

View-Based Tree-Language Rewritings

Laks Lakshmanan, Alex Thomo

University of British Columbia, Canada

University of Victoria, Canada

Importance of trees – XML

- Semi-structured textual formats are very popular.

```
<movie>
  <title>House of cards</title>
  <year>2013</year>

  <character>
    <name>Francis</name>
    <actor>Kevin Spacey</actor>
  </character>

  <character>
    <name>Claire</name>
    <actor>Robin Wright</actor>
  </character>
</movie>
```

XML (Multi TB) success stories:

1. Elsevier
 - Papers and books
2. JPMorgan Chase & Co
 - Stock research data
3. JetBlue Airways
 - Document management

Source: MarkLogic

XML Impacting the Enterprise Tapping into the Power of XML: Five Success Stories

Importance of trees – JSON

- Semi-structured textual formats are very popular.

```
"movie": {  
  "title": "House of cards",  
  "year": "2013",  
  "character": [  
    {  
      "name": "Francis",  
      "actor": "Kevin Spacey"  
    },  
    {  
      "name": "Claire",  
      "actor": "Robin Wright"  
    }  
  ]  
}
```

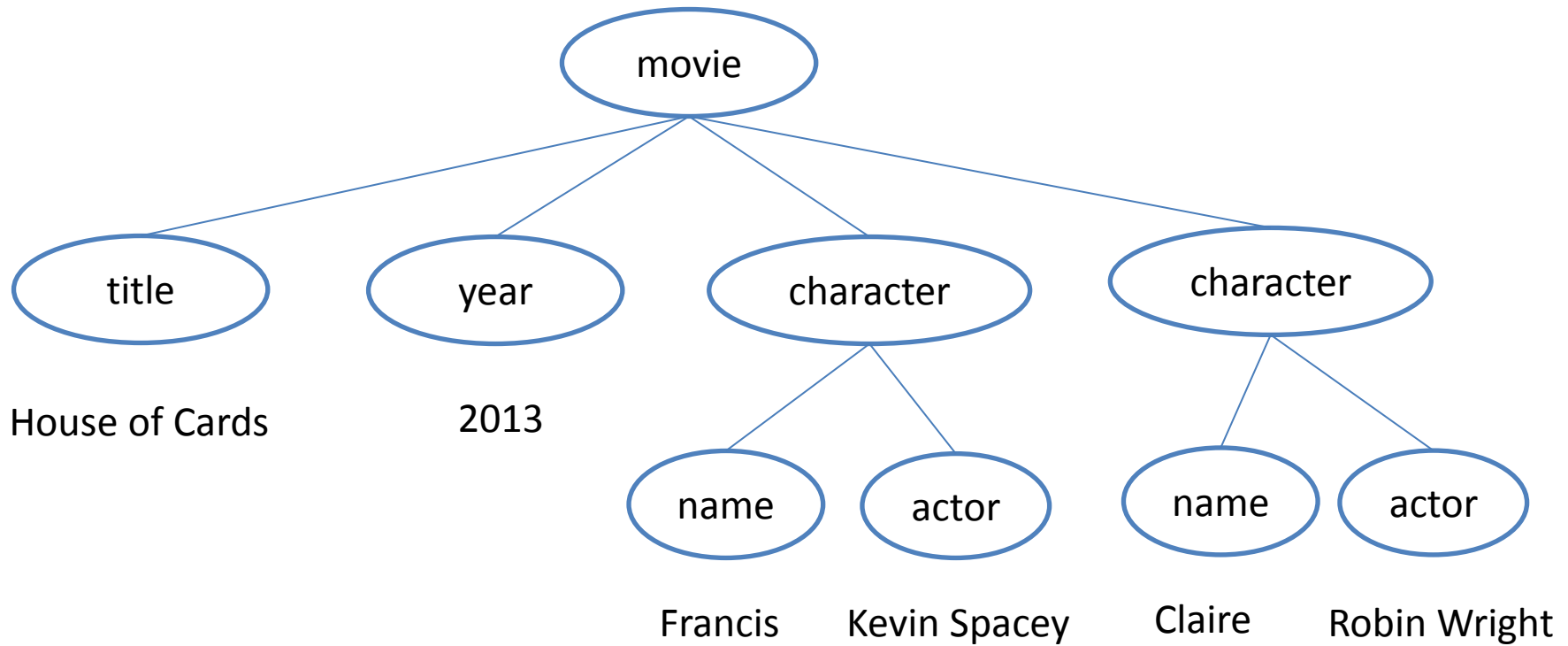
JSON (Multi TB) success stories:

