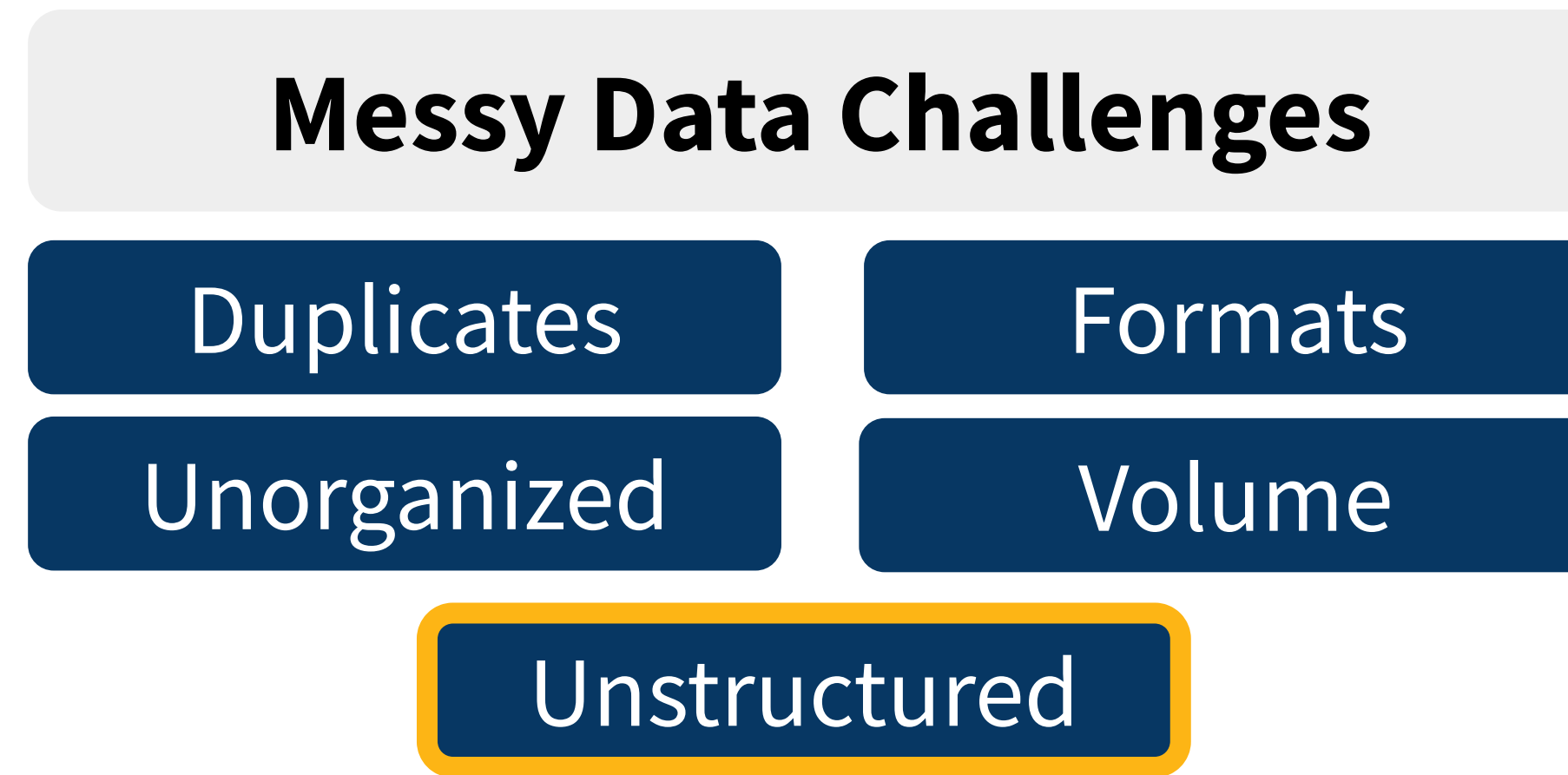
