# MetaPAD: Meta Pattern Discovery from Massive Text Corpora

Meng Jiang[1], Jingbo Shang[1], Talyor Cassidy[2], Xiang Ren[1],
Lance M. Kaplan[2], Timothy P. Hanratty[2], Jiawei Han[1]
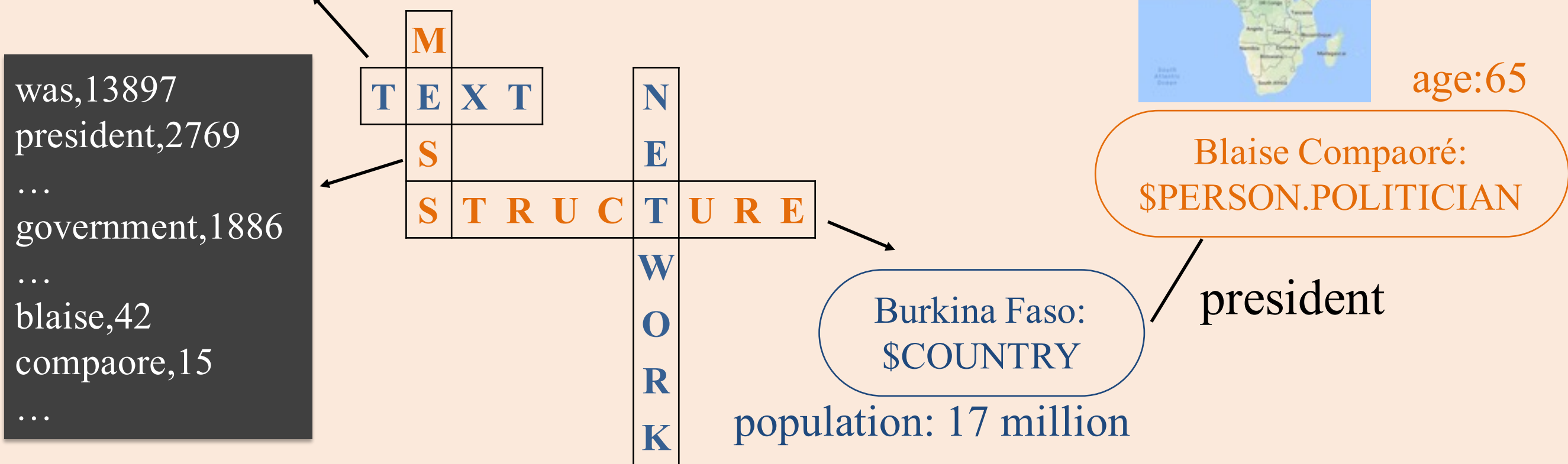[1] Department of Computer Science, University of Illinois at Urbana-Champaign, IL, USA
[2] Computational & Information Sciences Directorate, Army Research Laboratory, Adelphi, MD, USA

## Motivation

**Task:** *Fact extraction* from massive corpora (e.g., news, tweets, papers) to facilitate heterogeneous information network construction

Given a sentence "President **Blaise Compaoré**'s government of **Burkina Faso** was founded…", …

was,13897
president,2769
…
government,1886
…
blaise,42
compaore,15
…

Blaise Compaoré: $PERSON.POLITICIAN  age:65

Burkina Faso: $COUNTRY  president

population: 17 million

**Task 1:** *(entity, attribute name, attribute value)-tuple extraction*
(Burkina Faso, president, Blaise Compaoré)
(Burkina Faso, population, 17 million)
(Blaise Compaoré, age, 65)

**Task 2:** *(entity type, attribute name, value type)-tuple extraction*
($LOCATION.COUNTRY, president, $PERSON.POLITICIAN)
($LOCATION, population, $DIGIT $DIGITUNIT)
($PERSON, age, $DIGIT)