1. CouchDB
2. MongoDB
3. Jaql and Hive JSON SerDe for Hadoop

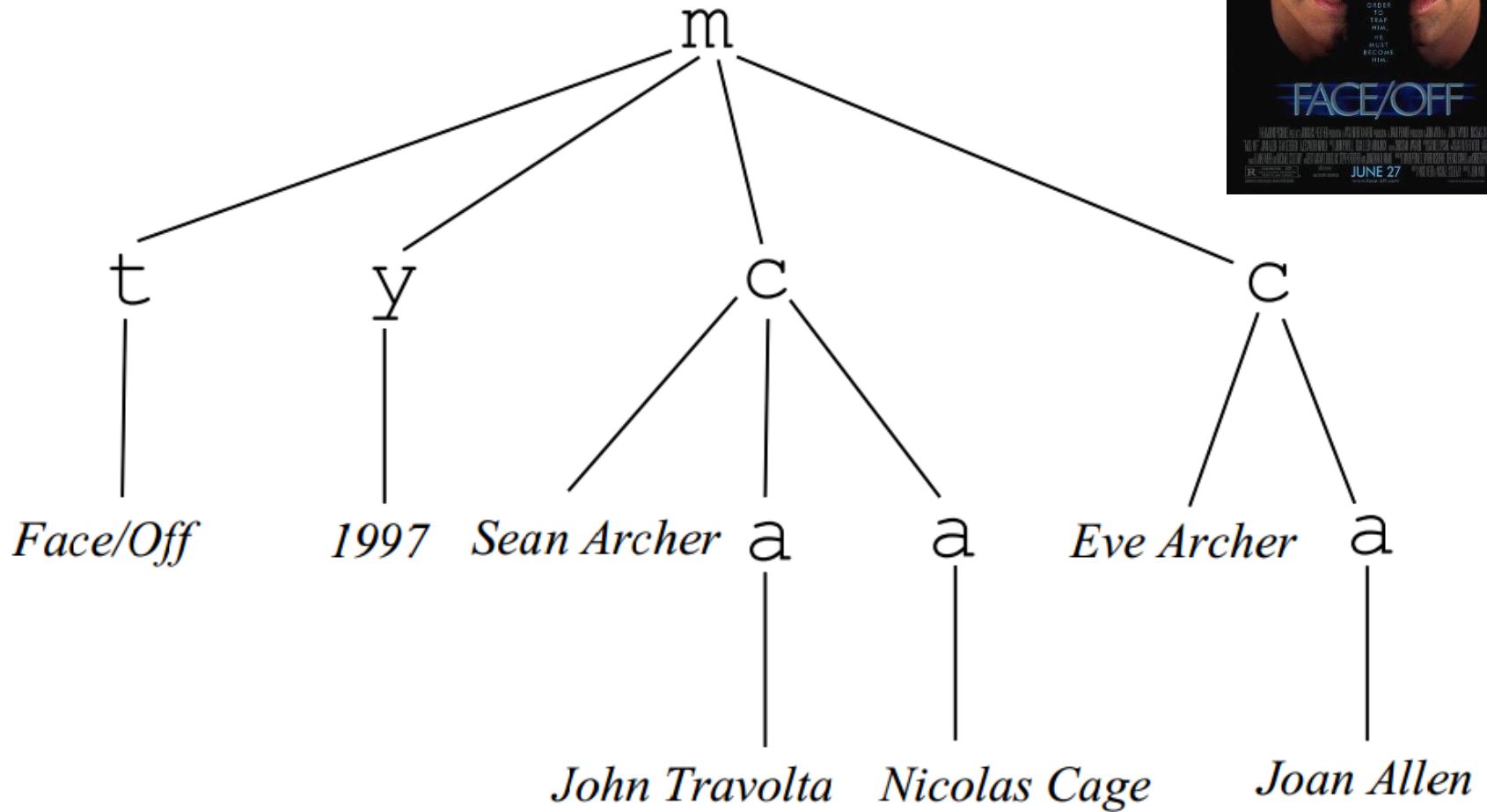
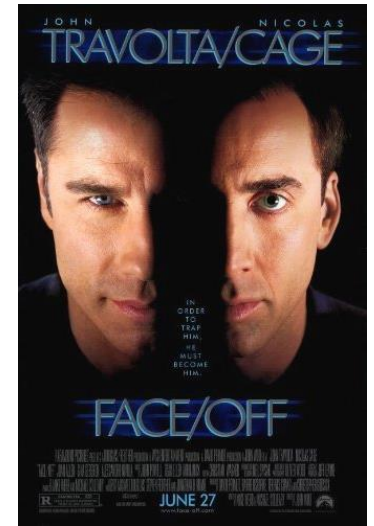
Mantra:

“Log first, ask questions later”

Trees visually



Another example



Importance of views (example)

- **Big database** of movies in a super-tree,
 - each movie being a sub-tree
- Query asks for all the movie sub-trees with a **MAC**.
 - small minority; number about 50.
 - Result **materialized into a view**.
- **Tremendous help** in answering new queries, e.g.
 - “find actors playing a MAC”.
 - Rewrite into: “find actors playing a MAC in a movie having a MAC”
 - answer it **on the materialized view**.

