

Dirty Data in the Newsroom

Comparing Data Preparation in Journalism and Data Science

ACM CHI Conference on Human Factors in Computing Systems
April 23-28, 2023, Hamburg, Germany



Stephen Kasica
University of British Columbia
Vancouver, Canada



Charles Berret
Linköping University
Norrköping, Sweden

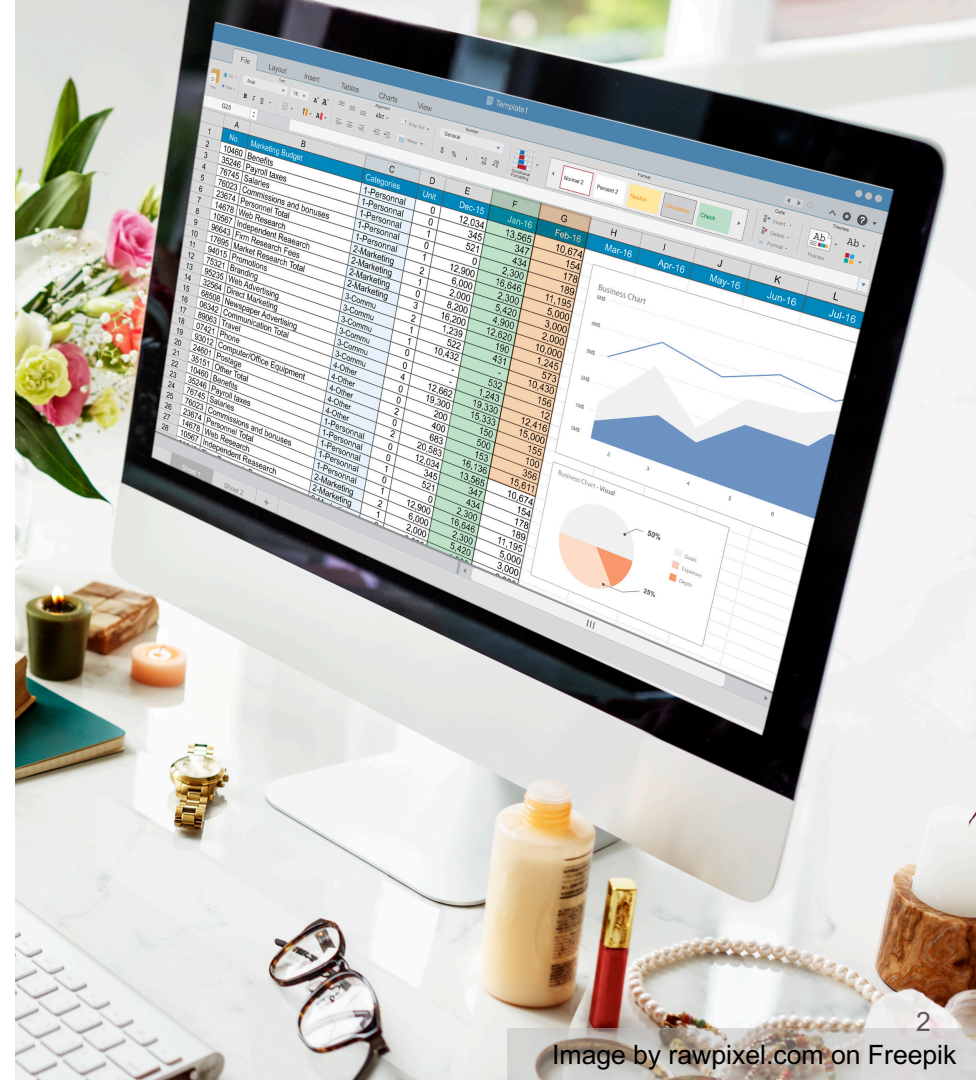


Tamara Munzner
University of British Columbia
Vancouver, Canada
@tamara@vis.social 
[@tamaramunzner](https://twitter.com/tamaramunzner) 



Data Preparation

- Getting data ready for analysis or visualization
 - Includes: wrangling, cleaning, munging, gathering, integrating, etc.
- Time-consuming process in data science
 - Up to 80% of someone's time

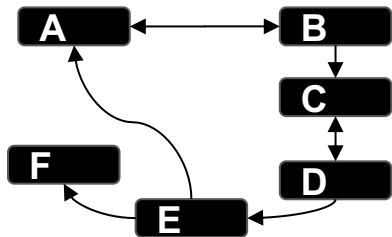






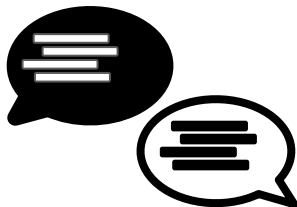
**How closely does research
on data scientists apply to data journalists,
with regards to data preparation?**