**Idea:** Discovering a group of **synonymous "meta patterns"** to find facts.

president $POLITICIAN's government of $COUNTRY
(e.g., President **Blaise Compaoré**'s government of **Burkina Faso**)
$COUNTRY president $POLITICIAN
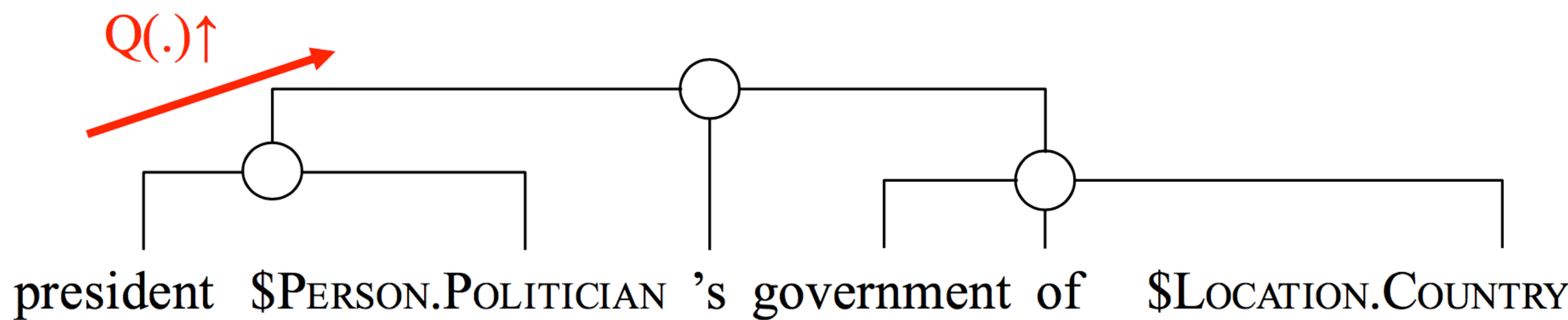$COUNTRY 's president $POLITICIAN
president $POLITICIAN of $COUNTRY
…

(Burkina Faso, president, Blaise Compaoré)
(U.S., president, Barack Obama)
…

## The MetaPAD Framework

(#1) "President Blaise Compaoré's government of Burkina Faso was founded …"
(#2) "President Barack Obama's government of U.S. claimed that…"
(#3) "U.S. President Barack Obama visited …"

1. Meta pattern generation via context-aware segmentation

Meta patterns:
⌈president $PERSON.POLITICIAN 's government of $LOCATION.COUNTRY⌋ was founded…
⌈$LOCATION.COUNTRY president $PERSON.POLITICIAN⌋ …

3. Adjust types for appropriate granularity

⟨$COUNTRY, {president}, $POLITICIAN⟩

Joint extraction

2. Group synonymous meta patterns

⟨Burkina Faso, {president}, Blaise Compaoré⟩
⟨U.S., {president}, Barack Obama⟩

*mutual enhancement:* Generating meta patterns: c(pattern) >> c(tuple)
Grouping synonymous meta patterns: #(patterns) ↑ ➔ #(tuples)↑

## Step 0: Preprocessing

"President Blaise Compaoré's government of Burkina Faso was founded …"

*Phrase mining (SegPhrase by Liu and Han et al. SIGMOD'15)*

"president blaise_compaoré 's government of burkina_faso was founded …"

*Entity recognition and typing with Distant Supervision (ClusType by Ren and Han et al. KDD'15)*

"president $PERSON 's government of $LOCATION was founded …"

*Fine-grained typing (PLE by Ren and Han et al. KDD'16)*

"president $PERSON.POLITICIAN 's government of $LOCATION.COUNTRY was founded …"

