

Quick Overview

- A major challenge in database integration is that of *schema matching* [10, 12], which seeks to determine the schema elements in different databases that correspond to the same real-world entity.
- There is interest in creating schema matching tools that can cope with a broad variety of schema definitions (e.g., relational, XML and OWL) [1, 9].
- Nowadays data is increasingly stored and exchanged in JSON, which is data-centric and often schema-less. To our knowledge, there is **no well-defined system** for matching the schema of JSON files.
- My contributions are: (1) create a novel (baseline) approach for matching JSON files using existing tools and methods, (2) evaluate the approach empirically, (3) formalize a few challenges specific to JSON schema matching. Experiments show that current tools are inadequate for JSON integration and work is needed to formally understand the underlying challenges and to create systems that natively support JSON

Schema Matching: Challenges

Missing, Misleading Info

Schema		Same Attribute?	Same Element?
Movies.year	~ Items.year	Yes	Yes
Movies.title	~ Items.name	No	Yes
Locations.name	≠ Items.name	Yes	No

Consider a matching system that looks at whether attributes are equivalent to determine if those attributes have the same real-world meaning. In the first example, this method is helpful, in the second it is not useful, and in the third it is misleading!

JSON Schema Matching: Additional Challenges

Schema	No Schema
<pre>{ "Director" : "Christopher Nolan", "Comments": [{ "ID" : 907, "Content" : "Love it" }, { "ID" : "909", "Content" : "Too scary" }] }</pre>	<pre>{ "Christopher Nolan" : [[907, "Love it"], [909, "Too scary"]] }</pre>

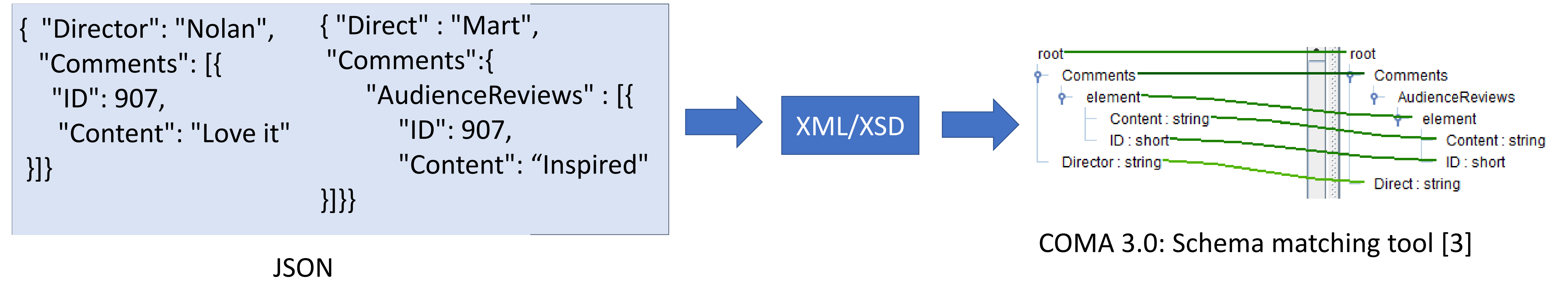
Challenges of matching JSON files includes,

- Different Nesting Structures:** Path to 907 is Comments -> ID in schema one, there is no clear path in the second one.
 - Data Entries Lacking Names and Data Types:** What the data represents is not defined in "No Schema" example and neither states the types of data
- This is an issue as similar nesting structures and data types/names are valuable clues for identifying matches

