

# Assessing and Improving the Accuracy of Estimators from Blended Data

Paul P. Biemer<sup>1,2</sup>,  
Ashley Amaya<sup>1</sup>

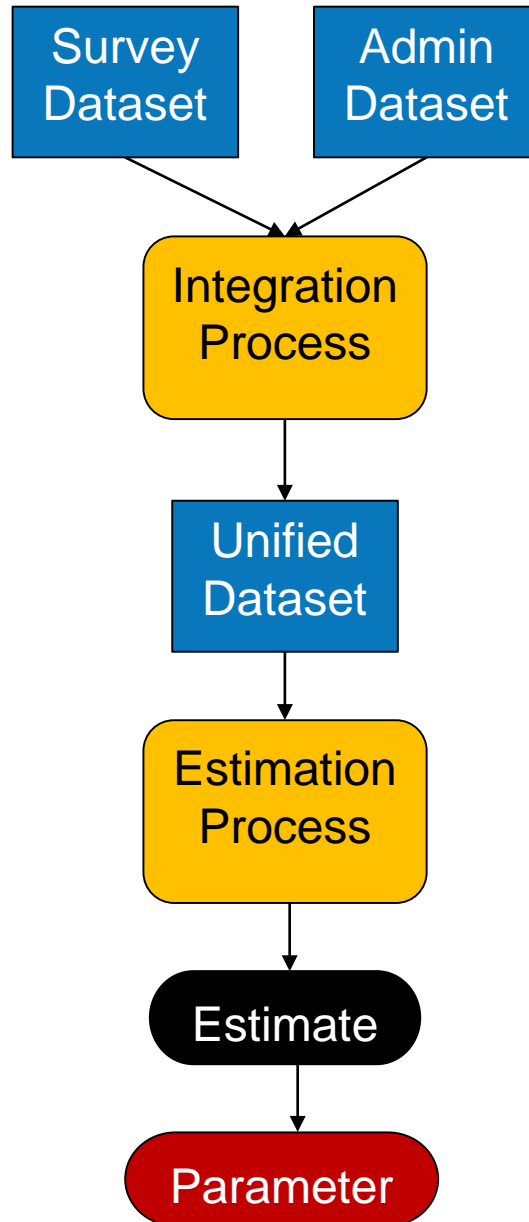
<sup>1</sup> RTI International; <sup>2</sup>University of  
North Carolina



# Outline

- Hybrid estimators
- An total error framework for datasets
- An total error framework for hybrid estimators
- Types of error risks
- Error risk profiles
- Illustration of the concepts