## Step 1: Meta pattern quality assessment and segmentation

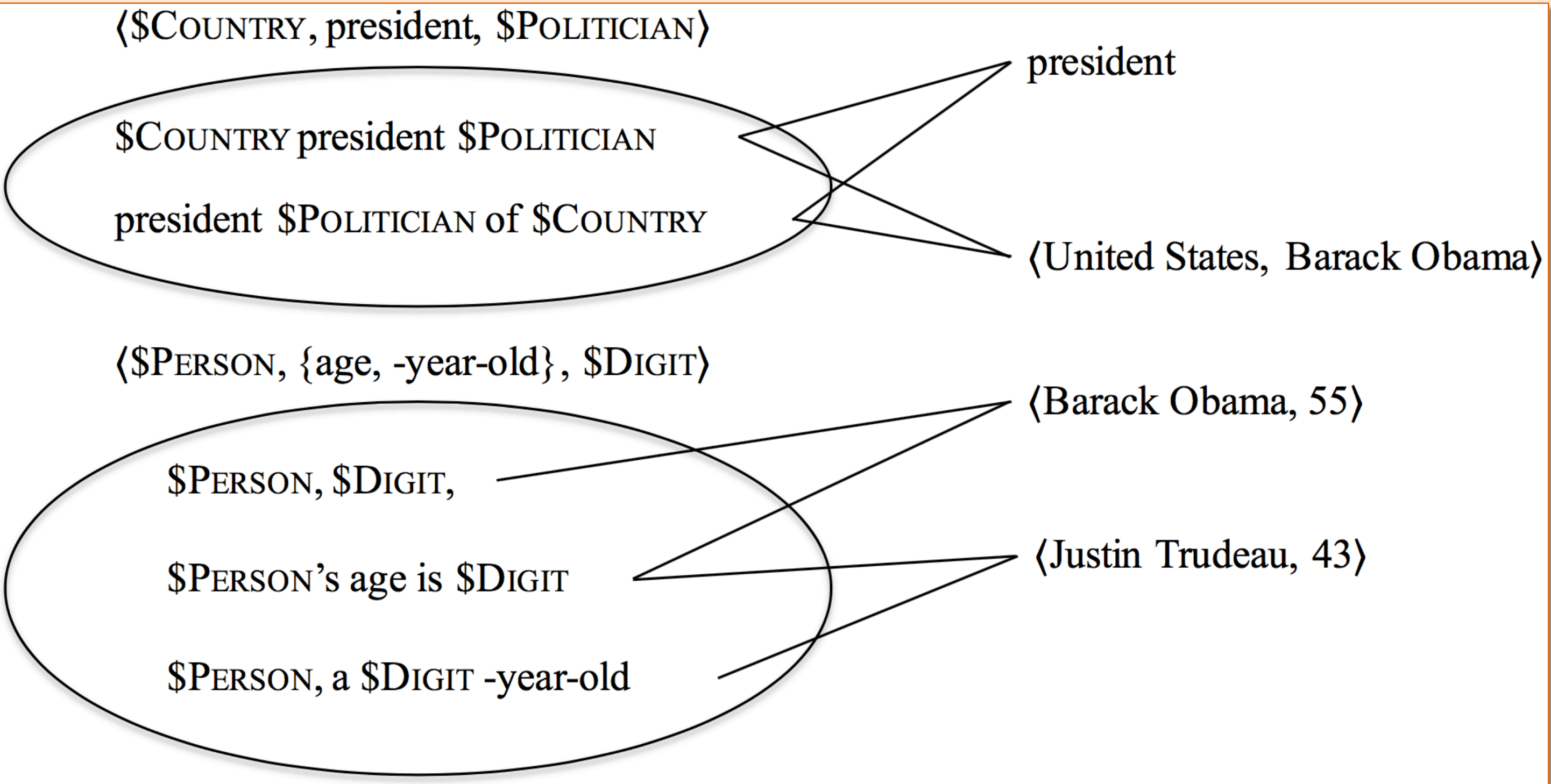A rich set of features:
- ✓ Frequency
- ✓ Concordance: "$PERSON 's wife"
- ✓ Completeness: "$COUNTRY president" vs "$COUNTRY president $POLITICIAN"
