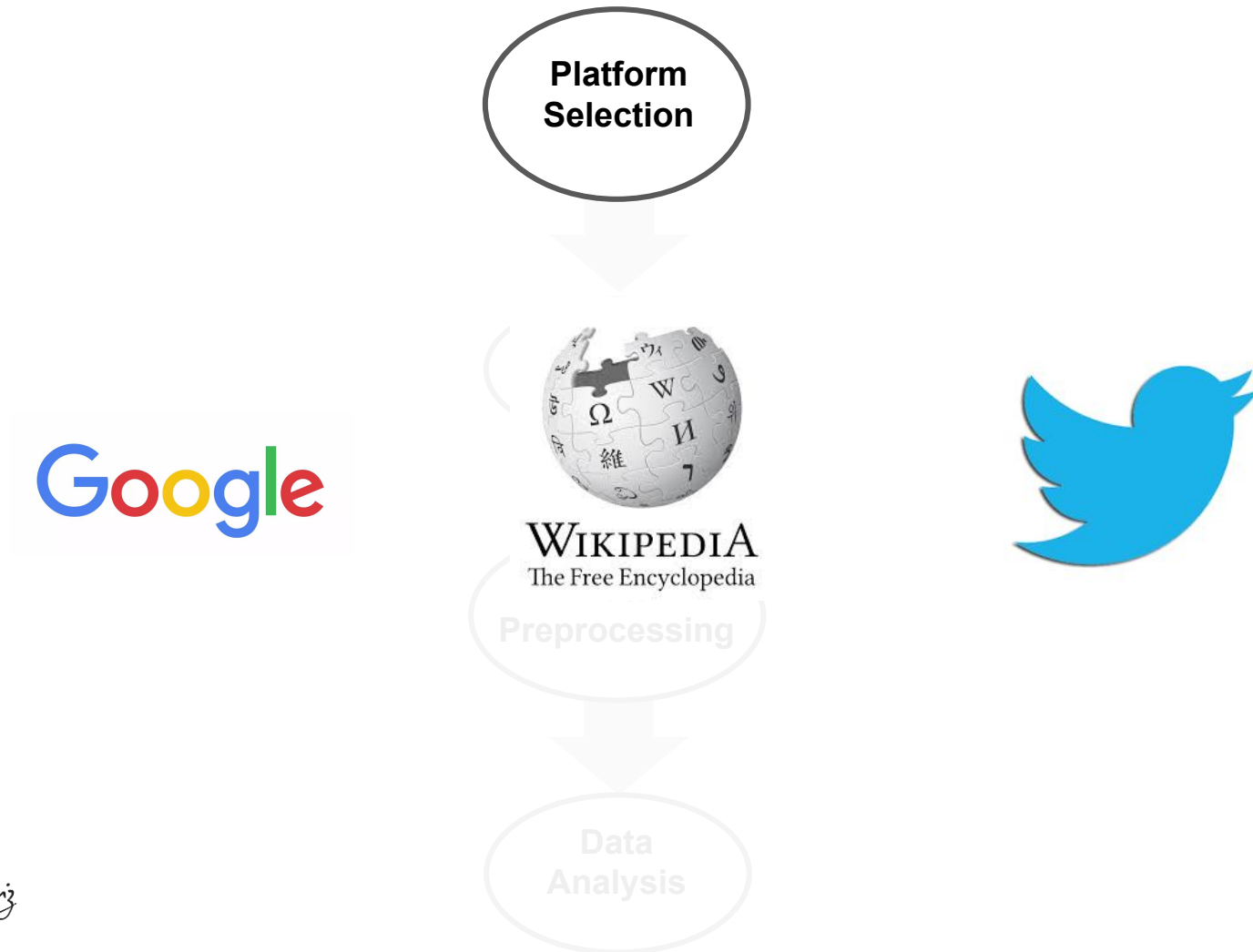**Platform Selection**

Data Collection

Data Preprocessing

Data Analysis

Construct
Definition

Platform
Affordances
Error

Platform
Selection

Distorting the
expression of
traces

Recording &
Sharing
traces

Preprocessing

Recommendations for search queries

Only aggregate queries

Session length for readers

Data
Analysis

280-character limit

Trending tweets

12

= The gap between the 'true' traces and traces distorted by platform affordances - technical/community standards, terms & conditions

**Platform Affordances Error**

**Platform(s) have affordances which distort traces**

**Platform Selection**

**TSE**
~ Reponse collection (collection mode): Response Error

Recommendations for search queries

Only aggregate queries

Session length for readers

280-character limit

Trending tweets

REPRESENTATION

Construct Definition

Platform Selection

Platform Coverage Error

Platform population = people who access Google and Wikipedia

Platform population = people with Twitter accounts

14

= The gap between the **target population (say US Population)** and the **platform population**

Platform Coverage Error

Platform Selection

Platform population = people who access Google and Wikipedia

Platform population = people with Twitter accounts

= The gap between the **target population (say US Population)** and the **platform population**

**Platform Coverage Error**

**Platform Selection**

**Platform(s) are not representative of national target population**

**TSE**
- Sampling Frame Selection: Coverage Error

Platform population = people who access Google and Wikipedia

Platform population = people with Twitter accounts

Construct Definition

Platform Selection



I'm open to correction, but it's my understanding that government already pays for all influenza vaccine that comes into the country (the contract with negotiated at the national level), so I don't think cost is a significant barrier to uptake amongst HCW

2:50 AM · Oct 9, 2019 · Twitter for iPhone

**Data Collection**

Data Preprocessing

**Retrieve tweets with keywords** related to Influenza: Avian influenza, Influenza Virus B, Centers for Disease Control and Prevention, Influenza Virus C, Common Cold, Vaccine, Flu(the Band), Influenza (English tweets only!)