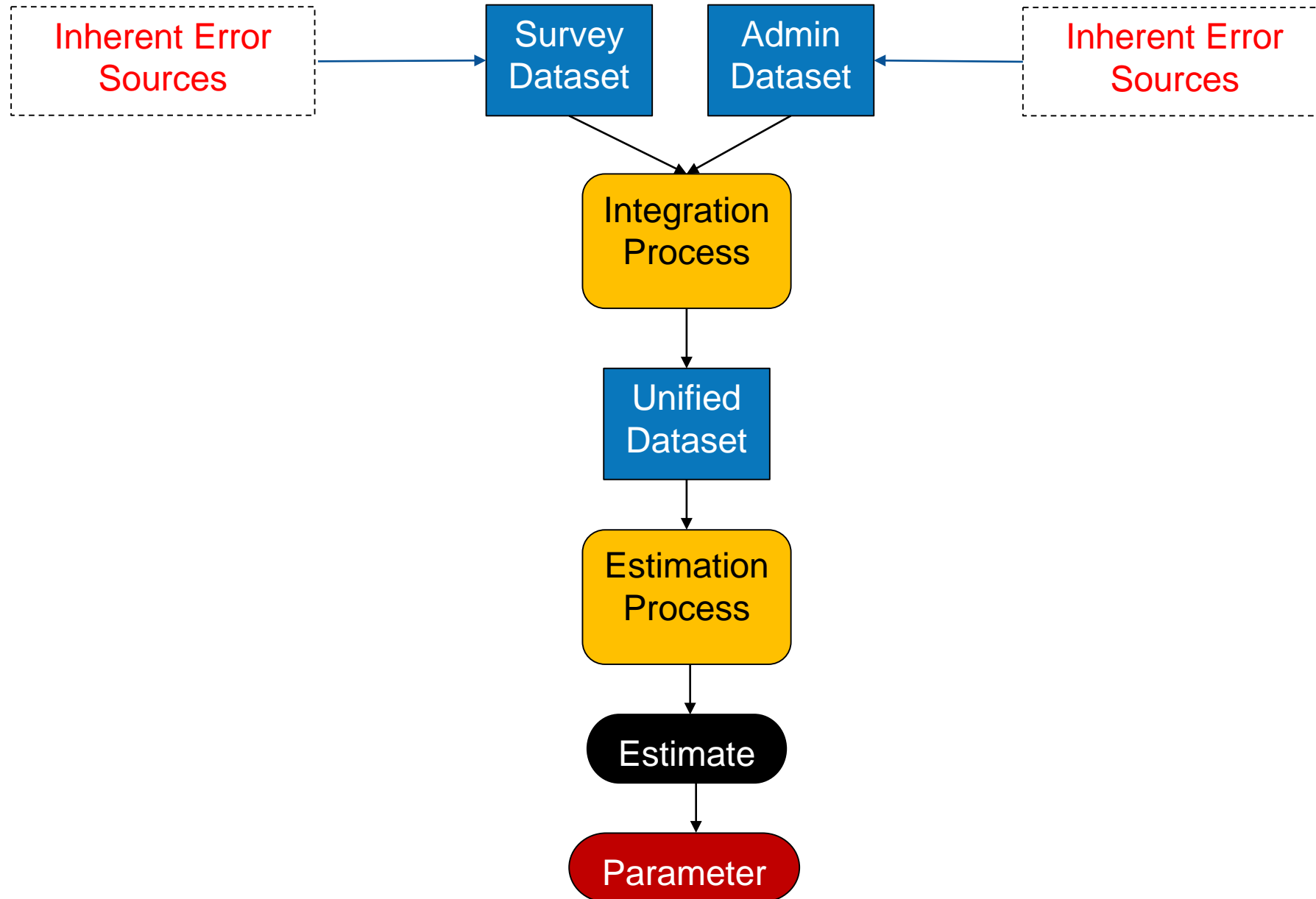
# The Hybrid Estimation Process



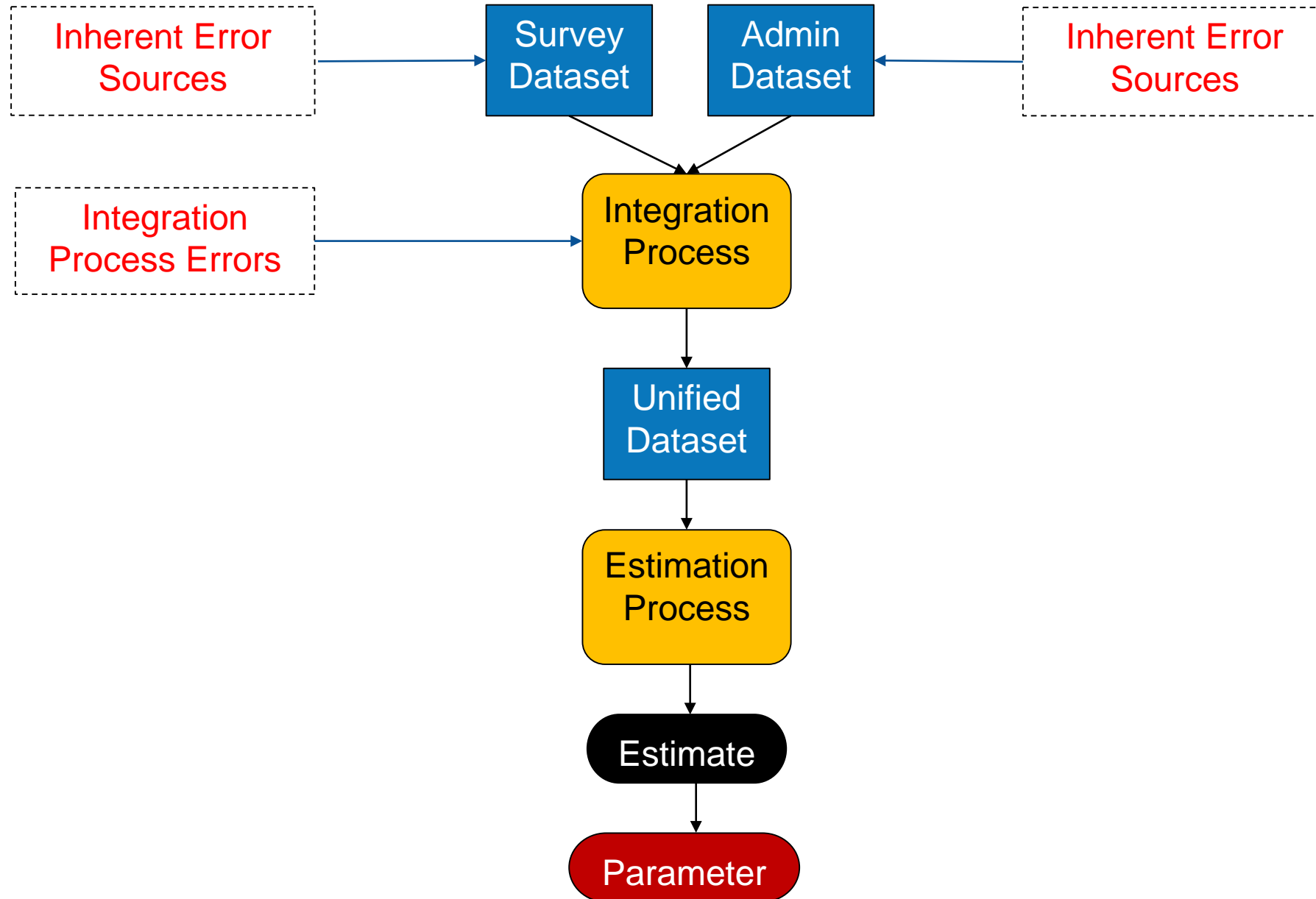
# The Hybrid Estimation Process



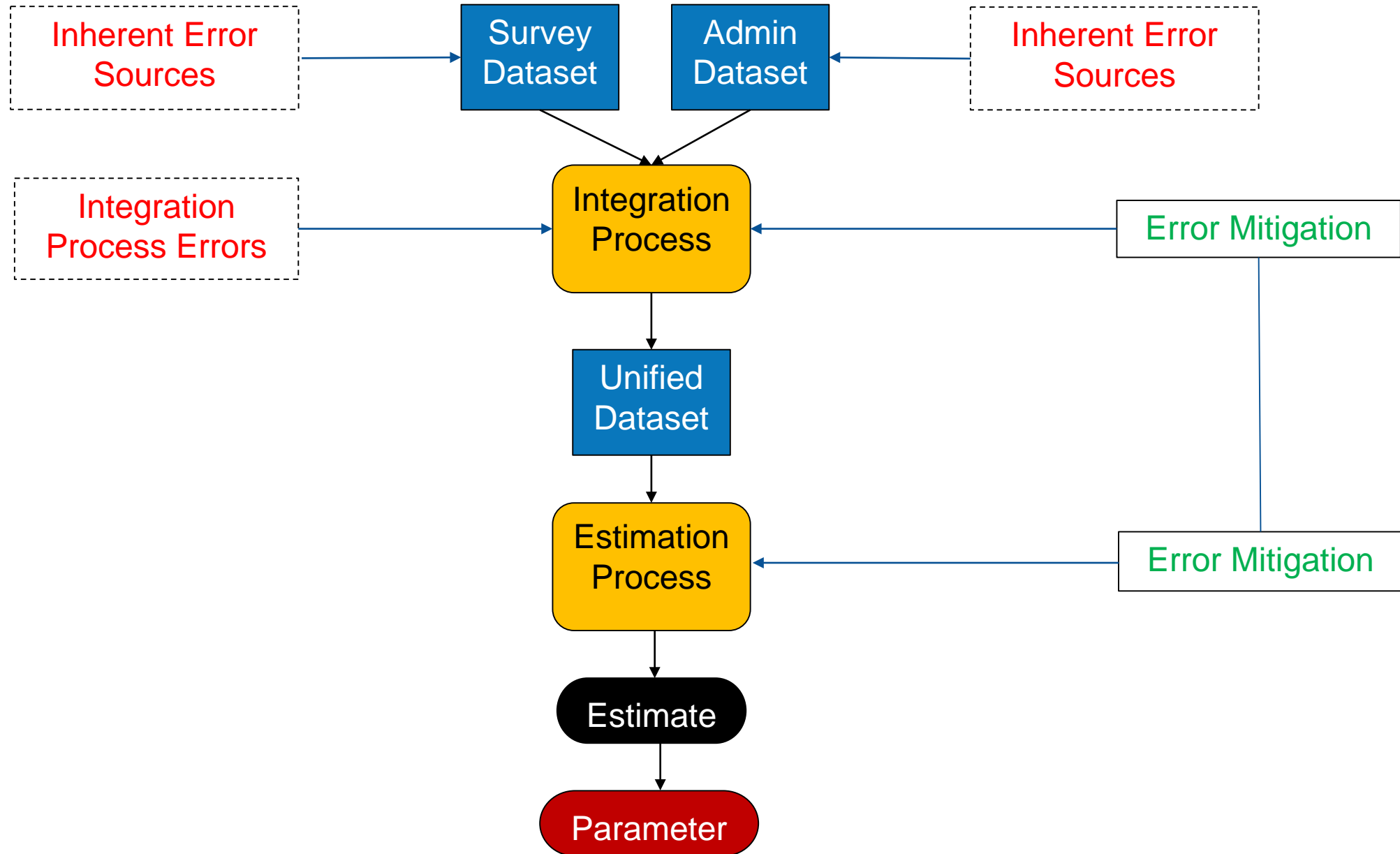
# The Hybrid Estimation Process



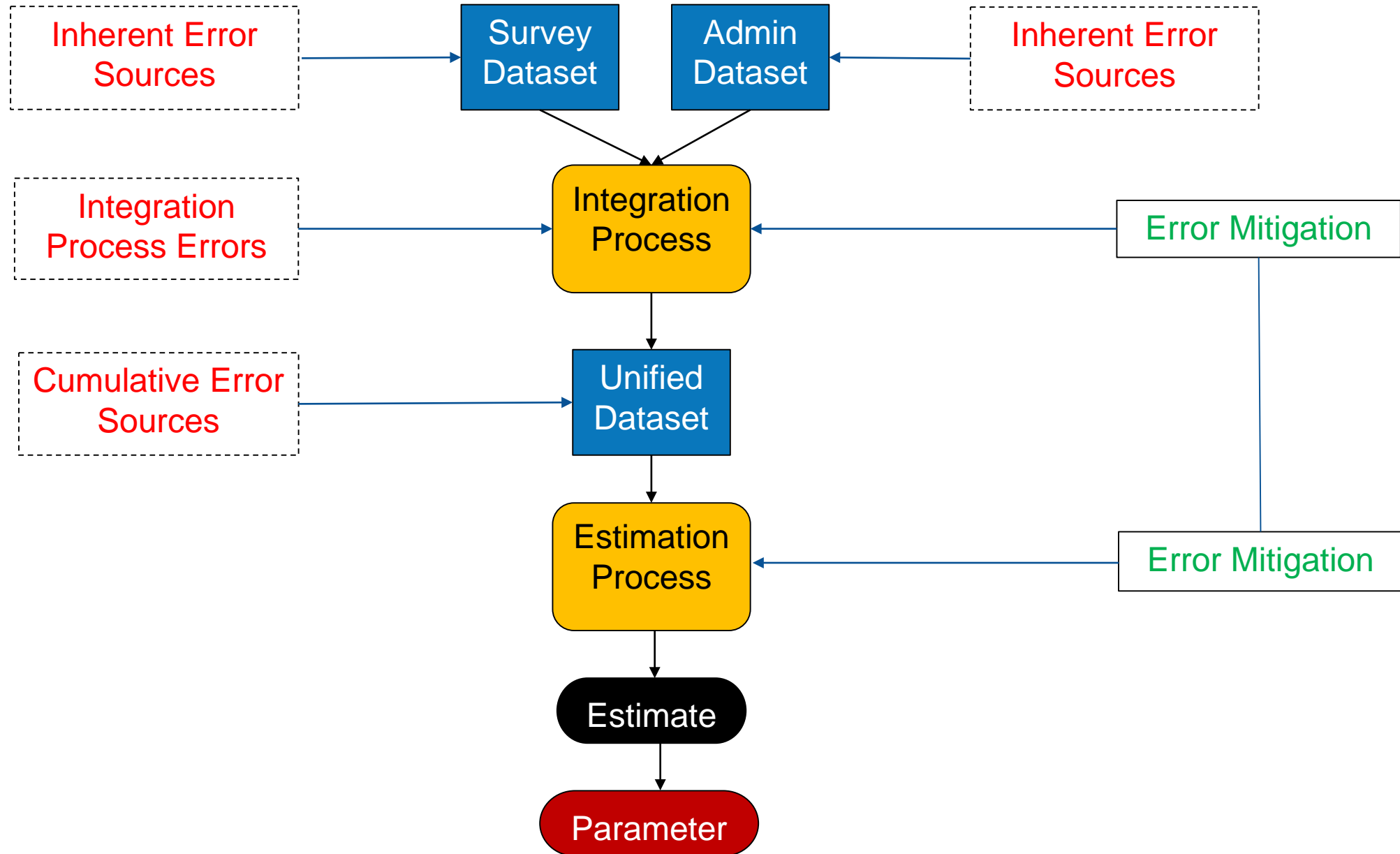
# The Hybrid Estimation Process



# The Hybrid Estimation Process

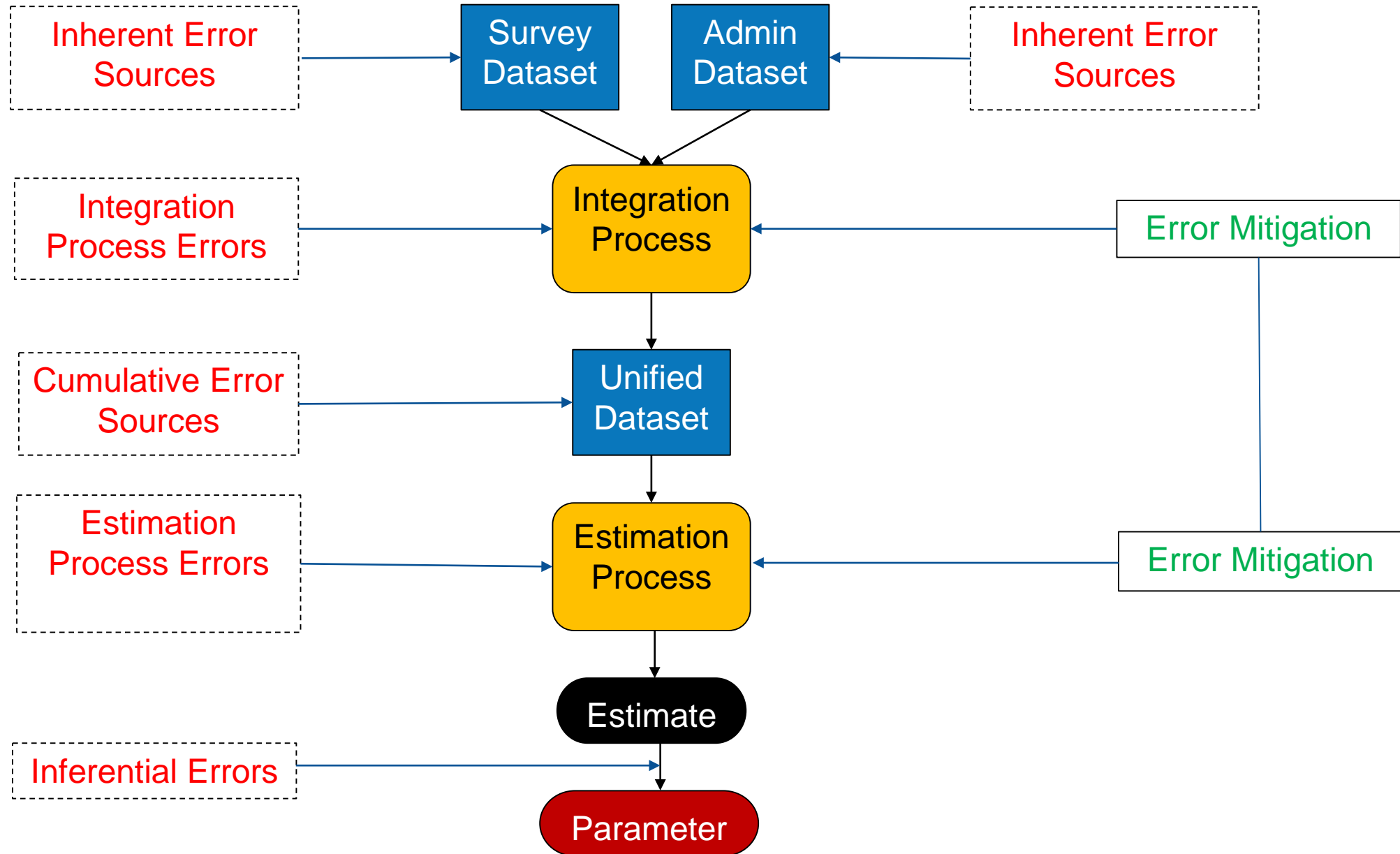