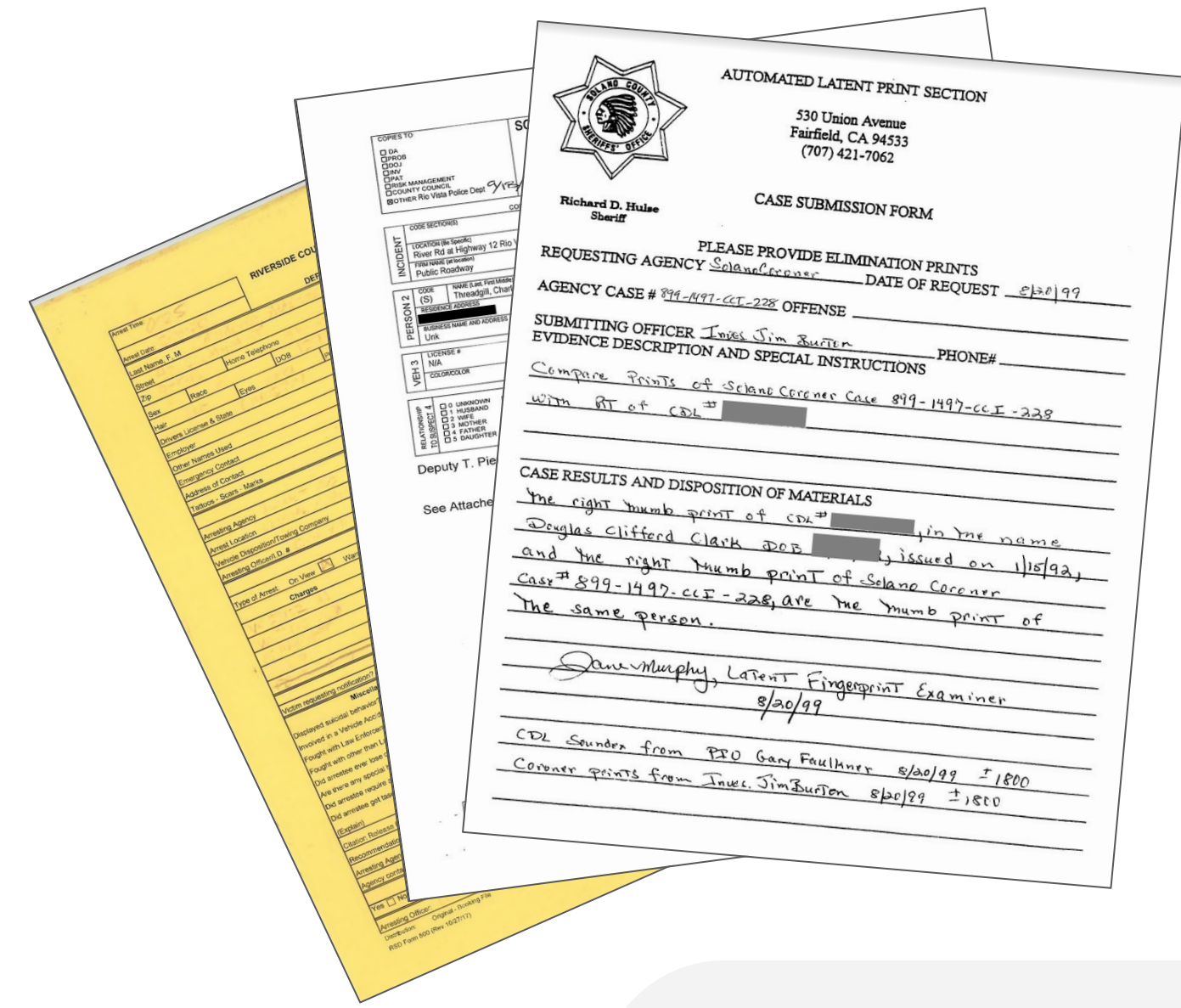