Phase 1



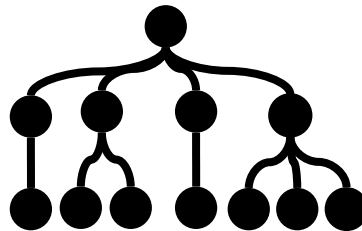
Data science
process papers (16)

Phase 2



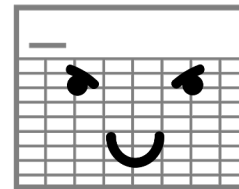
Data journalist
interviews (36)

Phase 3



Dirty data
Taxonomies (16)

Phase 4



Data preparation
nightmares (63)

Issues

11

15

Activities

30


13

Augmented model of
preparation activities

Data issues (60)


Model-discrepancy
taxonomy of
dirty data

Four challenges in
multi-table
data integration

 From data science

 From data journalism

 From database research

 Our contributions⁶

Contributions



1

Augmented model
of preparation
activities



2

Model-discrepancy
taxonomy of
dirty data



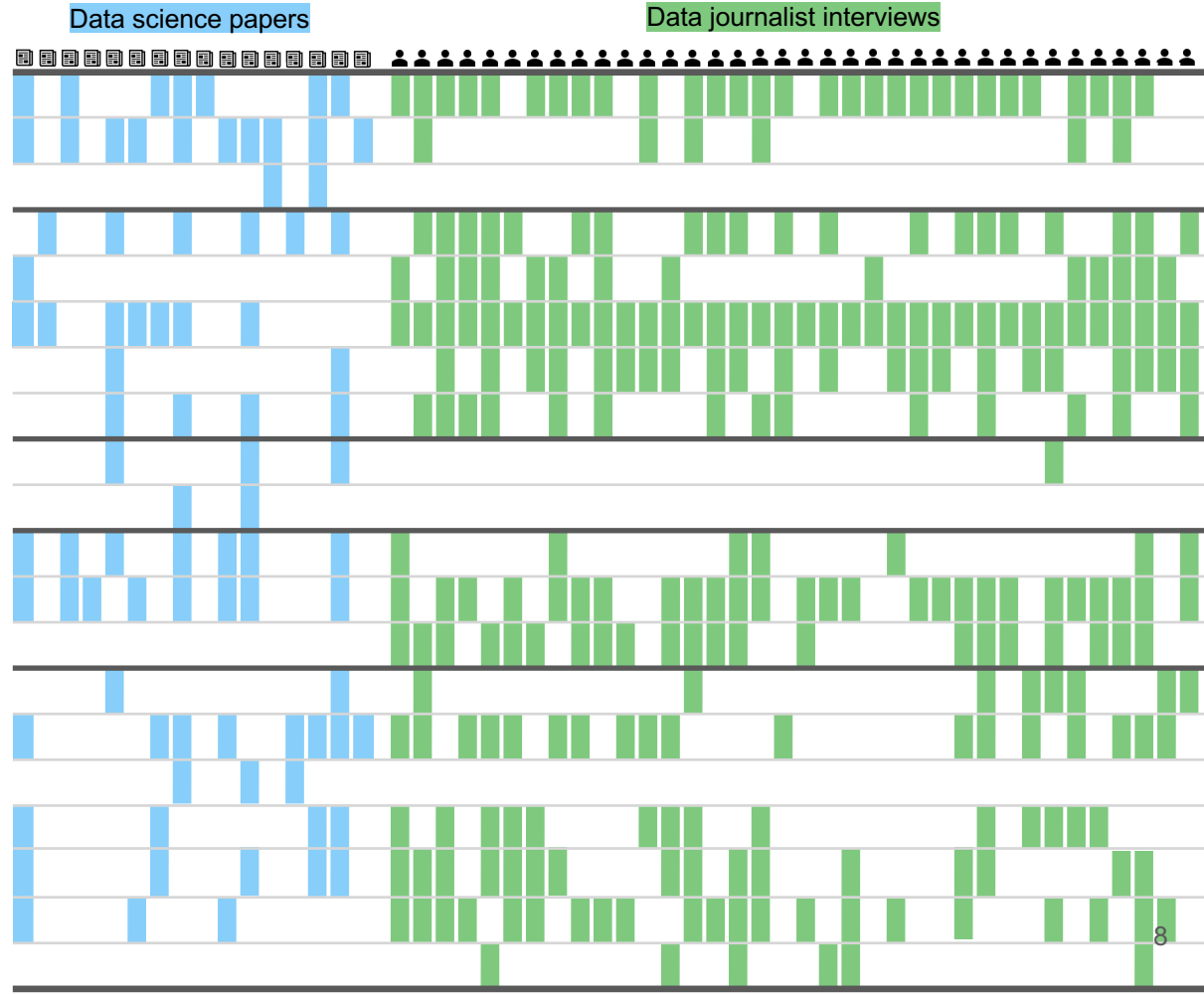
3

Challenges in
multi-table
data integration

Augmented model of prep. activities

Crisan Model	Our analysis
--------------	--------------

Prepare	Initiate	Establish goals
		Make a plan
		Test proof of concept
	Gather	Locate existing datasets
		Collect new data
		Integrate multiple datasets
		Parse documents
		Request datasets
	Create	Impute
		Synthesize
	Profile	Assess quality