# The Hybrid Estimation Process





# The Hybrid Estimation Process

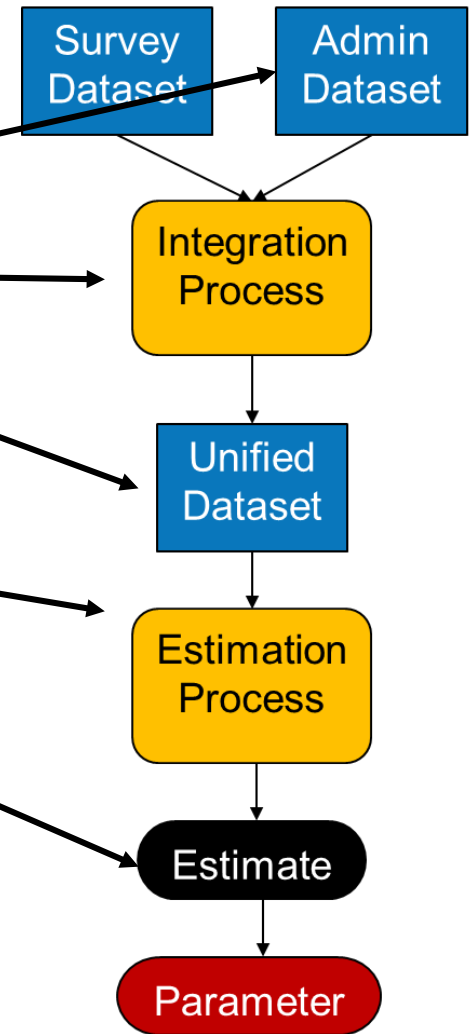


# Questions Regarding Hybrid Estimator Accuracy

- What error sources are associated with the unified dataset?
- Which of these pose the greatest *intrinsic* risks to data accuracy?
- Among the hybrid estimators that might be constructed from the unified dataset, which estimator minimizes the *total error risk*?
- What are the major *intrinsic* and *residual* error risks associated with the **hybrid** estimator?
- Which of these error risks could be further mitigated to maximally increase the accuracy of the hybrid estimator?