Analysis

Construct
Definition



= The gap between the ideal response and the measured trace

```
Trace
Selection Error
```
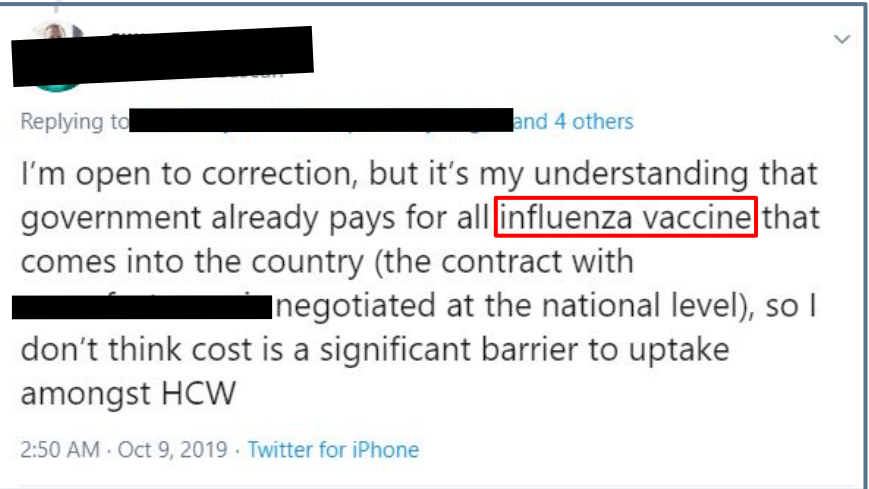
Data
Collection

Data
Preprocessing

**Retrieve tweets with keywords** related to Influenza: Avian influenza, Influenza Virus B, Centers for Disease Control and Prevention, Influenza Virus C, Common Cold, Vaccine, Flu(the Band), Influenza (English tweets only!)

Analysis

19

Construct
Definition

= The gap between the ideal
response and the measured
trace

I'm open to correction, but it's my understanding that government already pays for all influenza vaccine that comes into the country (the contract with negotiated at the national level), so I don't think cost is a significant barrier to uptake amongst HCW

2:50 AM · Oct 9, 2019 · Twitter for iPhone

**Trace
Selection Error**

**Data
Collection**

**Some tweets chosen
may not be relevant for
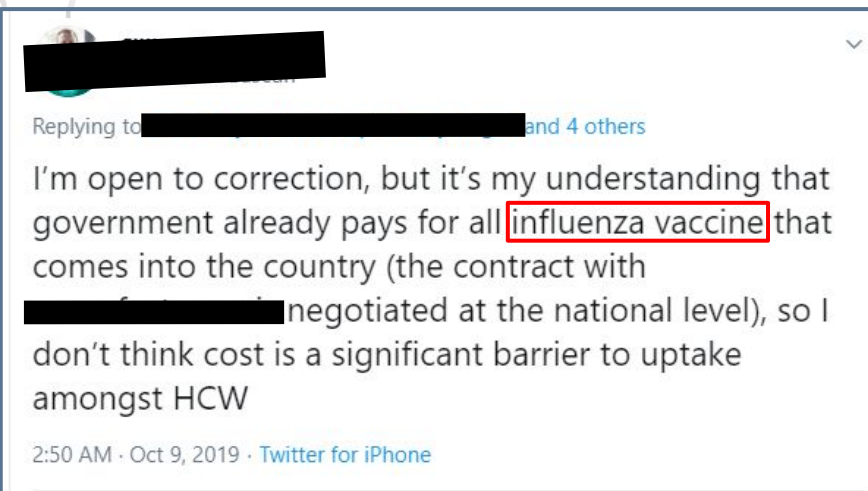influenza like illnesses**

**TSE**

~ Response collection:
response error

Data
Preprocessing

**Retrieve tweets with keywords** related to Influenza: Avian
influenza, Influenza Virus B, Centers for Disease Control
and Prevention, Influenza Virus C, Common Cold, Vaccine,
Flu(the Band), Influenza a (English tweets only!)

Analysis

Construct
Definition

Platform
Selection

**Data
Collection**

Data
Preprocessing

**Retrieve tweets with keywords** related to Influenza: Avian influenza, Influenza Virus B, Centers for Disease Control and Prevention, Influenza Virus C, Common Cold, Vaccine, Flu(the Band), Influenza (English tweets only!)