		Understand semantics
		Verify transformation
	Wrangle	Aggregate data
		Transform data schema
		Label data
		Normalize values
		Remove data
		Standardize values
		Identify items



Contributions



1

Augmented model
of preparation
activities



2

Model-discrepancy
taxonomy of
dirty data



3

Challenges in
multi-table
data integration

Model-discrepancy taxonomy of dirty data

- Consider data as a design artifact
 - Dirty data = discrepancy in mental models
- Extend issue analysis to incorporate database literature
 - Analyze 16 taxonomies on dirty data: cluster 330 issues → 45 DB issues
- Combine into synthesis set of 60 issues
- Categorize into new model-discrepancy taxonomy
 - Data qualities axis
 - Existing qualities: completeness, accuracy
 - New qualities: form, granularity, relation, semantics
 - Data objects axis: table, attribute, item, value
- More details in the paper

Contributions



1

Augmented model
of preparation
activities



2

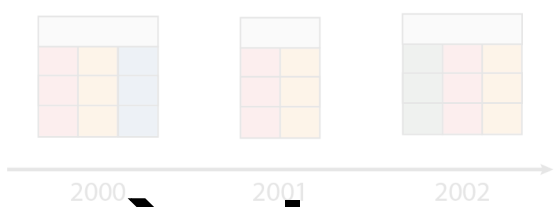
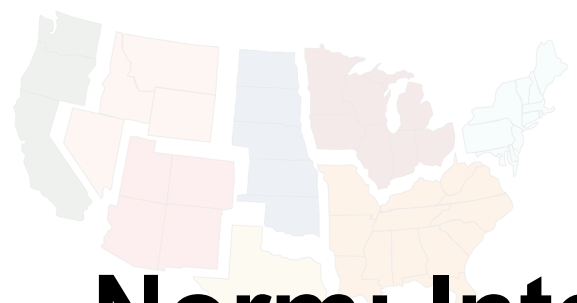
Model-discrepancy
taxonomy of
dirty data