# Discovering Structure in Messy Data

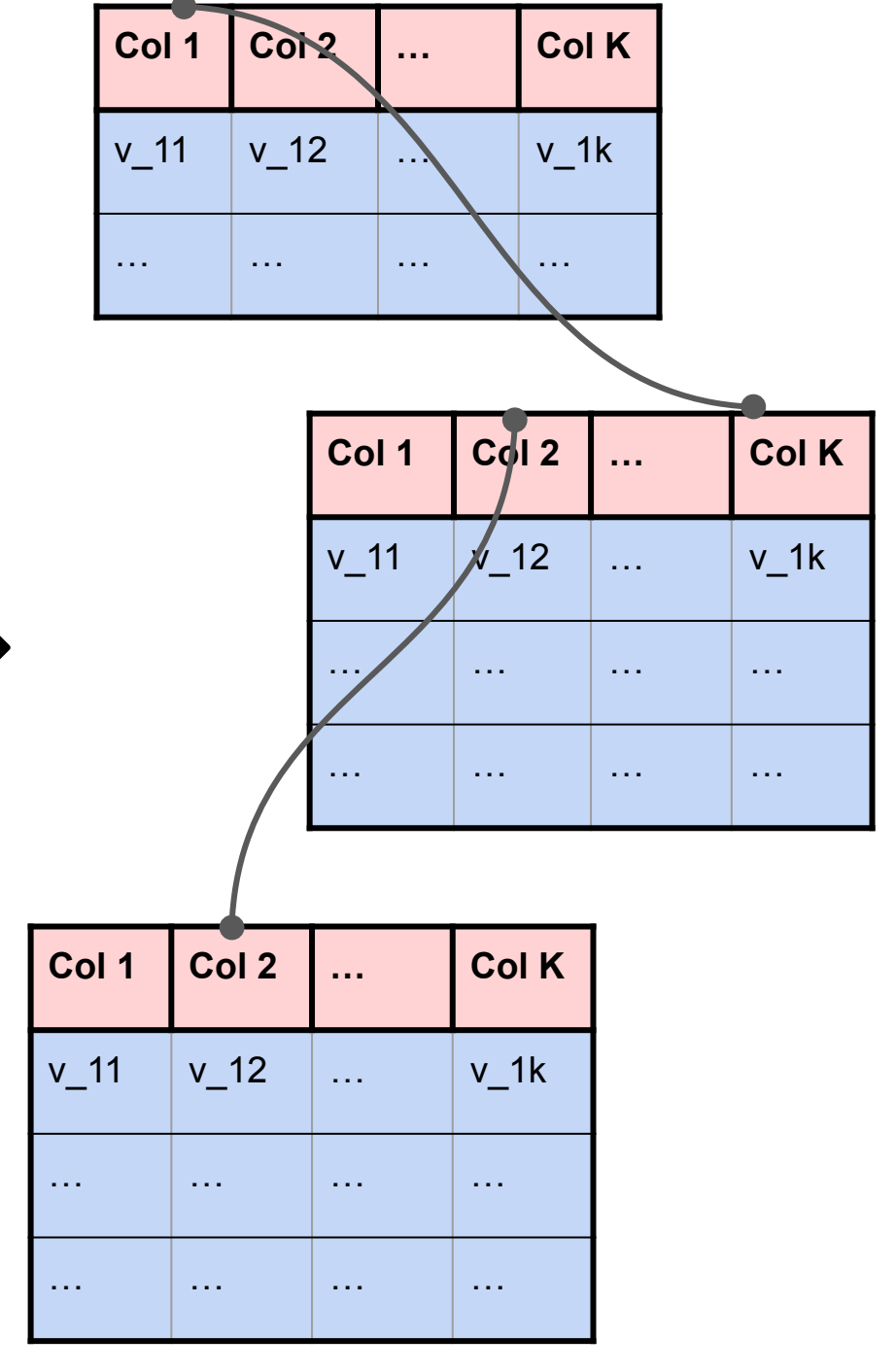