Analysis

Construct
Definition

= The difference between the platform population and the users chosen due to the query specification

Platform
Selection

**Data Collection** → **User Selection Error**

Data
Preprocessing

**Retrieve tweets with keywords** related to Influenza: Avian influenza, Influenza Virus B, Centers for Disease Control and Prevention, Influenza Virus C, Common Cold, Vaccine, Flu(the Band), Influenza (English tweets only!)
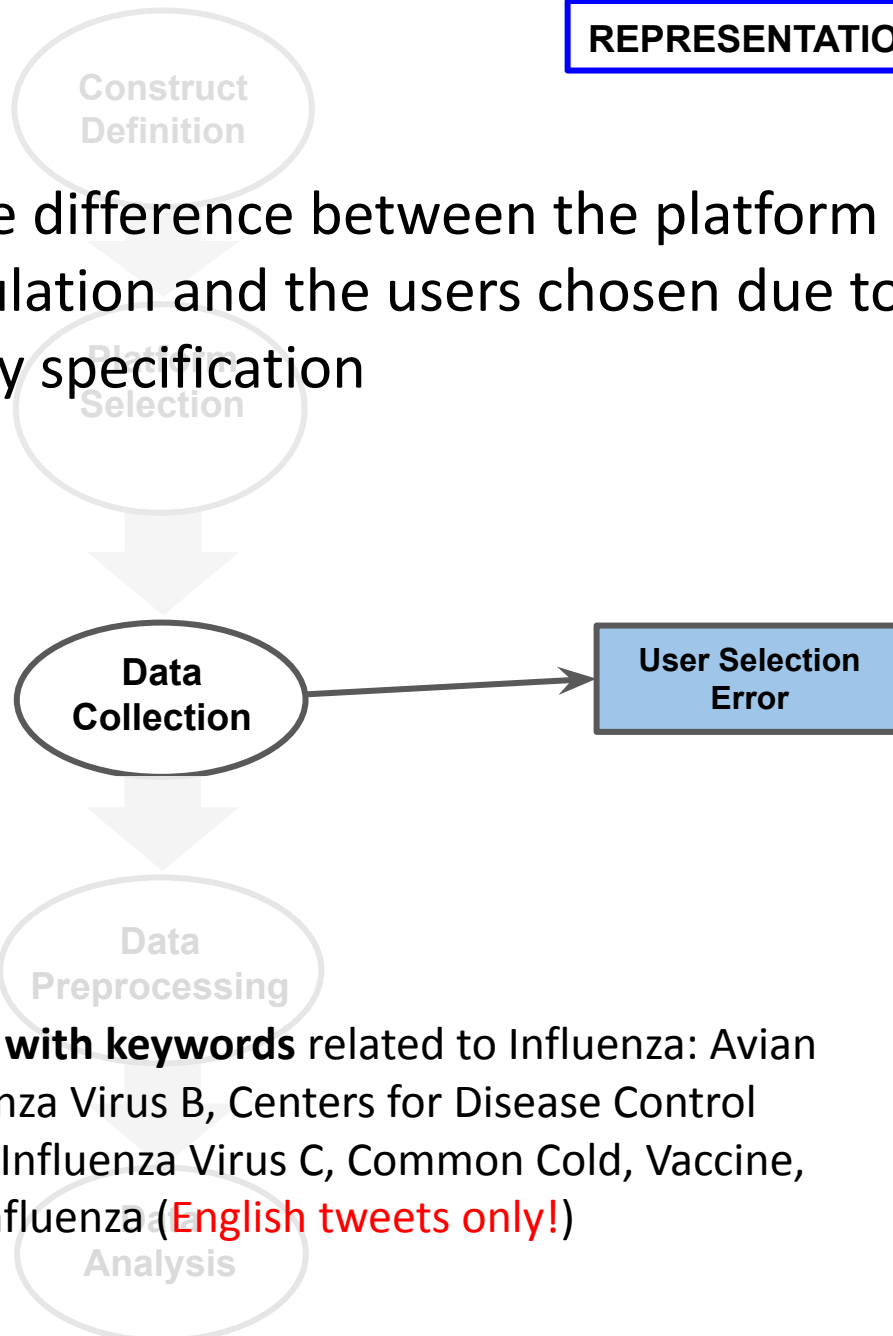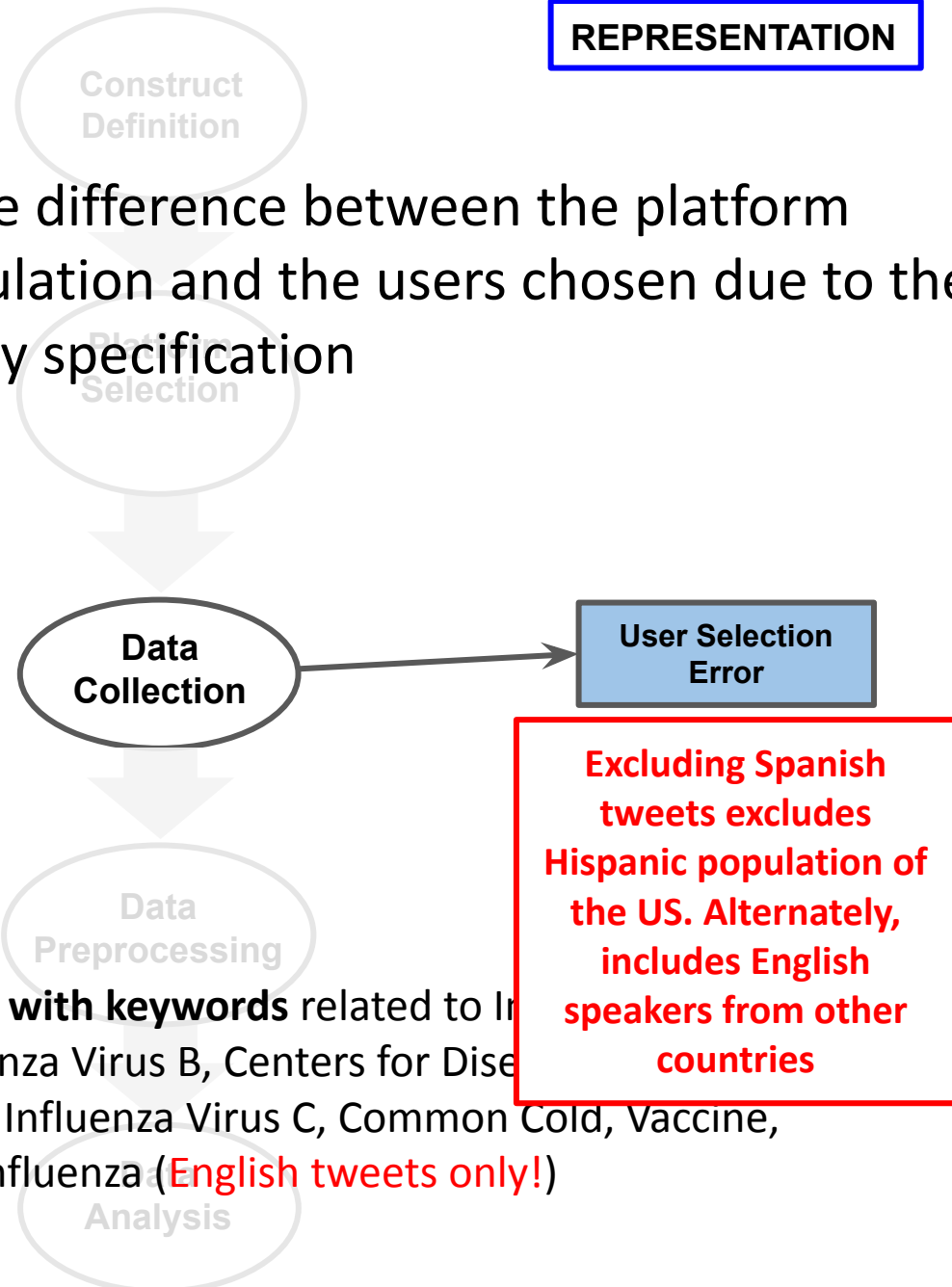
Analysis