3

Challenges in
multi-table
data integration

Four integration challenges



Norm: Integrate → clean

Regional

Diachronic

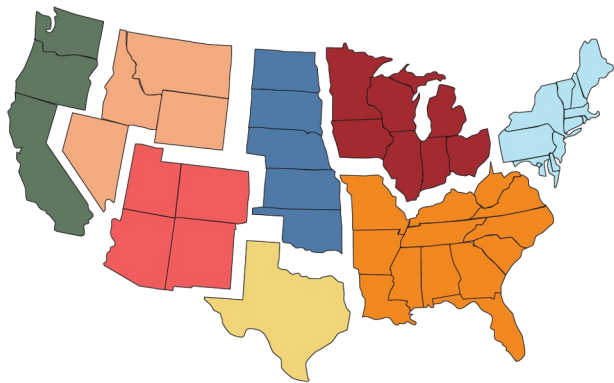
Findings: Clean → integrate



Fragmented

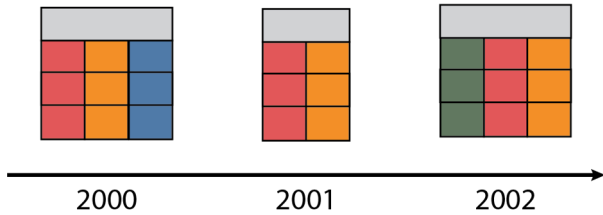


Disparate



Regional datasets

Tables with inconsistencies due to independent, spatially dispersed data sources



Diachronic datasets

Tables on the same phenomena that evolve over time

Diachronic: Economic data from Bureau of Labor Statistics

Computer analyst



Computer analyst



1970

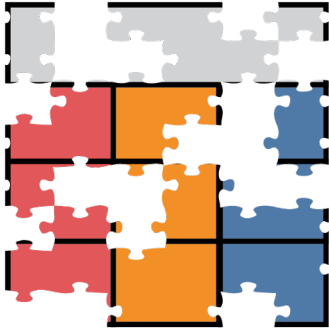
1980

1990

2000

2010

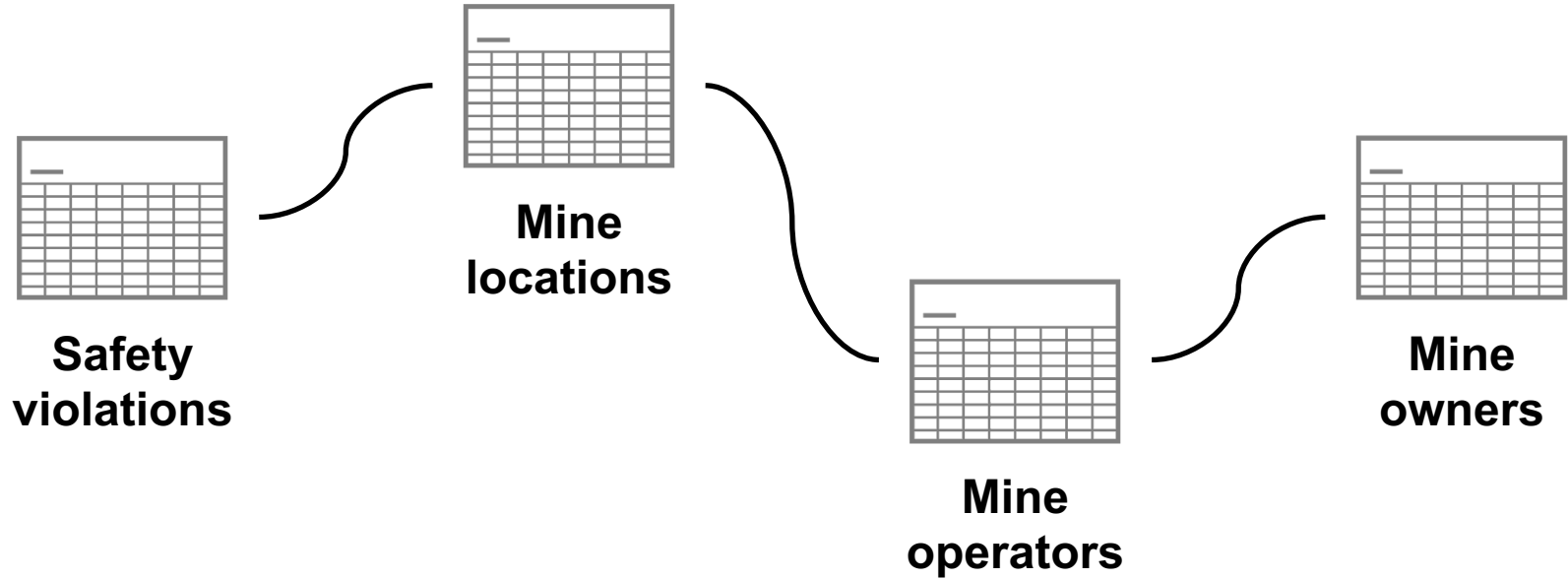
2020

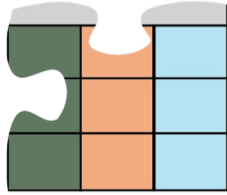


Fragmented datasets

Tables on a similar topic that contain different yet related items.

Fragmented: Unpaid mine safety violations



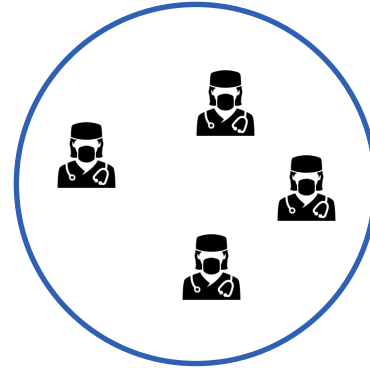


Disparate datasets

Tables that are topically dissimilar and seemingly unrelated.

Disparate: Opioid overdoses

Healthcare workers' death records from opioid overdoses



THE SPOKESMAN-REVIEW

Spokane, Washington

Est. May 19, 1883

Washington Idaho

NEWS > SPOKANE

Washington nurses, health care workers are dying of opioid overdoses

Sun., Feb. 4, 2018

Icons by
[Minh Do](#) and
[Sascha Elmers](#),
Noun Project

Dirty Data in the Newsroom

Comparing Data Preparation in Journalism and Data Science

ACM CHI Conference on Human Factors in Computing Systems
April 23-28, 2023, Hamburg, Germany

Contributions:

- Augmented model of preparation activities
- New model-discrepancy taxonomy of dirty data
- Four challenges in multi-table data integration



Stephen Kasica
University of British Columbia
Vancouver, Canada



Charles Berret
Linköping University
Norrköping, Sweden



Tamara Munzner
University of British Columbia
Vancouver, Canada
[@tamara@vis.social](mailto:tamara@vis.social) 
[@tamaramunzner](https://twitter.com/tamaramunzner) 