Regular Expressions and Automata

- Return all movie actors

$_ * m _ * c \hat{a}$

- Automaton

$$s_0 \xrightarrow{\Sigma} s_0$$

$$s_0 \xrightarrow{m} s_1$$

$$s_1 \xrightarrow{\Sigma} s_1$$

$$s_1 \xrightarrow{c} s_2$$

$$s_2 \xrightarrow{\hat{a}} s_3$$

A pattern



Reverse

- Return all movie actors

$_ * m _ * c \hat{a}$

- Automaton

$$s \xrightarrow{\hat{a}} s_a$$

$$s_a \xrightarrow{c} s_c$$

$$s_c \xrightarrow{\Sigma} s_c$$

$$s_c \xrightarrow{m} s_m$$

$$s_m \xrightarrow{\Sigma} s_m$$

A pattern



Bottom-up Tree Automata

- Return all movie **actors**

$_ * m _ * c \hat{a}$

- Automaton

$$S \xrightarrow{\hat{a}} S_a \quad S \xrightarrow{\Sigma} S$$

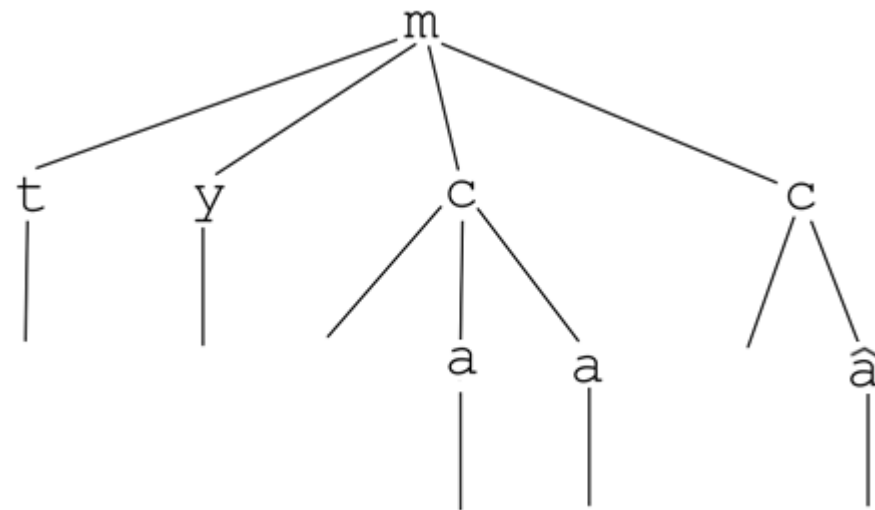
$$S * S_a S * \xrightarrow{c} S_c$$

$$S * S_c S * \xrightarrow{\Sigma} S_c$$

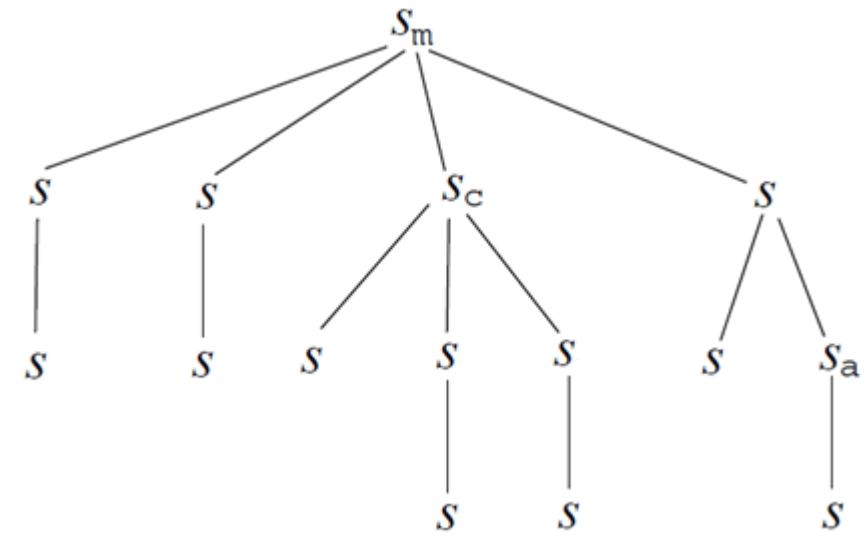
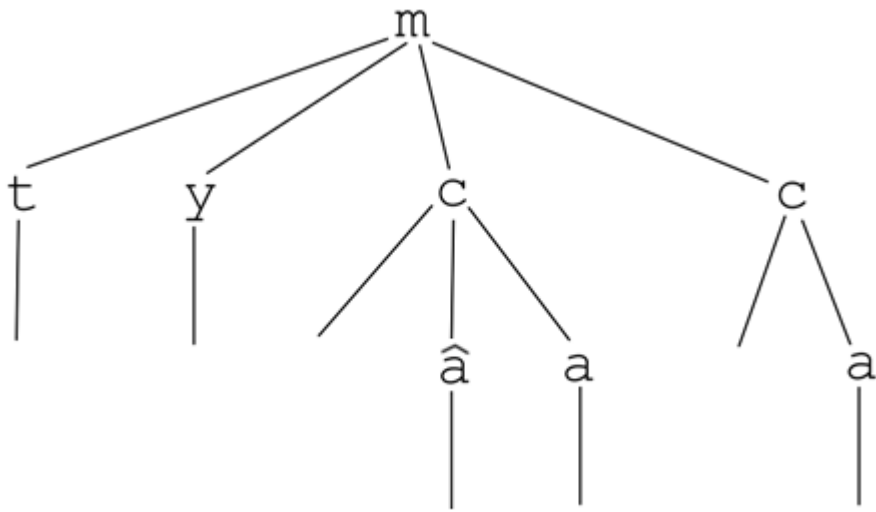
$$S * S_c S * \xrightarrow{m} S_m$$

$$S * S_m S * \xrightarrow{\Sigma} S_m$$

A pattern



Run



Bottom-up Tree Automata (II)

- Return all movie **actors of MACs**

$$_{}^* m _{}^* c_{MAC} \hat{a}$$

- Automaton

$$s \xrightarrow{\hat{a}} s_{\hat{a}} \quad s \xrightarrow{a} s_a \quad s \xrightarrow{\Sigma} s$$