# A Total Error Framework Can be Specified for Each Stage of the Process

- Inherent errors associated with survey dataset
- Inherent errors associated with administrative dataset
- Integration processing errors
- Cumulative errors associated with the unified dataset
- Estimation processing errors
- Final errors associated with the estimate



Error mitigation can occur at various stages of each the process

(see, for example, 3 stage framework of Reid, et al, 2017)

# In many cases it suffices to simply describe the errors in the final output

- Total error model for registers, frames and other datasets
- Total error model for survey point estimates
- Total error model for hybrid estimates
- Total error models for compilations such as the GDP and various price indexes

# A Total Error Framework for a Generic Dataset

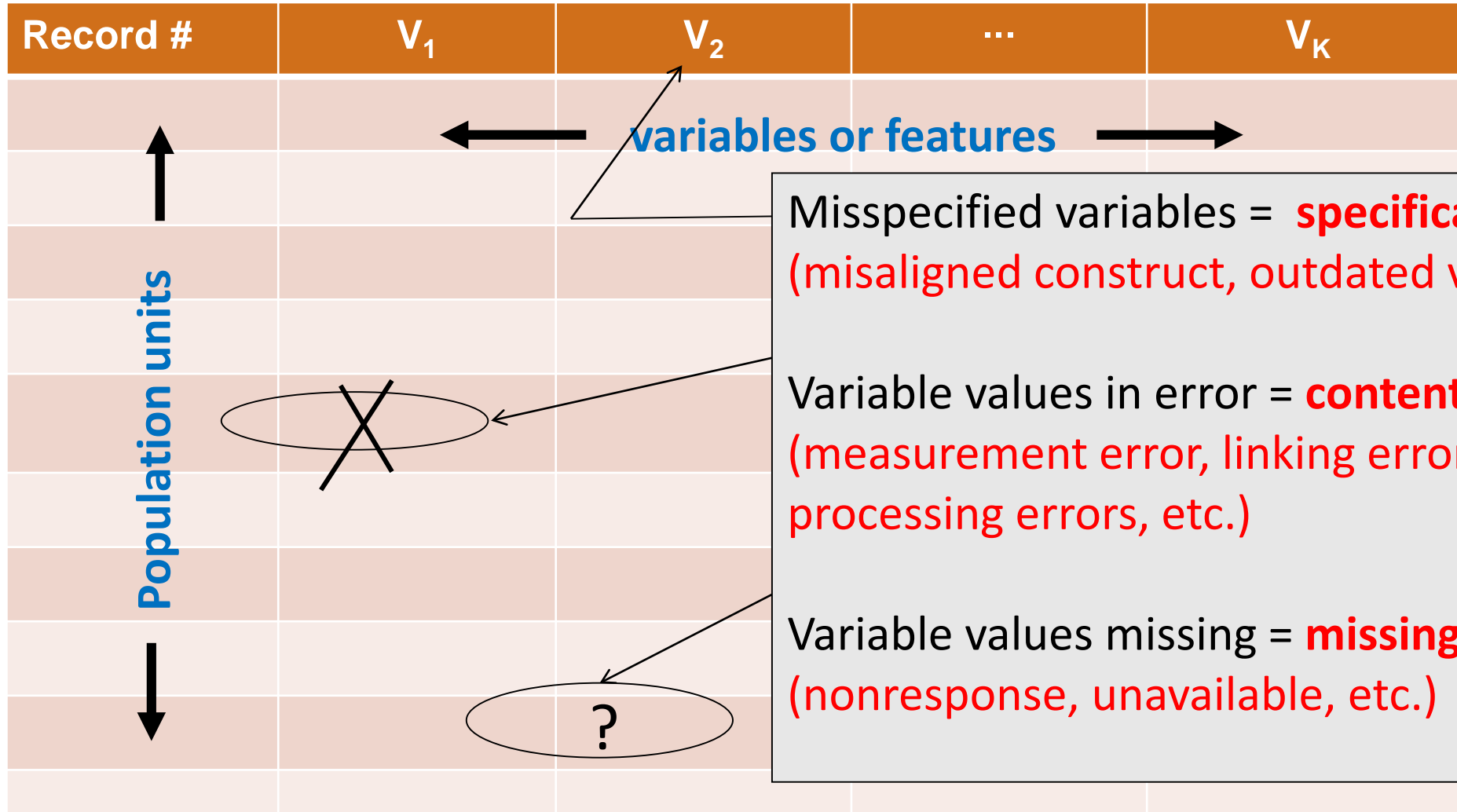
## Typical File Structure

Record #	$V_1$	$V_2$	...	$V_K$
	← variables or features →			



# Column and Cell Errors

## Typical File Structure



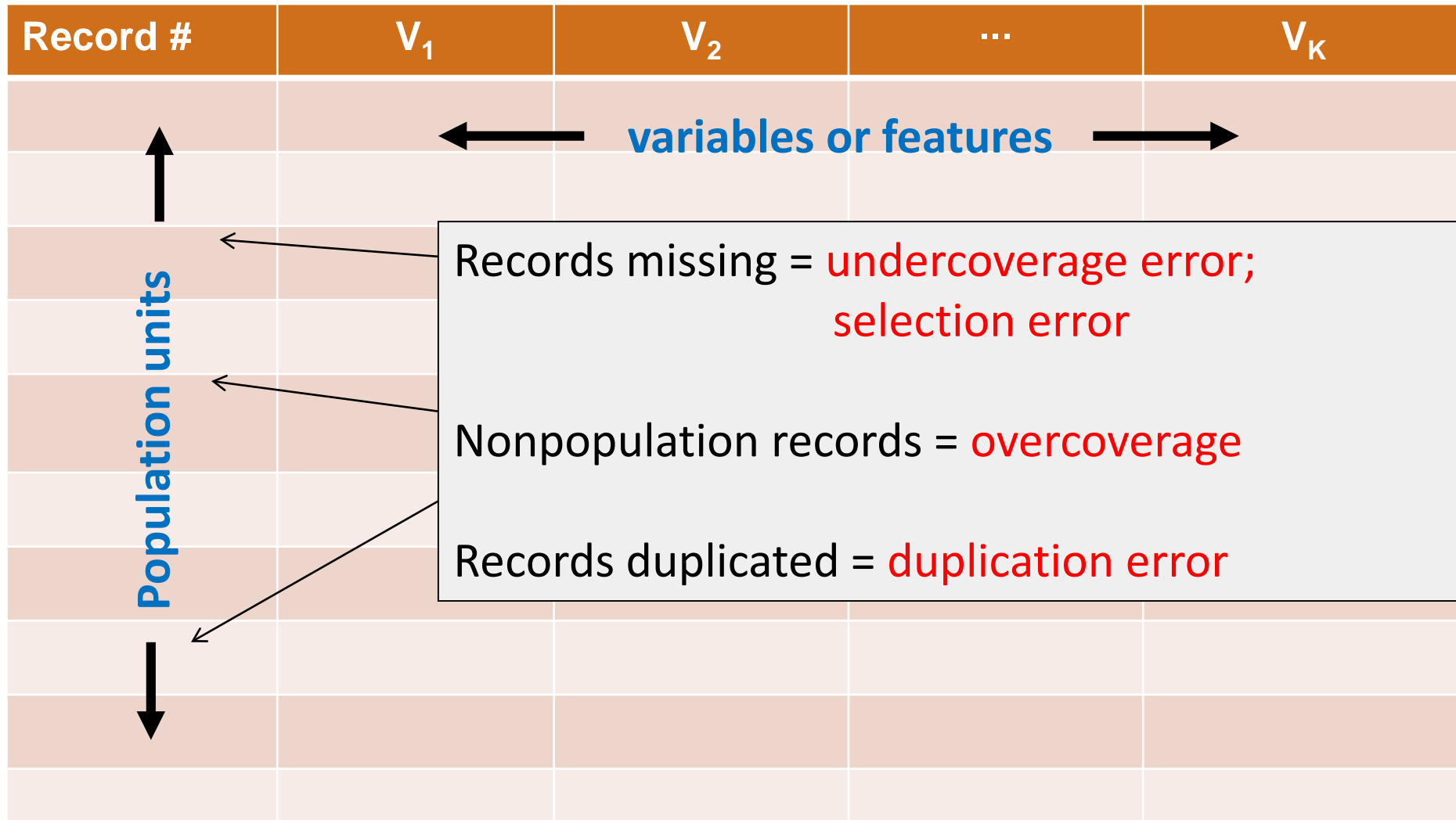
Misspecified variables = **specification error** (misaligned construct, outdated value, etc.)

Variable values in error = **content error** (measurement error, linking errors, data processing errors, etc.)

Variable values missing = **missing data** (nonresponse, unavailable, etc.)