22

**Construct Definition**

**Platform Selection**

= The difference between the platform population and the users chosen due to the query specification

**TSE**

~   Sampling: Sampling error

~   Coverage error

**Data Collection**

**User Selection Error**

**Excluding Spanish tweets excludes Hispanic population of the US. Alternately, includes English speakers from other countries**

**Data Preprocessing**

**Retrieve tweets with keywords** related to Influenza, Influenza Virus B, Centers for Disease and Prevention, Influenza Virus C, Common Cold, Vaccine, Flu(the Band), Influenza (English tweets only!)
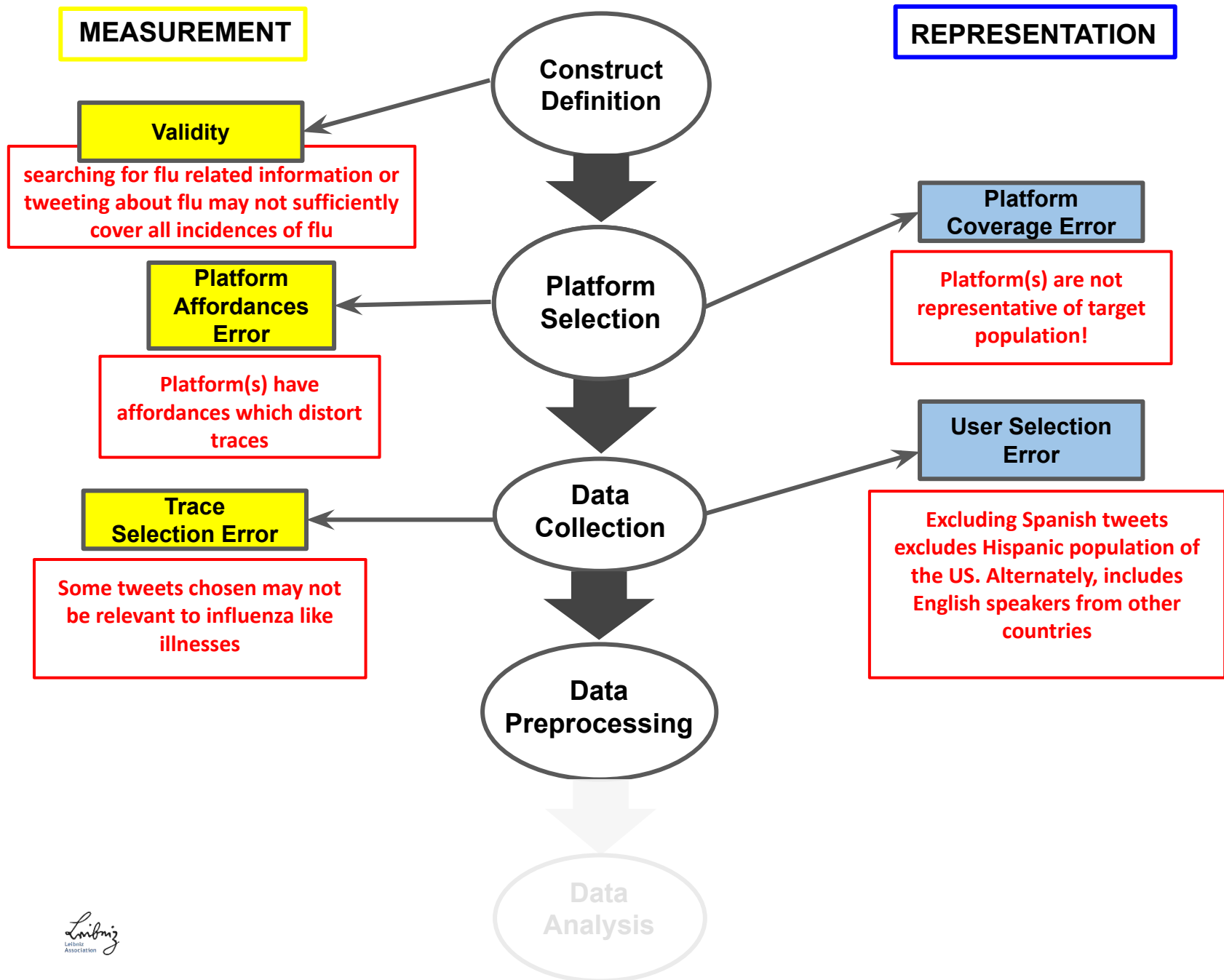
**Analysis**