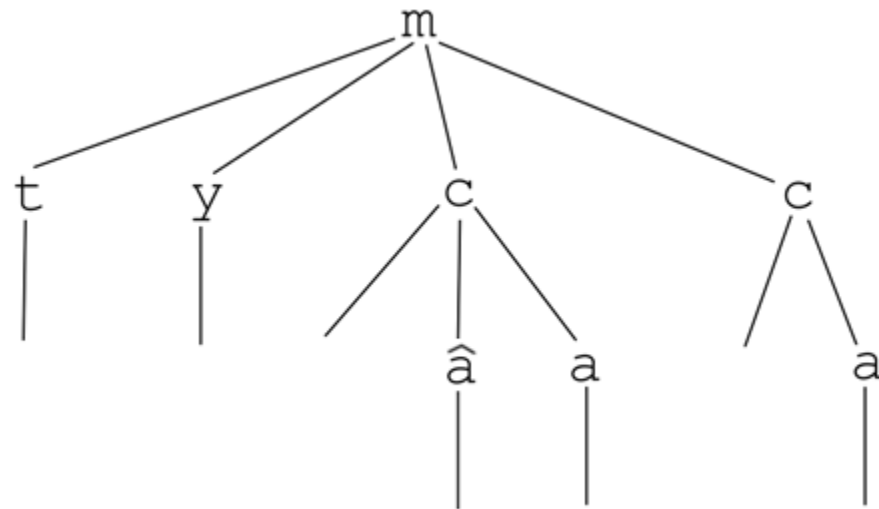
A pattern

$$s^* (s_{\hat{a}} s_a \mid s_a s_{\hat{a}}) s^* \xrightarrow{c} s_c$$

$$s^* s_c s^* \xrightarrow{\Sigma} s_c$$

$$s^* s_c s^* \xrightarrow{m} s_m$$

$$s^* s_m s^* \xrightarrow{\Sigma} s_m$$



Bottom-up Tree Automata (IV)

- Return all **movies** having some **MACs**

$$_ * \hat{m} _ * c_{MAC} a$$

- Automaton

$$s \xrightarrow{a} s_a \quad s \xrightarrow{\Sigma} s$$

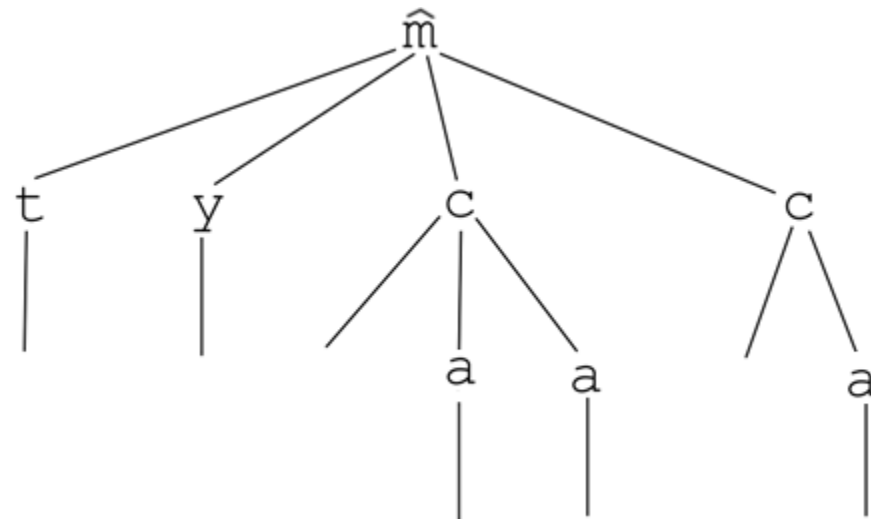
$$s * s_a s_a s * \xrightarrow{c} s_c$$

$$s * s_c s * \xrightarrow{\Sigma} s_c$$

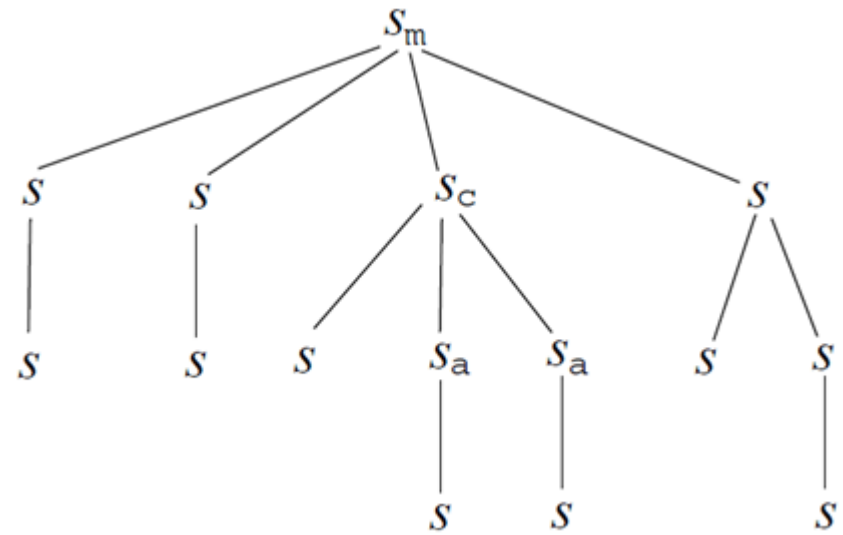
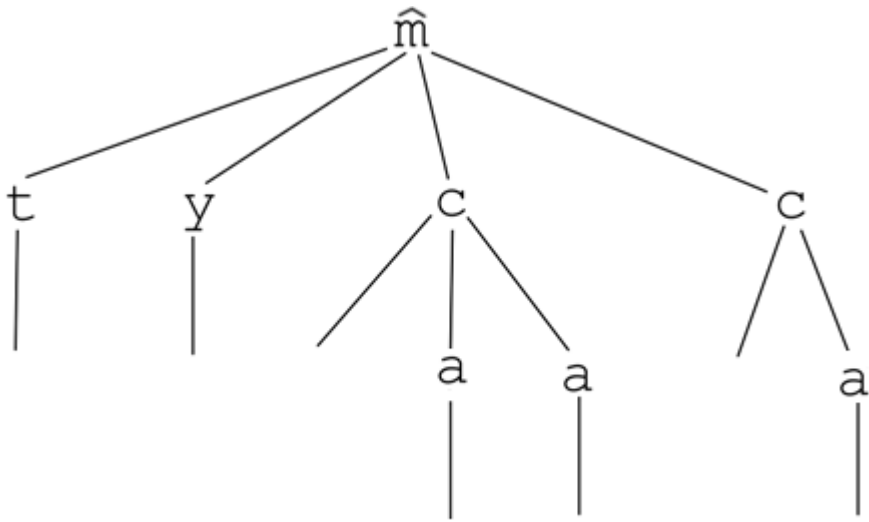
$$s * s_c s * \xrightarrow{\hat{m}} s_m$$

$$s * s_m s * \xrightarrow{\Sigma} s_m$$

A pattern



Run



Bottom-up Tree Automata (V)