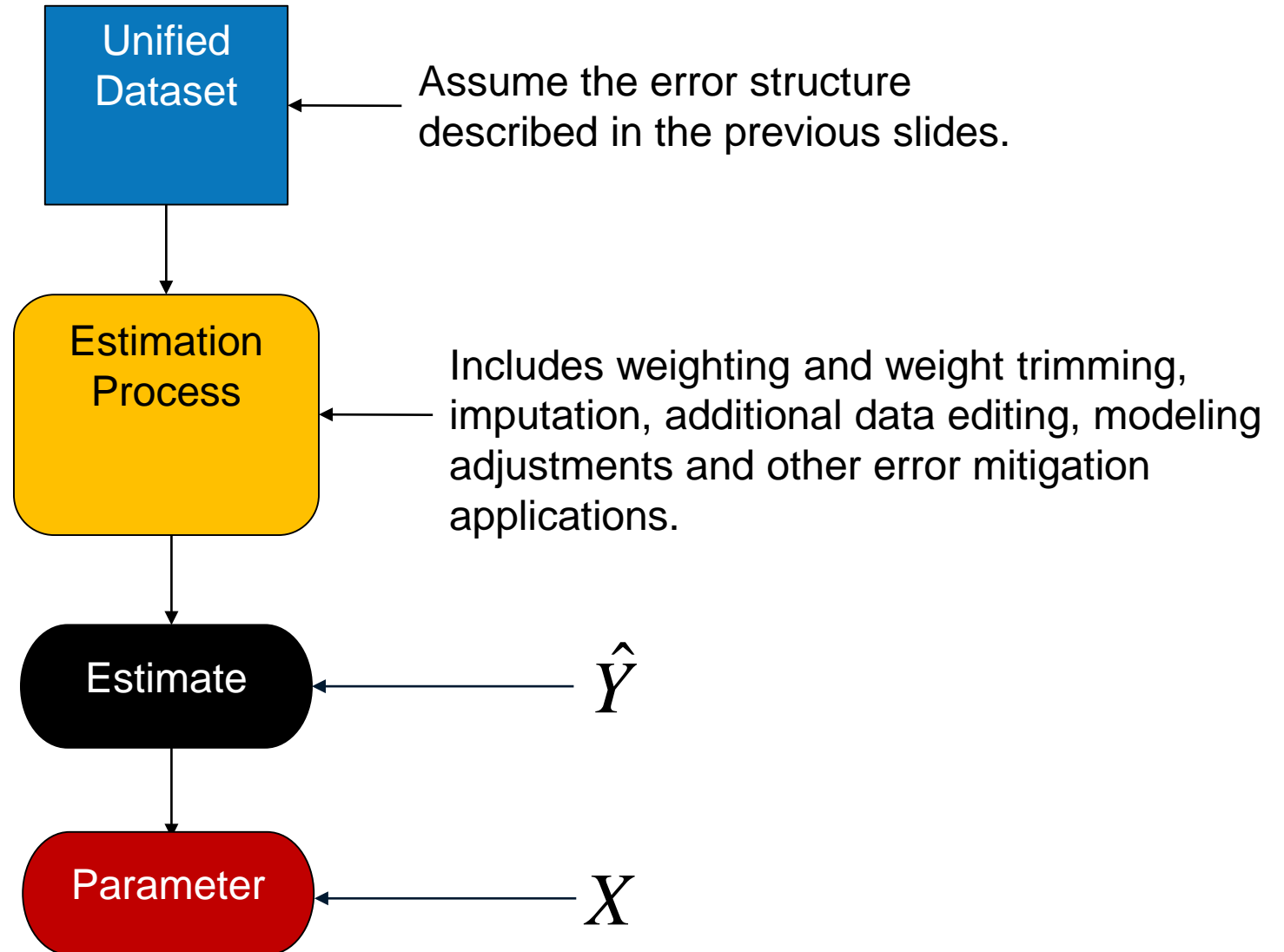
# Row errors

## Typical File Structure





# Errors Associated with the Hybrid Estimation Process



# Total Error Model for Hybrid Estimators

$$\underbrace{\hat{Y} - X}_{\text{total error}} = \underbrace{(\hat{Y} - Y)}_{(\varepsilon_1 + \mathbf{L} + \varepsilon_6)} + \underbrace{(Y - X)}_{\varepsilon_7}$$

$\varepsilon_1$  = Selection error

$\varepsilon_2$  = Coverage error (over-, under-, duplication)

$\varepsilon_3$  = Missing data error

$\varepsilon_4$  = Content error

$\varepsilon_5$  = Data processing error

$\varepsilon_6$  = Model/estimation error

$\varepsilon_7$  = Specification error

# Assessing Error Risk

## Types of Error Risks

- Intrinsic risk – risk that an error source poses if no steps are taken to reduce the error; error risk of “doing nothing.”
  - Example: The intrinsic risk of nonresponse bias in a linear estimator is

$$B_I = \frac{\text{cov}(y_i, \rho_i)}{\bar{\rho}}$$

- Residual risk – risk of error for a source that remains after mitigation strategies have been applied.
  - Example: After nonresponse weighting adjustments have been applied, the residual risk of bias is

$$B_R \leq B_I$$

# Risk Profile Comparing Survey, Administrative and Unified Datasets: Either Intrinsic or Residual Risks

<b><i>Error Sources</i></b>	<b><i>Survey Dataset</i></b>	<b><i>Administrative Dataset</i></b>	<b><i>Unified Dataset</i></b>