**MEASUREMENT**

**REPRESENTATION**

**Construct Definition**

**Validity**

searching for flu related information or tweeting about flu may not sufficiently cover all incidences of flu

**Platform Selection**

**Platform Coverage Error**

Platform(s) are not representative of target population!

**Platform Affordances Error**

Platform(s) have affordances which distort traces

**Data Collection**

**User Selection Error**

Excluding Spanish tweets excludes Hispanic population of the US. Alternately, includes English speakers from other countries

**Trace Selection Error**

Some tweets chosen may not be relevant to influenza like illnesses

**Data Preprocessing**

Data Analysis

**MEASUREMENT**



Construct
Definition

I'm ████████████ draw that since I feel sick
██████ Probably gonna go sleep now.

**Augment traces** with syntactic relations
indicating the speaker being afflicted by the flu

Data
Collection

**Data
Preprocessing**

Data
Analysis

25

**Construct Definition**

**Platform Selection**

**Remove** non 'first-person' tweets

> Faisal Hussein
>
> Wuhan Flu is dangerous and all, ▮▮▮▮▮▮ also dangerous for ▮▮▮▮▮▮
>
> Our tendency to believe & amplify fake news & unsubstantiated rumors. Fake news in general may cause unnecessary panic, create communal tensions etc. It can destroy the fabric of our society.

**Data Collection**

**Data Preprocessing**

**Data Analysis**

Construct
Definition

Selection



I'm ▮▮▮▮▮▮▮▮▮▮▮ draw that since I feel sick ▮▮▮▮▮▮ Probably gonna go sleep now.

Wuhan Flu is dangerous and all, ▮▮▮▮▮▮ also dangerous for ▮▮▮▮▮

Our tendency ▮ believe ▮ amplify fake news & unsubstantiate ▮▮▮▮▮. Fake news in general may cause unnecessar ▮ ▮ic, create communal tensions etc. It can destr ▮▮▮ bric of our society.

**Augment traces** with syntactic relations indicating the speaker being afflicted by the flu

**Remove** non 'first-person' tweets

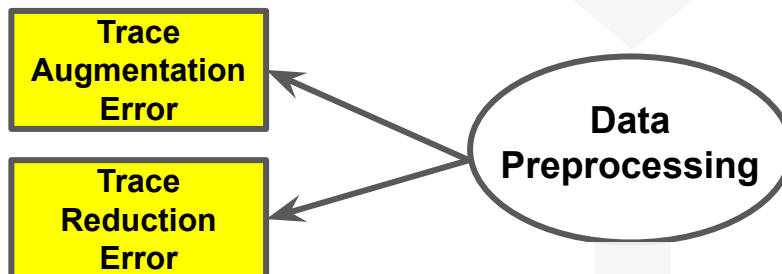= Errors due to **augmenting** or **filtering** out traces