- **Regular tree languages** (RTAs)
 - the sets of trees recognized by TAs.
 - closed under **intersection** and **complement**
- **Deterministic TA**
 - For any tree t , there can be **at most one** accepting run of A on t .
 - Power-wise, $TA = DTA$.
- **Complement** obtained from deterministic TA
- **Intersection** via a special construction preserves determinism.

Queries

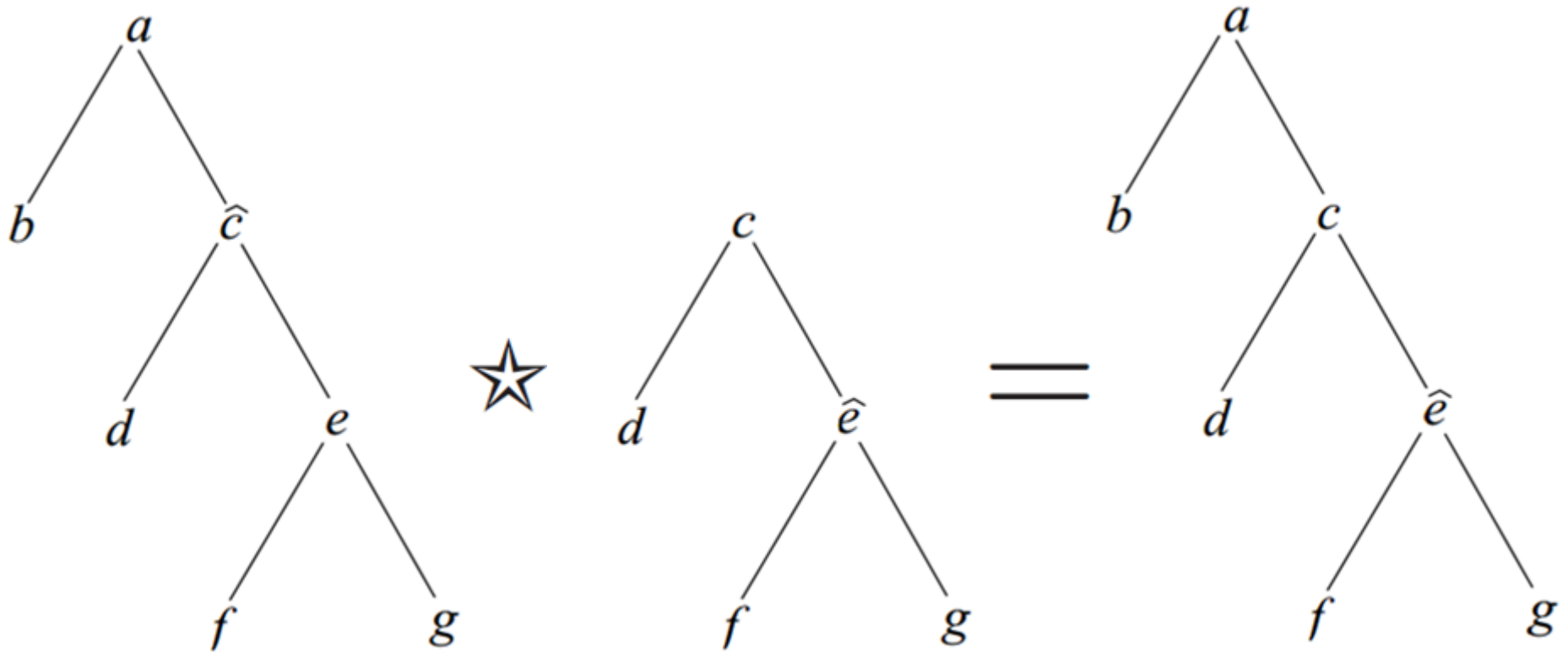
- Queries are **regular** sets of trees over