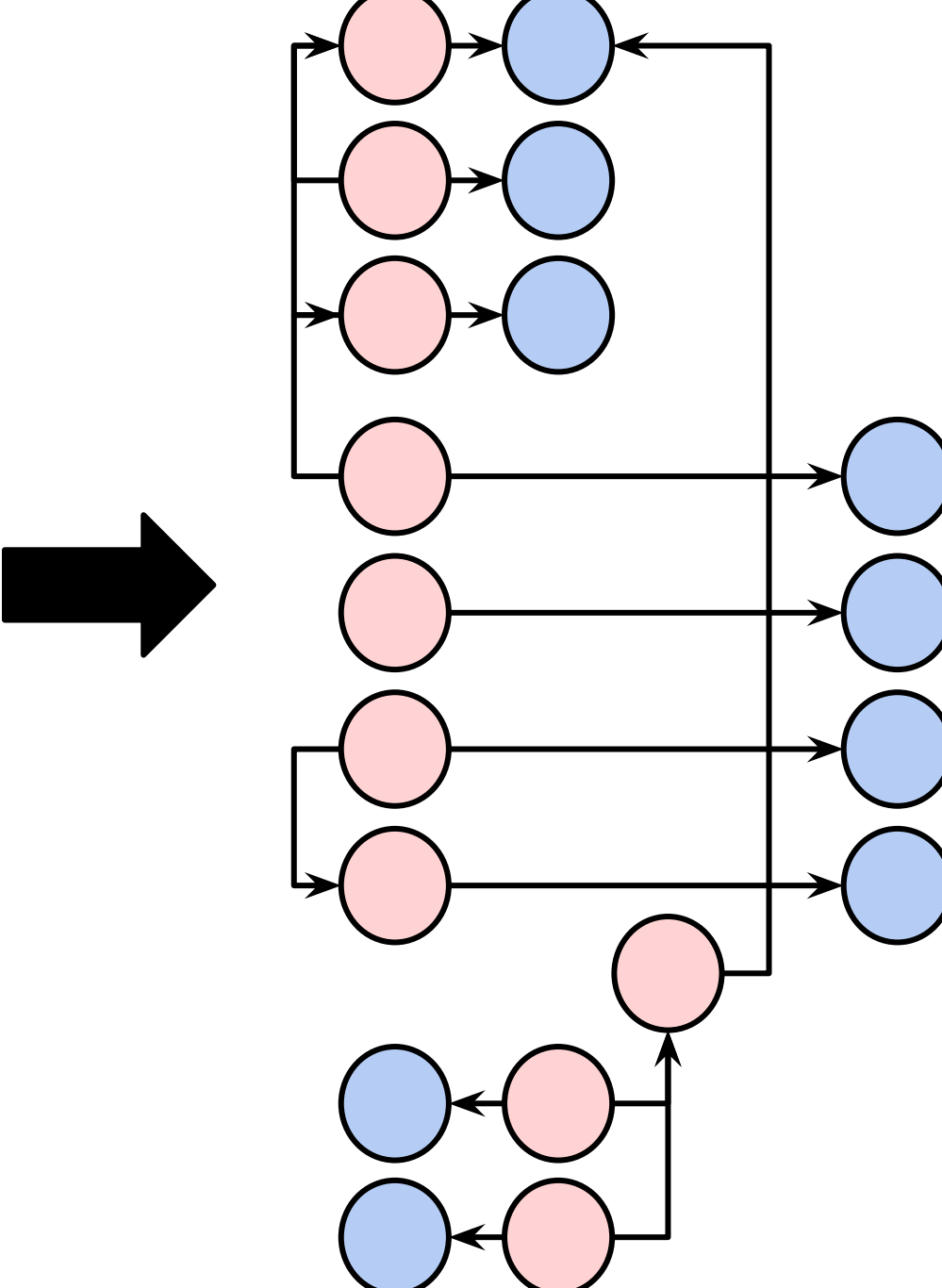
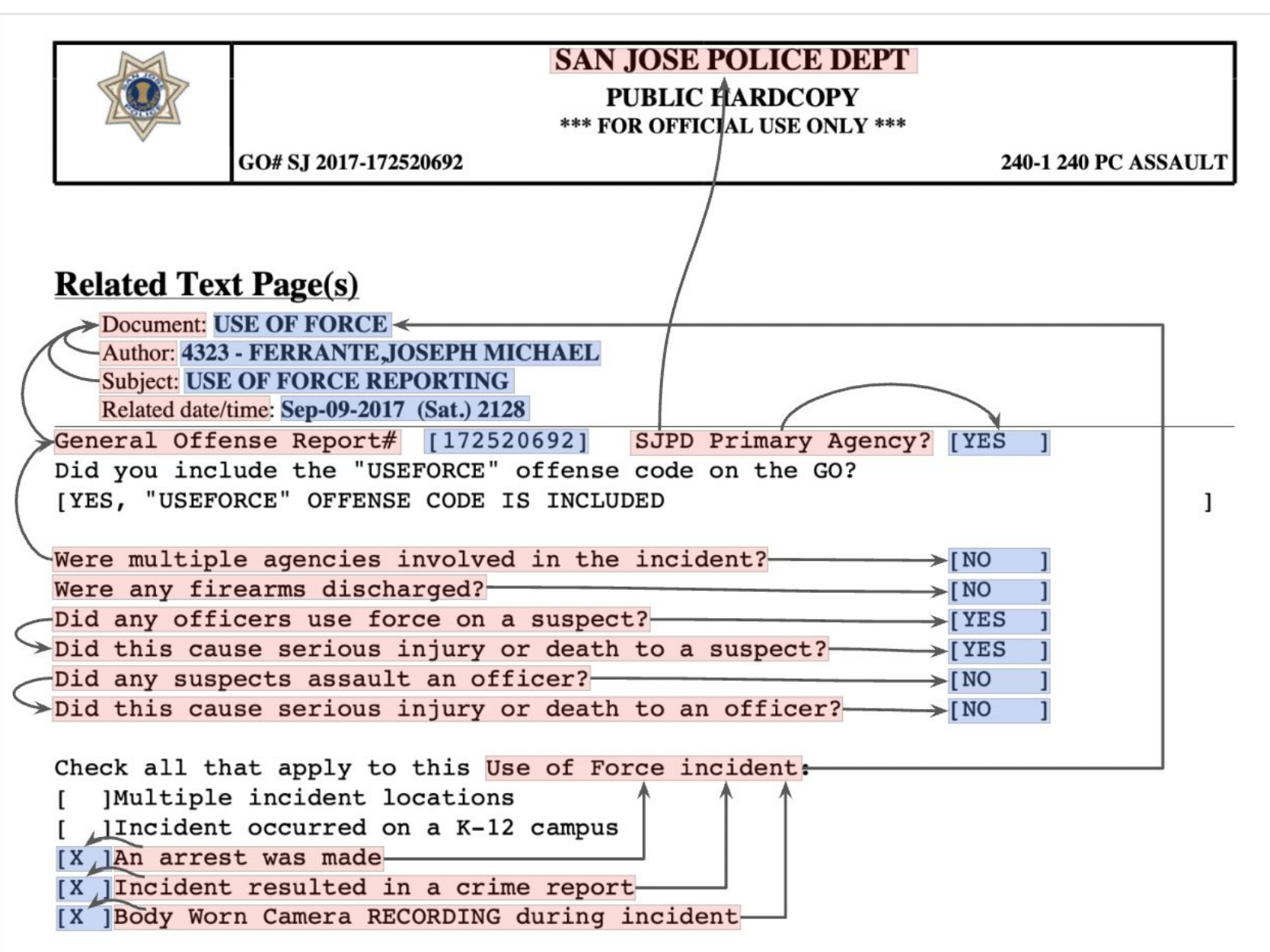
Manish Shetty, Aditya G. Parameswaran, Koushik Sen