Data Collection

| Trace Augmentation Error |
| Trace Reduction Error |

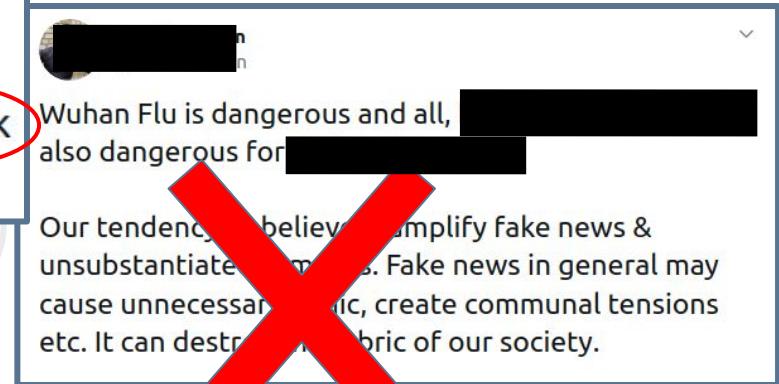Data Preprocessing

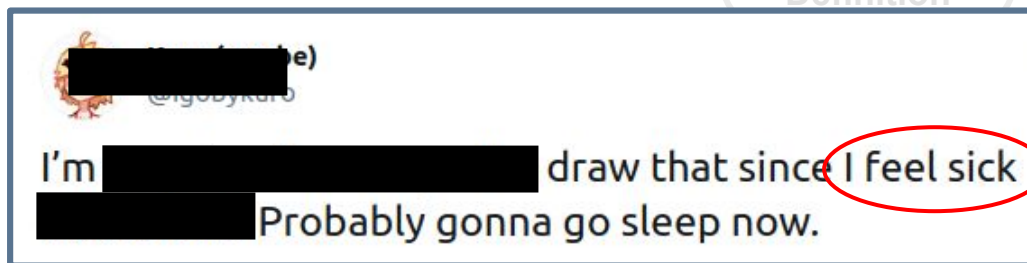Data Analysis

Construct
Definition

Selection

**Augment traces** with syntactic relations
indicating the speaker being afflicted by the flu

**Remove** non 'first-person' tweets

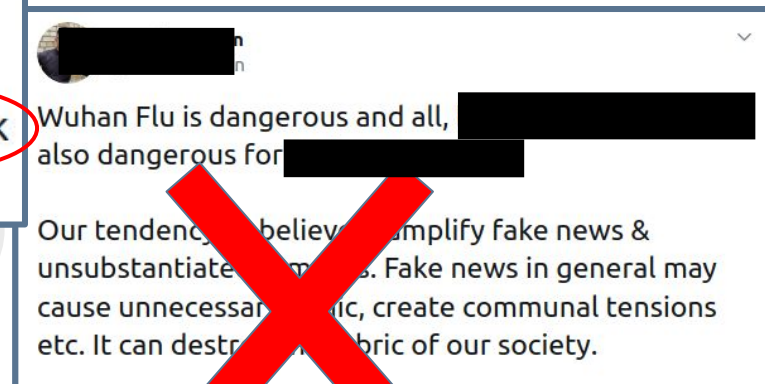= Errors due to **augmenting** or
**filtering** out tra

Data
Collection

I'm ███████████████████ draw that since I feel sick ███████ Probably gonna go sleep now.

Wuhan Flu is dangerous and all, ████████ also dangerous for ███████████

Our tendency ██ believ██ ██mplify fake news & unsubstantiate██ ██m██s. Fake news in general may cause unnecessar██ ██ic, create communal tensions etc. It can destr██ ██ bric of our society.