$$\Sigma \cup \hat{\Sigma}$$

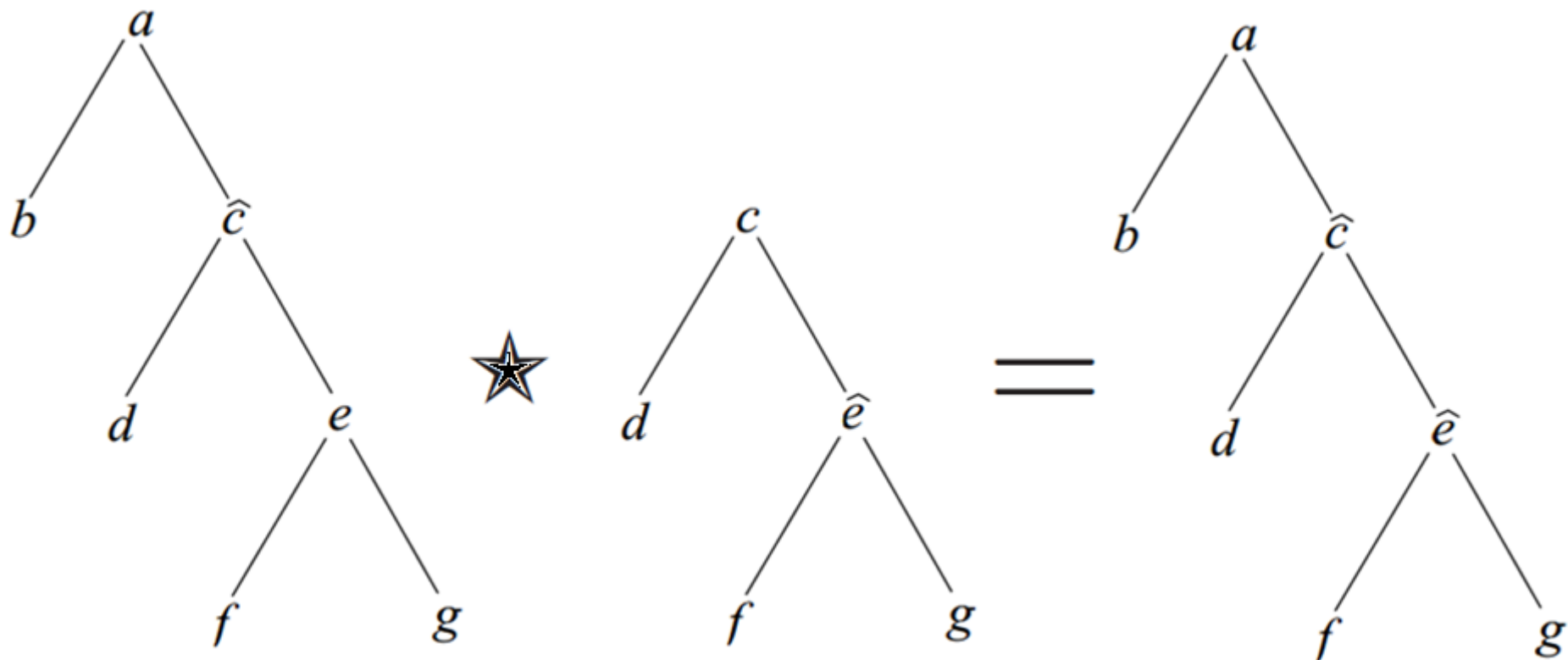
- Containment Lemma

$$Q_1 \subseteq Q_2 \quad \text{implies} \quad \text{ans}(Q_1) \subseteq \text{ans}(Q_2)$$

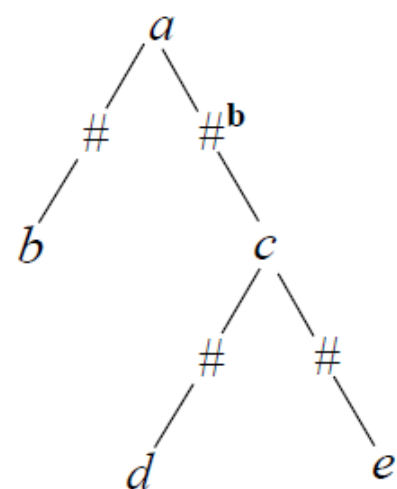
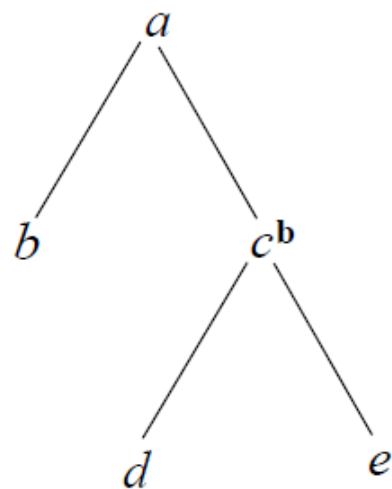
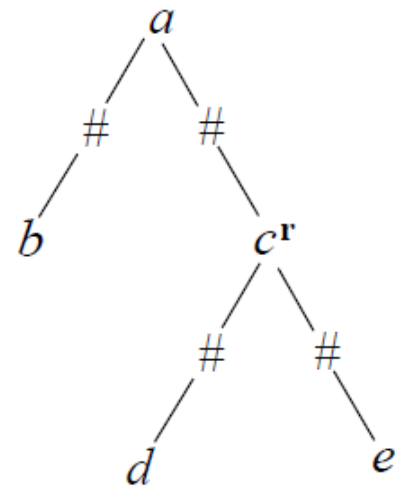
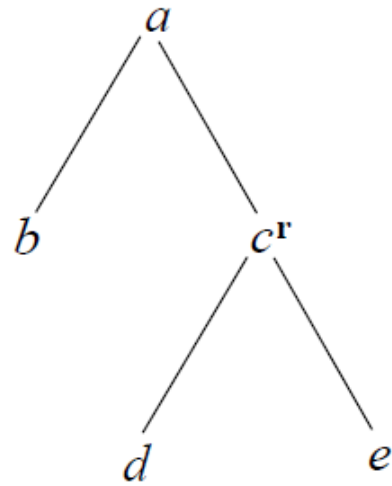
Star Operation



Filled Star Operation



Transformation for avoiding marker overlap



Rewriting, and two sets

Maximally contained **rewriting**:

$$R = \{\xi \in \mathcal{T}^x : V \star \xi \subseteq Q\}$$

The **bad** set:

$$X = \{\xi \in \mathcal{T}^x : \text{there exists } v \in V \text{ such that } v \star \xi \in Q^c\}$$

The **promising** set:

$$Y = \{\xi \in \mathcal{T}^x : \text{there exists } v \in V \text{ such that } v \star \xi \in Q\}$$