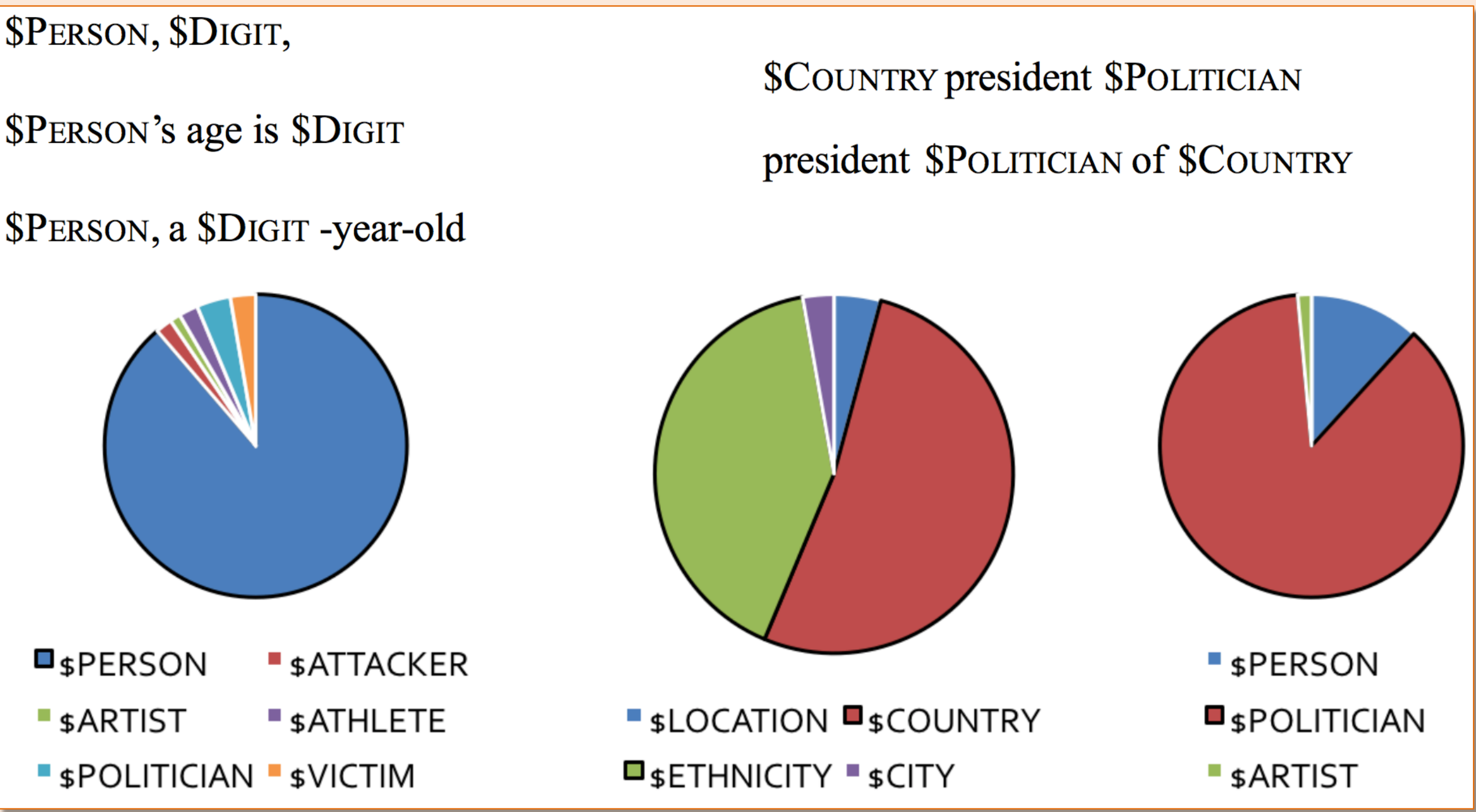- ✓ Informativeness: "$PERSON and $PERSON " vs "$PERSON 's wife, $PERSON"
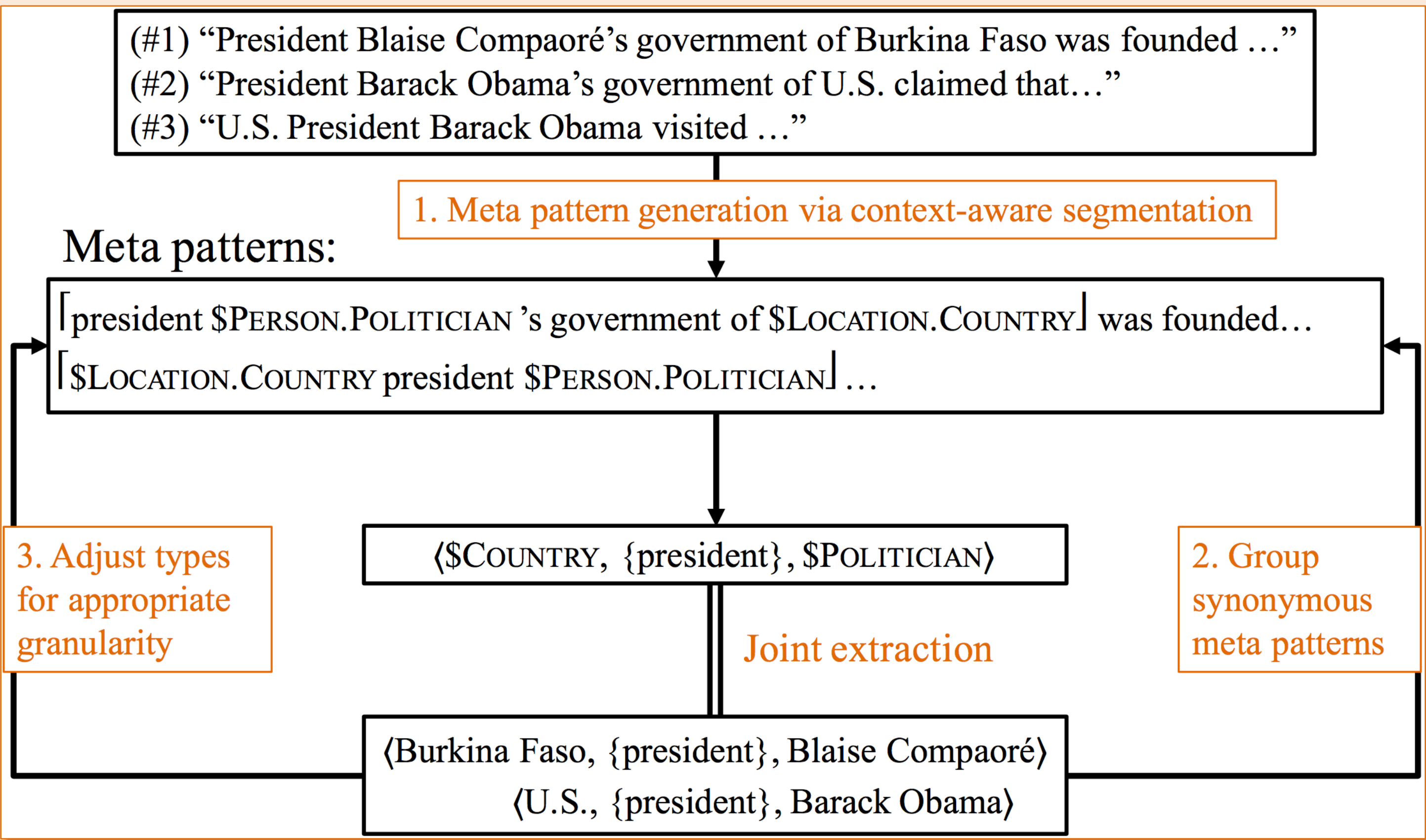
Regression Q(.): random forest with only 300 labels