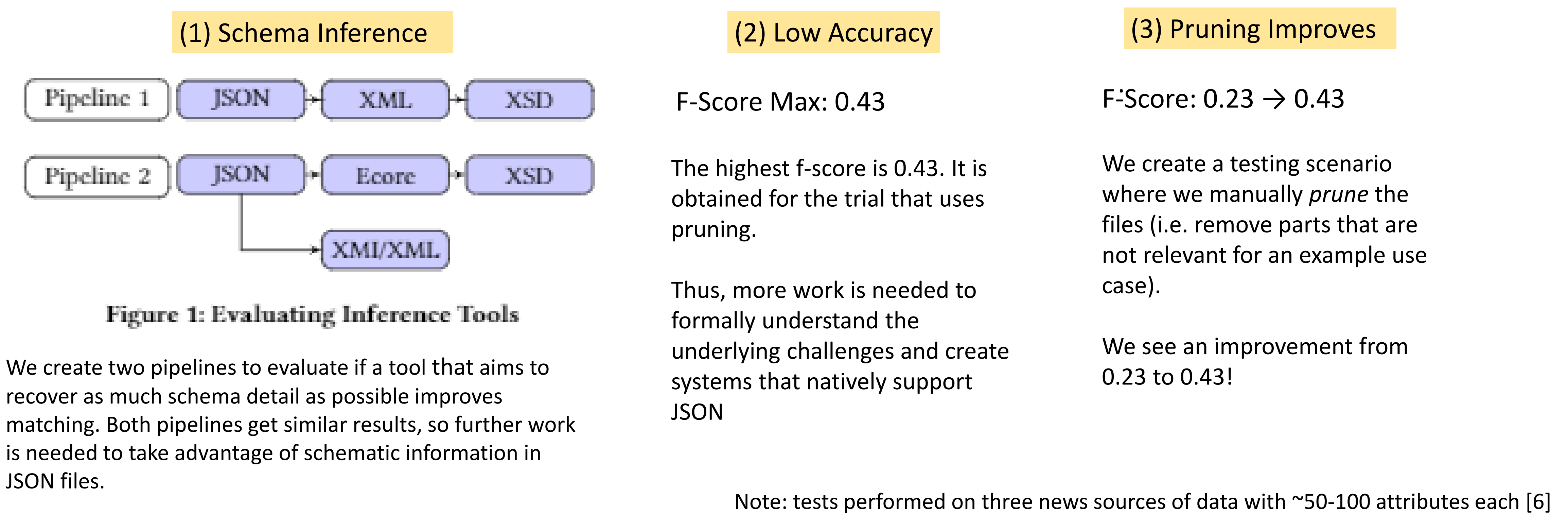
References

- Algergawy, Cheatham, Faria, et al. 2018. Results of the Ontology Alignment Evaluation Initiative 2018. In 13th International Work-shop on Ontology Matching co-located with the 17th ISWC (OM 2018), Vol. 2288. 76–116
- Baazizi, Colazzo, Ghelli, et al. 2019. Schemas and types for json data: From theory to practice. In Proceedings of the 2019 International Conference on Management of Data. ACM, 2060–2063.
- Hong-Hai Do and Erhard Rahm. 2002. COMA: a system for flexible combination of schema matching approaches. In Proceedings of the 28th international conference on Very Large Data Bases. VLDB Endowment, 610–621.
- AnHai Doan, Alon Halevy, and Zachary Ives. 2012. Principles of data integration. Elsevier.
- LiHong He, Chao Han, Arjun Mukherjee, Zoran Obradovic, and Eduard Dragut. 2019. On the dynamics of user engagement in news comment media. Wiley Interdisciplinary Reviews: Data Mining Knowledge Discovery (11 2019).
- Lorena Otero-Cerdeira, Francisco J Rodríguez-Martínez, and AlmaGómez-Rodríguez. 2015. Ontology matching: A literature review. Expert Systems with Applications 42, 2 (2015), 949–971.
- M Tamer Özsu and Patrick Valduriez. 2011. Principles of distributed database systems. Springer Science & Business Media.
- Erhard Rahm and Philip A Bernstein. 2001. A survey of approaches to automatic schema matching. the VLDB Journal 10, 4 (2001), 334–350.

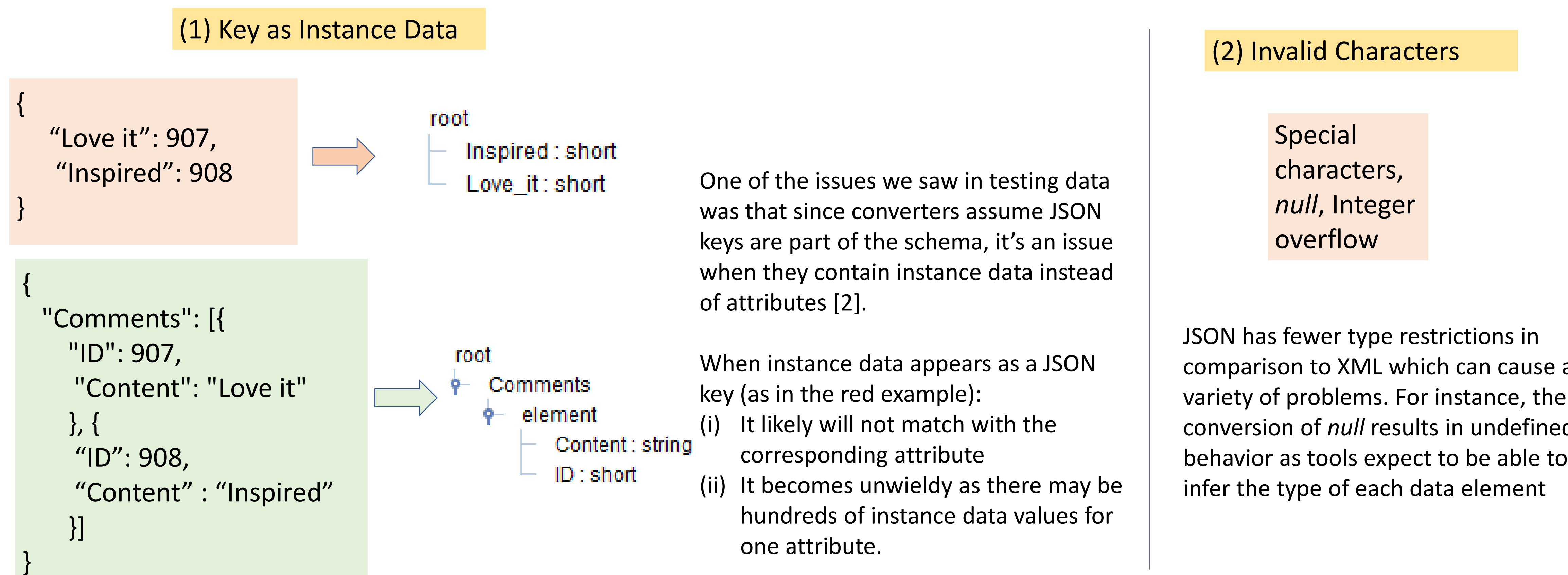
Proposed (Baseline) Approach



Empirical Results



Formalized Challenges



Acknowledgements

I would like to thank Dr. Eduard Dragut for his mentoring on this work and Dr. AnHai Doan for his valuable feedback. Research grant support was provided by NSF CNS-1757533 and BigData-1838145