Proposition. $R = Y \setminus X$

Example with chains

$$Q = \{aac\hat{c}de, bbc\hat{c}de, ccc\hat{c}de\}$$

$$V = \{aa\hat{c}de, bb\hat{c}de\}$$

$$c\hat{c}de \in R$$

$$V \star c\hat{c}de = \{aac\hat{c}de, bbc\hat{c}de\} \subseteq Q$$

Example with chains (II)

$$Q = \{aac\hat{c}de, ccc\hat{c}de\}$$

$$V = \{aa\hat{c}de, bb\hat{c}de\}$$

$$c\hat{c}de \notin R$$

$$bbc\hat{c}de \in V \star c\hat{c}de$$

$$bbc\hat{c}de \notin Q$$

$$c\hat{c}de \in X$$

$$c\hat{c}de \in Y$$

Inverse of the star operation

$$v \star v' = \begin{cases} \xi & \text{if } v' = v \star \xi \\ \text{undefined} & \text{otherwise} \end{cases}$$

Proposition.

$$X = V \star Q^c \qquad Y = V \star Q$$

Compute $K = J \star J'$ where J and J' are RTQ

Colored Alphabets

- Markers will be colors
 - Blue for J
 - Red for J'

$$\Sigma^{\mathbf{r}} = \{a^{\mathbf{r}} : a \in \Sigma\}$$

$$\Sigma^{\mathbf{b}} = \{a^{\mathbf{b}} : a \in \Sigma\}$$

Colored Languages

γ^b set of all trees having one node blue

γ^r set of all trees having one node red

$\gamma^{b,r}$ set of all trees having one node blue and another red as descendant of the blue node

Φ^r set of all trees having all nodes black, except root which is red

$\Phi^{b,r}$ set of all trees having all nodes black, except for the root which is blue and another node which is red.

Colored Languages (II)

$$J \subseteq \underline{\gamma^b}$$

$$J' \subseteq \underline{\gamma^r}$$

$$K \subseteq \underline{\gamma^r} \quad K \subseteq \underline{\gamma^r} \setminus \Phi^r$$

$$J \star K \subseteq \underline{\gamma^r} \quad J \star K \subseteq \underline{\gamma^{b,r}}$$

Colored Languages (III)

p over $\Sigma \cup \Sigma^r \cup \Sigma^b$

p^{-b} same as p , but with blue nodes turned black

p^{-r} same as p , but with red nodes turned black

L over $\Sigma \cup \Sigma^r \cup \Sigma^b$

$$L^{-b} = \{p^{-b} : p \in L\}$$

$$L^{-r} = \{p^{-r} : p \in L\}$$

Colored Languages (IV)

$$B_L = \{p \in \Upsilon^{b,r} : p^{-b} \in L\} \text{ for } L \subseteq \Upsilon^r$$

$$\mathcal{A} = (S, \Sigma \cup \Sigma^r, F, \Delta) \text{ automaton for } L$$

$$\mathcal{B} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F, \Delta_{\mathcal{B}})$$

$$\Delta_{\mathcal{B}} = \Delta \cup \{H \xrightarrow{a^b} s : H \xrightarrow{a} s \text{ in } \Delta \text{ and } a \in \Sigma\}$$

$$L(\mathcal{B}) \cap \Upsilon^{b,r} = B_L$$

Similarly:

$$B'_L = \{p \in \Upsilon^{b,r} : p^{-r} \in L\} \text{ for } L \subseteq \Upsilon^b$$

$$C_L = \{p \in \Phi^{b,r} : p^{-b} \in L\} \text{ for } L \subseteq \Upsilon^r \setminus \Phi^r$$

Rewriting Algorithm $K = J \star J'$

Compute: $B_J \quad B'_{J'}$

$$B_J \cap B'_{J'}$$

$$\mathcal{D} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F, \Delta) \quad L(\mathcal{D}) \subseteq \Upsilon^{b,r}$$

$$\mathcal{E} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F_{\mathcal{E}}, \Delta) \quad L(\mathcal{E}) \subseteq \Phi^{b,r}$$

$$F_{\mathcal{E}} = \{s \in S : \text{there exists } H \xrightarrow{a^b} s \text{ in } \Delta\}$$

Theorem. $L(\mathcal{E}) = C_K \quad C_K^{-b} = K$

Rewriting Algorithm (II)

$$X = V \star Q^c \qquad Y = V \star Q$$

$$R = Y \setminus X$$

Complexity

- **Proposition.** $K = J \star J'$ can be computed in **polynomial time**.
- **Theorem.** The MCR of Q using V can be computed in **exponential time**.
- **Theorem.** Computing the MCR of Q using V is **EXPTIME-hard**.

Final Notes

- Query automata formalism used is equivalent in power to MSO (**golden standard**)
 - For specifying node-selecting queries.
 - Colors correspond to Boolean markings
 - J. Niehren, L. Planque, J.-M. Talbot, and S. Tison. N-ary queries by tree automata. DBPL, 2005
- XPath rewriting is NP-hard.
- XPath is a subclass of our formalism.
 - Our automata-based algorithm can be used as well for rewriting XPath queries.