<b><i>Specification</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Coverage: Undercoverage</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Coverage: Overcoverage</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Coverage: Duplication</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Selection</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Content</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<b><i>Missing Data</i></b>	Risk level (1, 2, 3)	Risk level (1, 2, 3)	Risk level (1, 2, 3)

# Intrinsic Error Risk Profile Comparing Survey and Hybrid Estimates

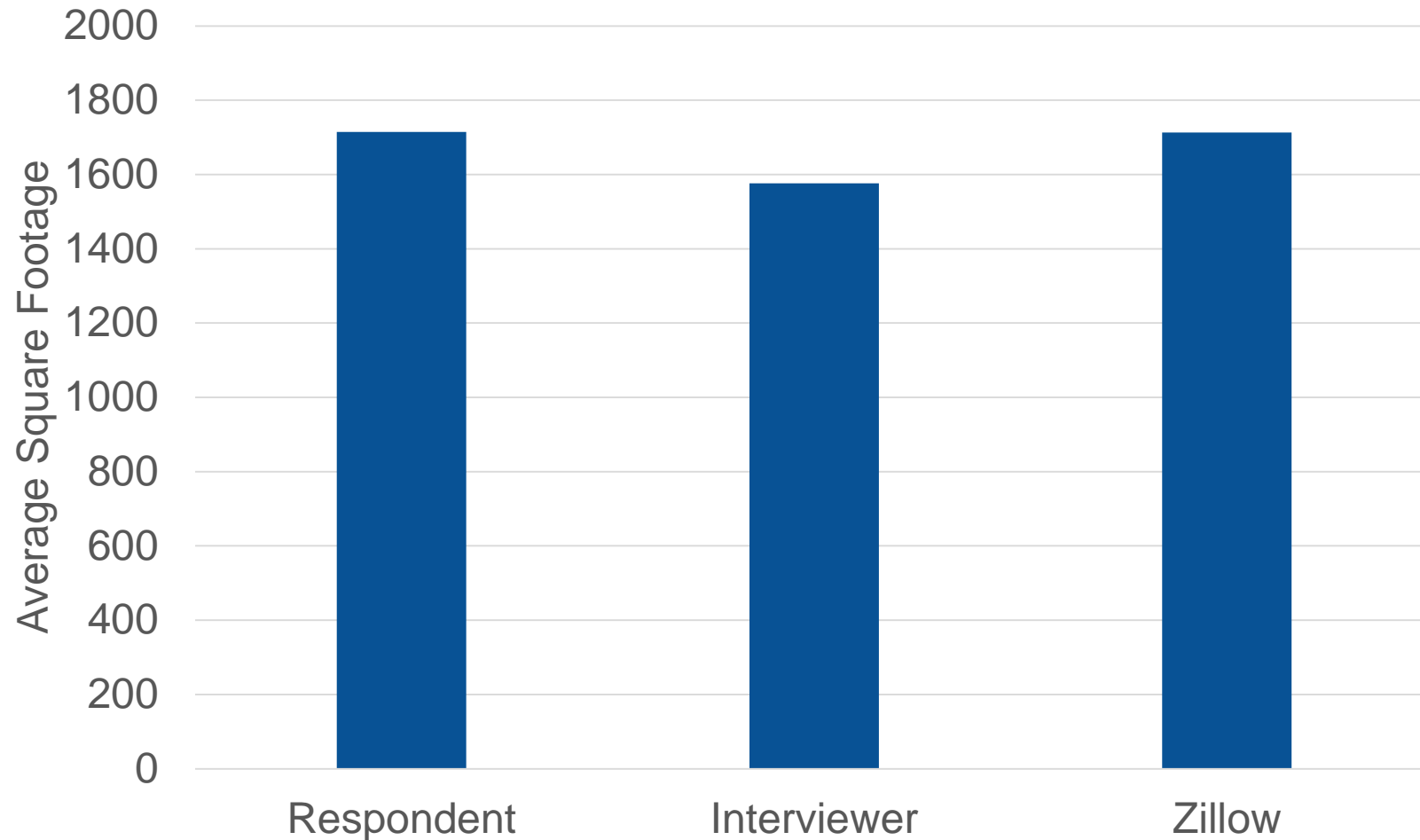
<i><b>Error Sources</b></i>	<i><b>Survey Estimator</b></i>	<i><b>Hybrid Estimator</b></i>
<i><b>Specification</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Coverage: Undercoverage</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Coverage: Overcoverage</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Coverage: Duplication</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Sampling/Selection</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Measurement/Content</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Data Processing</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Nonresponse/Missing data</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)
<i><b>Modeling/estimation</b></i>	Risk level (1, 2, 3)	Risk level (1, 2, 3)