Q(.)↑

president $PERSON.POLITICIAN 's government of $LOCATION.COUNTRY

## Step 2: Grouping synonymous meta patterns

($COUNTRY, president, $POLITICIAN)
$COUNTRY president $POLITICIAN
president $POLITICIAN of $COUNTRY
→ president
→ ⟨United States, Barack Obama⟩

($PERSON, {age, -year-old}, $DIGIT)
$PERSON, $DIGIT,
$PERSON's age is $DIGIT
$PERSON, a $DIGIT -year-old
→ ⟨Barack Obama, 55⟩
→ ⟨Justin Trudeau, 43⟩

## Step 3: Adjusting types for appropriate granularity

$PERSON, $DIGIT,
$PERSON's age is $DIGIT
$PERSON, a $DIGIT -year-old

$COUNTRY president $POLITICIAN
president $POLITICIAN of $COUNTRY

Legend (left pie): $PERSON, $ATTACKER, $ARTIST, $ATHLETE, $POLITICIAN, $VICTIM
Legend (middle pie): $LOCATION, $COUNTRY, $ETHNICITY, $CITY
Legend (right pie): $PERSON, $POLITICIAN, $ARTIST

| Meta patterns | Entity | Attribute value |
|---|---|---|
| $COUNTRY president $POLITICIAN | United States | Barack Obama |
| $COUNTRY's president $POLITICIAN | Russia | Vladimir Putin |
| president $POLITICIAN of $COUNTRY | France | Francois Hollande |
| … | … | … |
| president $POLITICIAN's government of $COUNTRY | Burkina Faso | Blaise Compaoré |

| Meta patterns | Entity | Attribute value |
|---|---|---|
| $COMPANY CEO $PERSON | Apple | Tim Cook |
| $COMPANY chief executive $PERSON | Facebook | Mark Zuckerberg |
| $PERSON, the $COMPANY CEO, | Hewlett-Packard | Carly Fiorina |
| … | … | … |
| $COMPANY former CEO $PERSON | Infor | Charles Phillips |
| $PERSON, the $COMPANY former CEO, | Afghan Citadel | Roya Mahboob |

| Meta patterns | Entity | Attribute value |
|---|---|---|
| $BACTERIA was resistant to $ANTIBIOTICS | corynebacterium striatum BM4687 | gentamicin |
| $BACTERIA are resistant to $ANTIBIOTICS | methicillin-susceptible S aureus | vancomycin |
| $BACTERIA is the most resistant to $ANTIBIOTICS | multidrug-resistant enterobacteriaceae | gentamicin |
| $BACTERIA, particularly those resistant to $ANTIBIOTICS | | |

| Meta patterns | Entity | Attribute value |
|---|---|---|
| $TREATMENT was used to treat $DISEASE | zoledronic acid therapy | Paget's disease of bone |
| $DISEASE using the $TREATMENT | bisphosphonates | osteoporosis |
| $TREATMENT has been used to treat $DISEASE | calcitonin | Paget's disease of bone |
| $TREATMENT of patients with $DISEASE | | |
| … | | |