- State-of-the-art**
1. Not robust to format changes 📄
  2. Do not go beyond identifying isolated entities 🌴
  3. Do not understand repeated structures 🔄
  4. Do not completely utilize spatial and semantic cues 🌐
  5. Require significant manual labeling of data 🧑‍🔬

💡 **Most unstructured documents originate from structured data** = 📊 📄 🗄️

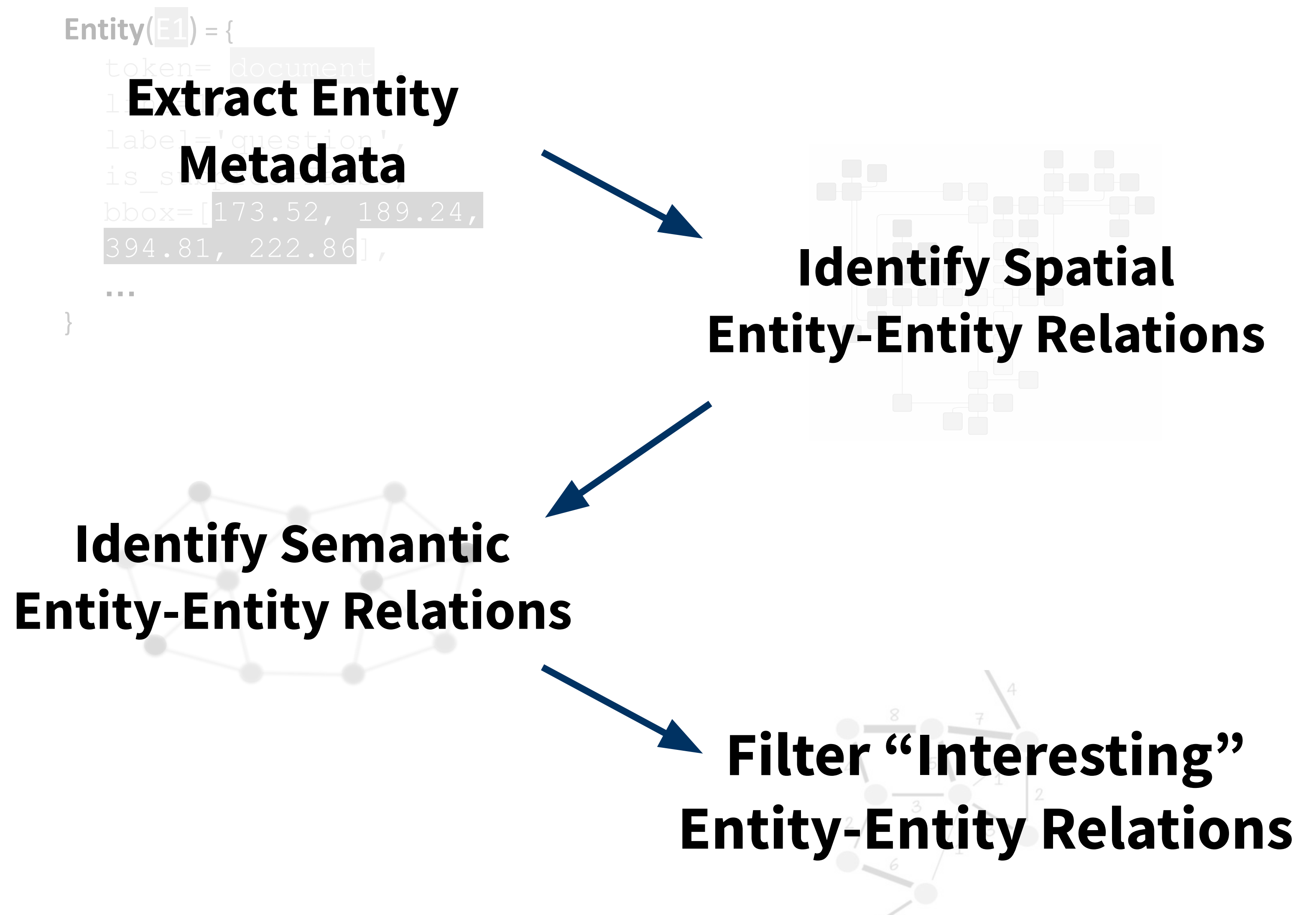
*Can we automatically identify structure in unstructured data?*



**Messy PDF Document**

**Entity Graph**

**Structured Tables**



**Domain Specific Language for PDF data extraction**

```
@start Graph<> G := ConstructGraph(t, ..., t)
Triad <Entity, Entity, string>[] t := AllEntityPairs()
| LeftSibling(e)
| RightSibling(e)
| VerticallyAbove(e)
| VerticallyBelow(e)
| select

Triad <Entity, Entity, string>[] select := IsQnA(t)
| IsLabel(lb, lb, t)
| IsMaxDist(d, t)
| IsFixedLoc(x1, x2, y1, y2, t)
| IsRelatedGPT3(t)
| Conjunction(select, select)
```

Entity e  
float x<sub>1</sub>, x<sub>2</sub>, y<sub>1</sub>, y<sub>2</sub>  
string s, lb  
int d

@input PDFEntityMatrix

**Predictive Synthesis to alleviate annotation overhead!!**