# Case Study: Error Mitigation for Energy Use Survey Square Footage Data using Unified Data

- Data sources
  - Survey data: 2015 Residential Energy Consumption Survey (RECS)
    - $n \approx 2,400$  completed cases
  - Big Data (data pulled from various sources)
    - Zillow
    - Acxiom
    - CoreLogic
- Variable of interest: housing unit square footage
- **Goal: Integrate the external data sources with the survey data to improve and/or evaluate the accuracy of survey square footage data**

# Evidence of Nonsampling Error from the RECS

## RECS Average Reported Square Footage



# More Evidence of Intrinsic Error Risks

	<b>Survey (R)</b>	<b>Zillow</b>
<b>NR/Missing</b>		
Unit NR/missing rate	58.2%	22.0%
Item NR rate	19.2%	
<b>Overcoverage rate</b>	11.8%	~0
<b>Undercoverage rate</b>	~0	15.3%
<b>Reliability</b>	50%	68%



# Intrinsic Error Risk Profile for the RECS, Zillow and Unified Datasets

<i>Error Sources</i>	<i>RECS</i>	<i>Zillow</i>	<i>RECS U Zillow</i>
<i>Specification</i>	2	2	2
<i>Coverage: Undercoverage</i>	1	2	1
<i>Coverage: Overcoverage</i>	2	1	1
<i>Coverage: Duplication</i>	2	1	1
<i>Selection</i>	3	1	3
<i>Content</i>	3	3	3
<i>Missing Data</i>	3	2	1
<i>Average</i>	2.3	1.7	1.7

# Intrinsic Error Risk Profile for the RECS, Zillow and Unified Datasets

<i>Error Sources</i>	<i>RECS</i>	<i>Zillow</i>	<i>RECS U Zillow</i>
<i>Specification</i>	2	2	2
<i>Coverage: Undercoverage</i>	1	2	1
<i>Coverage: Overcoverage</i>	2	1	1
<i>Coverage: Duplication</i>	2	1	1
<i>Selection</i>	3	1	3
<i>Content</i>	3	3	3
<i>Missing Data</i>	3	2	1
<i>Average</i>	2.3	1.7	1.7

Unified data offers no advantage to Zillow only dataset.

# Intrinsic Error Risk Profile RECS and RECS/Zillow Hybrid Estimates

<i><b>Error Sources</b></i>	<i><b>RECS Estimator</b></i>	<i><b>RECS/Zillow Hybrid Estimator</b></i>
<i><b>Specification</b></i>	1	2
<i><b>Coverage: Undercoverage</b></i>	2	1
<i><b>Coverage: Overcoverage</b></i>	1	1
<i><b>Coverage: Duplication</b></i>	2	1
<i><b>Sampling/Selection</b></i>	3	1
<i><b>Measurement/Content</b></i>	3	2
<i><b>Data Processing</b></i>	2	2
<i><b>Nonresponse/Missing data</b></i>	3	1
<i><b>Modeling/estimation</b></i>	3	3
<i><b>Average</b></i>	<b>2.2</b>	<b>1.6</b>

Initial evaluations suggest a total error reduction with the hybrid estimator even before error mitigation efforts have been fully exploited.

## Illustration 2 – Market Research Client

- Currently conducting a large scale survey to evaluate market share for its customers products about 300 markets
- Quarterly estimates tend to be unstable in some markets
- Various administrative sets have been identified that would improve estimator stability, but each brings with it other error risks that have been fully investigated

### **Question:**

Can a hybrid estimator be constructed having greater stability than the current survey estimator without increasing total error?

# Intrinsic Risk for the Hybrid Estimator Compared to Estimators Based on the Survey and Administrative Data

Quality Component	Survey	Admin Data	Hybrid Estimator
<b>Specification</b>	<b>3</b>	<b>2</b>	<b>3</b>
<b>Coverage</b>	<b>2.7</b>	<b>2.0</b>	<b>2.7</b>
Undercoverage	3	1	3
Duplication	3	3	3
Within unit	2	2	2
<b>Selection</b>	<b>2.5</b>	<b>3</b>	<b>1.5</b>
Sample size	2	N/A	1
Weight variation	3	3	2
<b>Nonresponse/Missing Data</b>	<b>1.5</b>	<b>2</b>	<b>2</b>
Unit	1	3	2
Item	2	1	2
<b>Measurement</b>	<b>2</b>	<b>3</b>	<b>3</b>
<b>Data processing</b>	<b>2</b>	<b>1</b>	<b>2</b>
Keying/editing	1	1	1
Design weighting	3	N/A	3
<b>Estimation/modeling</b>	<b>1</b>	<b>1</b>	<b>1</b>
<b>Analysis</b>	<b>1</b>	<b>3</b>	<b>1</b>
<b>Overall assessment</b>	<b>2.0</b>	<b>2.1</b>	<b>2.0</b>

# Summary

- A total error framework decomposes total error so that key subcomponents can be identified and addressed.
- A unified error risk framework facilitates comparisons across individual and unified data sources.
- An error risk profile can provide insights regarding the quality implications of unified datasets
  - Assesses intrinsic risks by error source
  - Helps determine whether residual risk can be reduced by data unification.

**Thank you!**

Please direct inquiries to:

Paul Biemer: [ppb@rti.org](mailto:ppb@rti.org)

Ashley Amaya: [aamaya@rti.org](mailto:aamaya@rti.org)