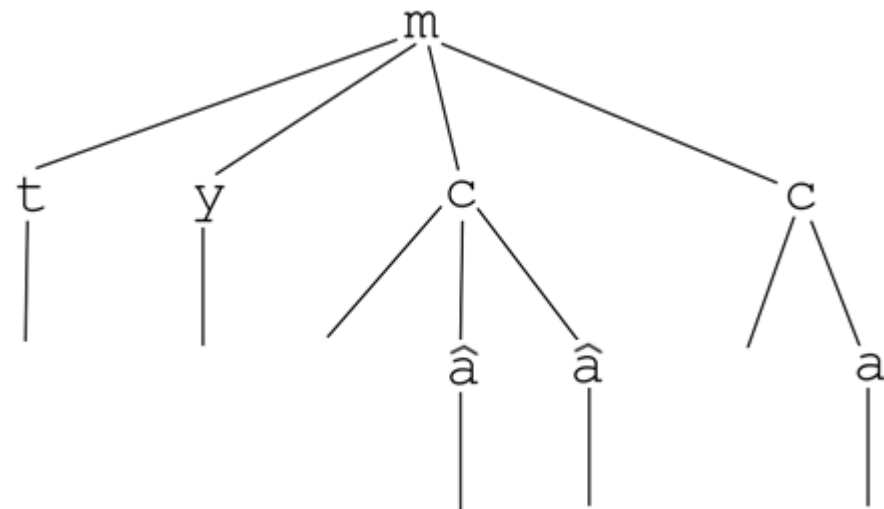
K-ary queries

- **Example:** Find the **2-forests of actor tree pairs** for actors who have played the same character together in some movie.

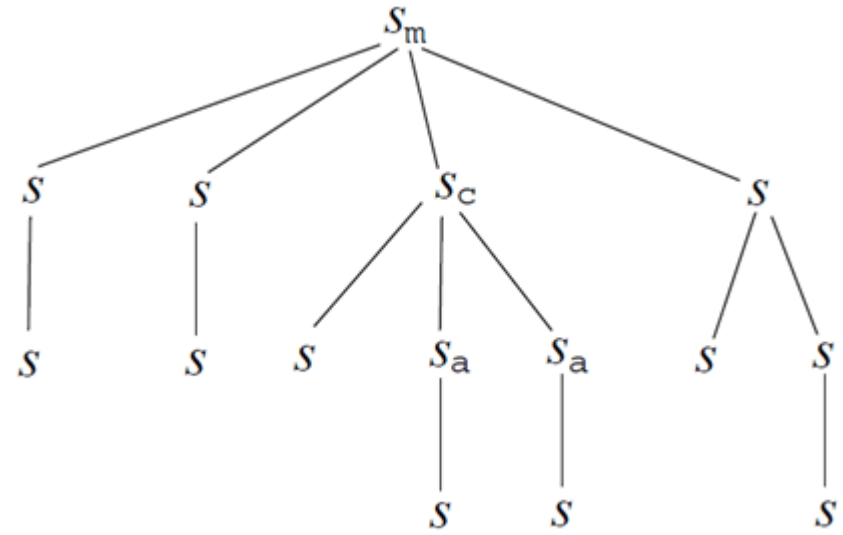
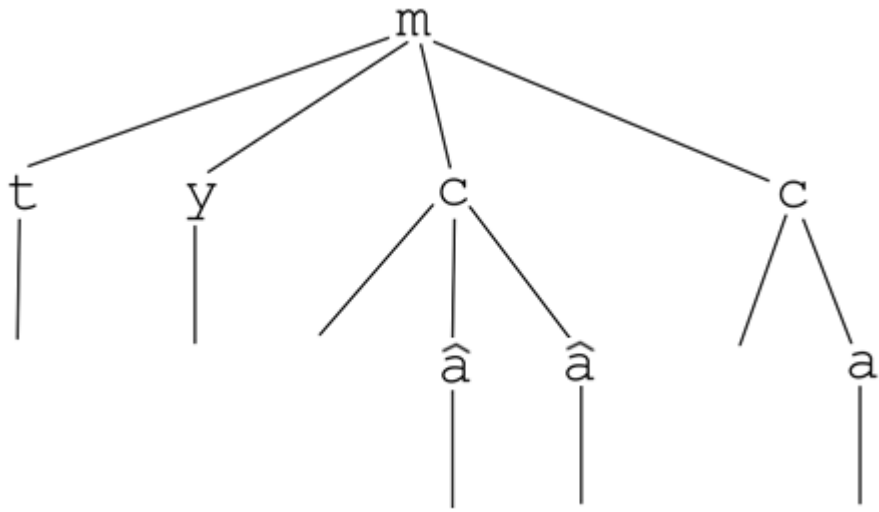
- Automaton

$$\begin{aligned}
 & S \xrightarrow{\hat{a}} S_{\hat{a}} & S \xrightarrow{\Sigma} S \\
 & S^* S_{\hat{a}} S_{\hat{a}} S^* \xrightarrow{c} S_c \\
 & S^* S_c S^* \xrightarrow{\Sigma} S_c \\
 & S^* S_c S^* \xrightarrow{m} S_m \\
 & S^* S_m S^* \xrightarrow{\Sigma} S_m
 \end{aligned}$$

A pattern



Run



Why is rewriting K-ary queries challenging

- It has been shown that k-ary queries can be encoded by unary queries
 - T. Schwentick. On diving in trees. In MFCS, 2000.
 - Done by going through MSO formulas.
 - Going from a k-ary query to an MSO encoding and then back to automata **incurs non elementary complexity**.
- Therefore we need a another algorithm for rewriting k-ary queries
 - that doesn't go via MSO formulas

Conclusions

- Characterized **view-based rewriting** as **solving a lang. equation**
 - Defined appropriate tree operators
- Defined **colored languages**
 - Gave automata constructions
- Computed rewriting as a series of **operations on automata**
- Characterized the **complexity** of computing rewriting
 - **Tight lower bound** provided
- Extended the results to **k-ary queries**
 - Common in XQuery

Thank You