**error rate of
augmentation method**

**Trace
Augmentation
Error**

**Data
Preprocessing**

**Trace
Reduction
Error**

**criteria used for
filtering tweets may be
error-prone**

Data
Analysis

**TSE**
- Preprocessing:
  Processing Error

28

**Augment** users with demographic information

Construct ion

Data Collection

Data Preprocessing

Data Analysis

29

Construct
Definition

**Remove** bots or organizations
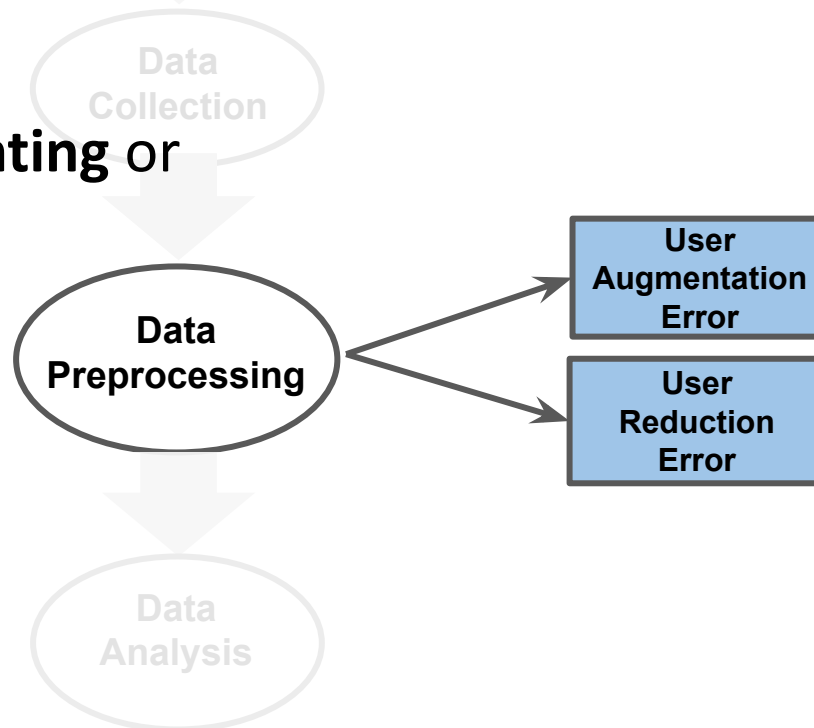
Data
Collection

**Data
Preprocessing**

Data
Analysis

30

**Augment users** with demographic information
Remove bots or organizations

= Errors due to **augmenting** or
**filtering** out users



Data
Collection

Data
Preprocessing

User
Augmentation
Error

User
Reduction
Error

Data
Analysis

31

## Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*

**Joy Buolamwini**                     JOYAB@MIT.EDU
*MIT Media Lab 75 Amherst St. Cambridge, MA 02139*

**Timnit Gebru**              TIMNIT.GEBRU@MICROSOFT.COM
*Microsoft Research 641 Avenue of the Americas, New York, NY 10011*

**Augment users** with demographic information
Remove bots or organizations

Construct

Data Collection

= Errors due to **augmenting** or **filtering** out users

**error rate of augmentation method**

**User Augmentation Error**

**TSE**
~  (Preprocessing: Processing Error)

**Data Preprocessing**

**User Reduction Error**

**criteria used for filtering users may be error-prone**

Data Analysis

32

**MEASUREMENT**

**REPRESENTATION**

Construct Definition

Validity

searching for flu related information or tweeting about flu may not sufficiently cover all incidences of flu

Platform Affordances Error

Platform(s) have affordances which distort traces

Platform Selection

Platform Coverage Error

Platform(s) are not representative of target US population!

Trace Selection Error

Some tweets chosen may not be relevant to influenza like illness

error rate of augmentation method

Trace Augmentation

Trace Reduction Error

criteria used for filtering tweets may be error-prone

Data Collection

User Selection Error

Excluding Spanish tweets excludes Hispanic population of the US. Alternately, includes English speakers from other countries

error rate of augmentation method

Data Preprocessing

User Augmentation Error

User Reduction Error

criteria used for filtering users may be error-prone

Data Analysis

33

Construct
Definition

Platform

Develop machine learning model to
**measure rates of flu** based on tweeting
activity



Data
Analysis

34

**MEASUREMENT**

Construct
Definition

Platform
???

Develop machine learning model to
**measure rates of flu** based on tweeting
activity

= Errors due to the choice
of modeling or aggregation
used by the researcher as well
as how the traces are mapped
to the users



Prepr???

**Trace
Measurement
Error**

**Data
Analysis**

Construct
Definition

Platform

Develop machine learning model to
**measure rates of flu** based on tweeting
activity

= Errors due to the choice
of modeling or aggregation
used by the researcher as well
as how the traces are mapped
to the users

I feel sick 🤢

I still have the flu. Someone @ me ███████ entertainment 😒😌💜

Alright, I probably have a fever. I'm going under ████████████ and
hope to ████████ to this tweet filled with wholesome fan art of good
waifus.

Prepr

<span style="color:red">**error due to how
traces are included in
the model**</span>

**Trace
Measurement
Error**

Data
Analysis

**TSE**
~ Analysis
~ Questionnaire design:
validity

36

Construct
Definition

Platform
Selection

Use demographic data (say, location) to
correct for coverage errors using post-
stratification

Data
Collection

Data
Preprocessing

Data
Analysis



Explore flu trends across the U.S.

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity in your state up to two weeks faster than traditional systems. Read more »

United States flu activity: Low          Entire United States

2008-2009   Past years          Minimal  Low  Moderate  High  Intense
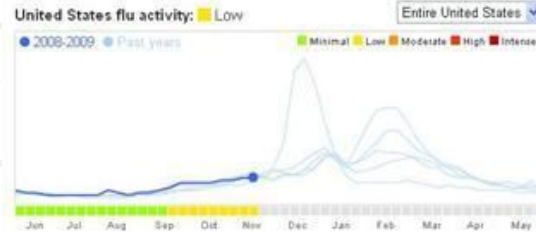
Construct
Definition

Platform
Selection

Use demographic data (say, location) to **correct for coverage errors** using post-stratification

Data
Collection

### Explore flu trends across the U.S.

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity in your state up to two weeks faster than traditional systems. Read more »

United States flu activity: Low

Entire United States ▾

● 2008-2009  ● Past years

Minimal  Low  Moderate  High  Intense

Jun  Jul  Aug  Sep  Oct  Nov  Dec  Jan  Feb  Mar  Apr  May

= Errors due to correcting for representation errors through reweighting

Data
Analysis → **Adjustment Error**

Construct
Definition

Platform
Selection

Use demographic data (say, location) to correct for coverage errors using post-stratification

**TSE**
- Adjustment error

Data
Collection

= Errors due to correcting for representation errors through reweighting



Explore flu trends across the U.S.

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity in your state up to two weeks faster than traditional systems. Read more »

United States flu activity: Low          Entire United States

2008-2009   Past years          Minimal  Low  Moderate  High  Intense

**error due to choice of reweighting method or scaling used**

Data
Analysis → Adjustment Error

**MEASUREMENT**

**REPRESENTATION**

Construct Definition

Validity

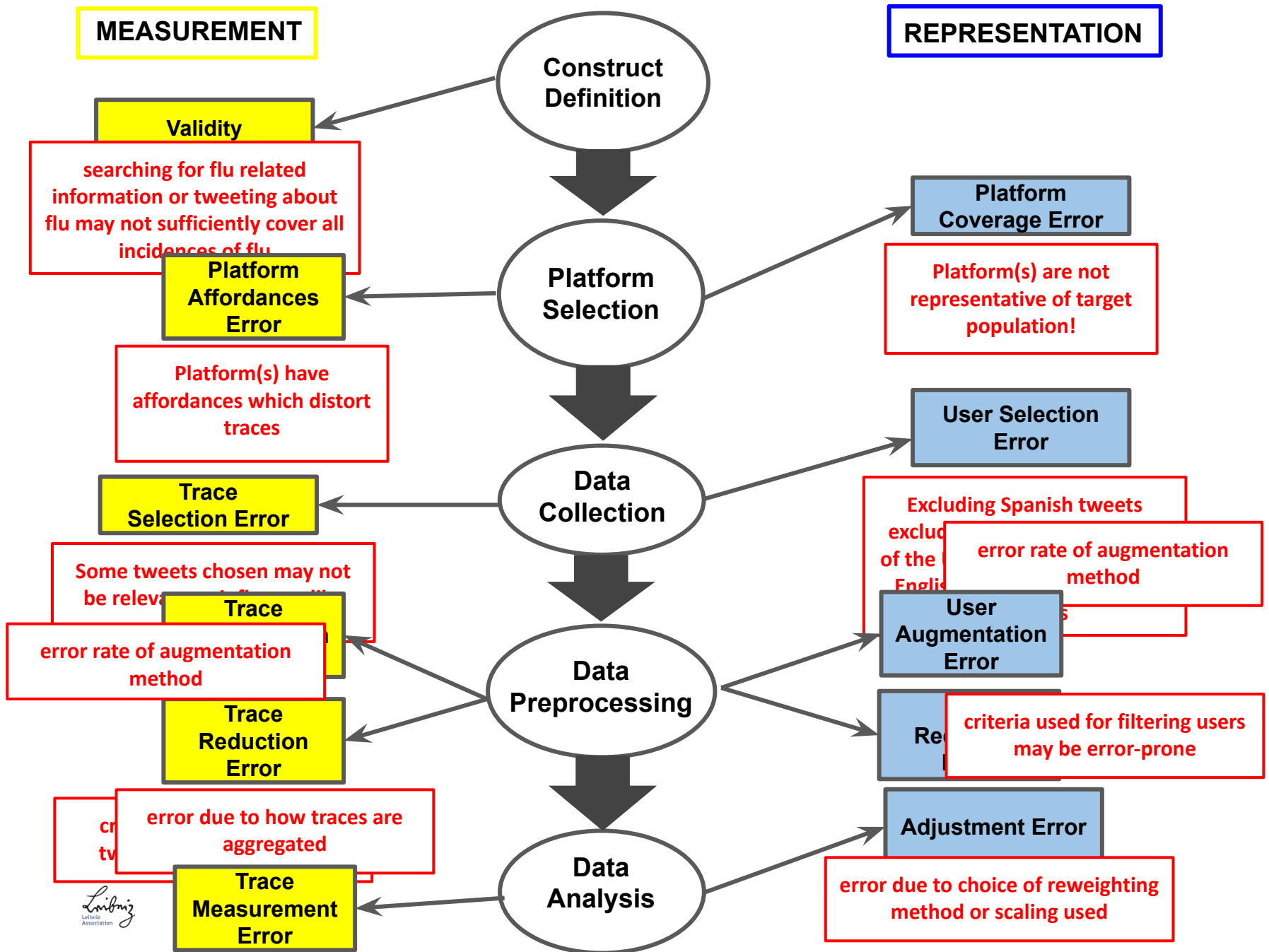searching for flu related information or tweeting about flu may not sufficiently cover all incidences of flu

Platform Affordances Error

Platform(s) have affordances which distort traces

Platform Selection

Platform Coverage Error

Platform(s) are not representative of target population!

Data Collection

User Selection Error

Trace Selection Error

Some tweets chosen may not be relevant or influenced by...

error rate of augmentation method

Trace

Excluding Spanish tweets exclud... of the ... Englis...

error rate of augmentation method

User Augmentation Error

Trace Reduction Error

Data Preprocessing

Re...

criteria used for filtering users may be error-prone

error due to how traces are aggregated

cr... tw...

Trace Measurement Error

Data Analysis

Adjustment Error

error due to choice of reweighting method or scaling used

40

# Time as factor

- Every error has a time component, e.g.:
  - System drift → Platform affordances
  - Population drift → Coverage
  - Behavioral drift → Trace selection
  - Behavioral drift, System drift → User augmentation
  - …

# TED-On: A Total Error Framework for Digital Traces of Humans



Icons used in this image have been designed by Becris, EliasBikbulatov and Pixel perfect from www.flaticon.com