

Crowdsourcing Taxonomies

Network-Centric
NetCInS Systems
Information

Dimitris Karampinas & Peter Triantafillou

*Computer Engineering and Informatics
Department - University of Patras*

INTRODUCTION

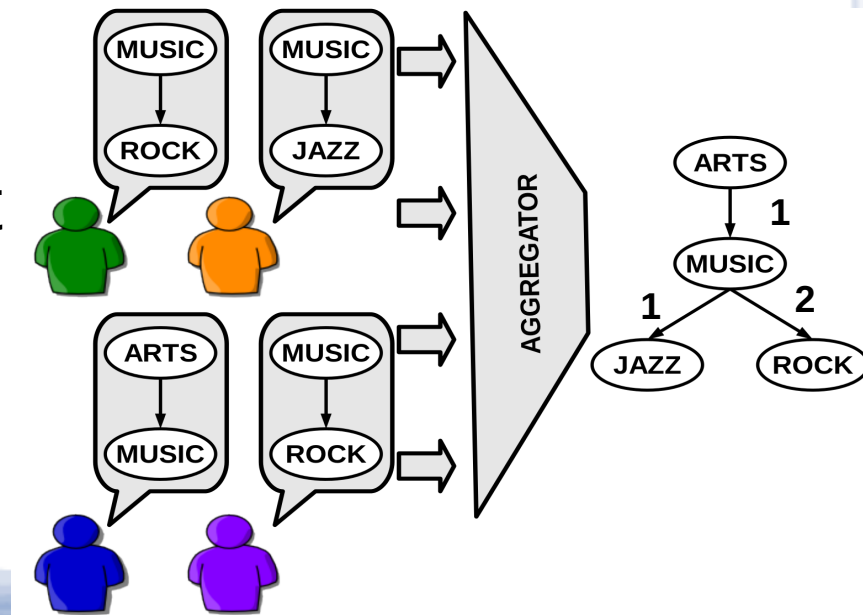
- Social Media & Tagging
 - Users organize their own data
 - Users fulfill the need of communication
- Crowdsourcing (Crowd + Outsourcing)
 - Massive collaboration
 - Collective Intelligence
- Taxonomies
 - Index and search web content
 - Need for experts (**maintenance overhead!**)

ROADMAP

- Introduction
- **Model & Theory**
- Algorithm
- Experimentation
- Conclusions

MODEL

- Extended Tagging: 'Vote' ($\text{tag}_a \rightarrow \text{tag}_d$)
- Vote Cardinality – Weights
- Challenges & Goals
 - Individuals are prone to errors
 - Incomplete (structural) information
 - Incremental process / Maintainability
- Vote Types
 - Ancestor \rightarrow Descendant
 - Descendant \rightarrow Ancestor
 - Crosslinks



THEORY

- Invariants
 - Integrity of Tree Properties
 - Maximum Votes Satisfiability
- Trivial Maximum Spanning Tree algorithms do not suffice

- **Problem Definition:**

INPUT : Complete graph $G = (V, E)$, weight $w(e) \in \mathbb{Z}^0$ for each $e \in E$

OUTPUT : A spanning tree T for G such that, if $W(u, v)$ denotes the sum of weights of the edges on the path joining u and v in T , then find B where:

$$B = \max \left(\sum_{u, v \in V} W(u, v) \right)$$

(B represents the maximum number of votes that is satisfied)

THEORY

- Invariants
 - Integrity of Tree Properties
 - Maximum Votes Satisfiability
- Trivial Maximum Spanning Tree algorithms do not suffice
- Problem Definition:

INPUT : Complete graph $G = (V, E)$, weight $w(e) \in \mathbb{Z}^0$ for each $e \in E$

OUTPUT : A spanning tree T for G such that, if $W(u, v)$ denotes the sum of weights of the edges on the path joining u and v in T , then find B where:

$$B = \max \left(\sum_{u, v \in V} W(u, v) \right)$$

(B represents the maximum number of votes that is satisfied)

NP-Hard !

Optimum Communication Spanning Tree (Garey & Johnson)

ROADMAP

- Introduction
- Model & Theory
- **Algorithm**
- Experimentation
- Conclusions

Vote Classification Routine

For all incoming votes (tag_anc→tag_des)

Node u ← hash(tag_anc)

Node v ← hash(tag_des)

if (only u exists)

 attach_new_child(u, v);

else if (none exists)

 create_new_tree(u, v);

else if (only v exists)

 merge(u, v);

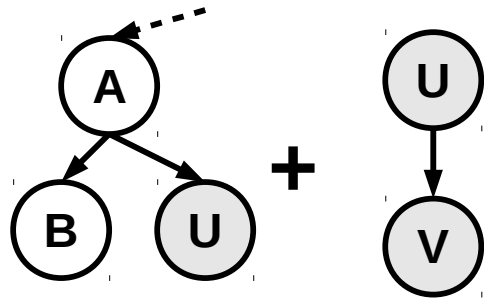
else (both u AND v exist)

 LeastCommonAncestor(u, v);

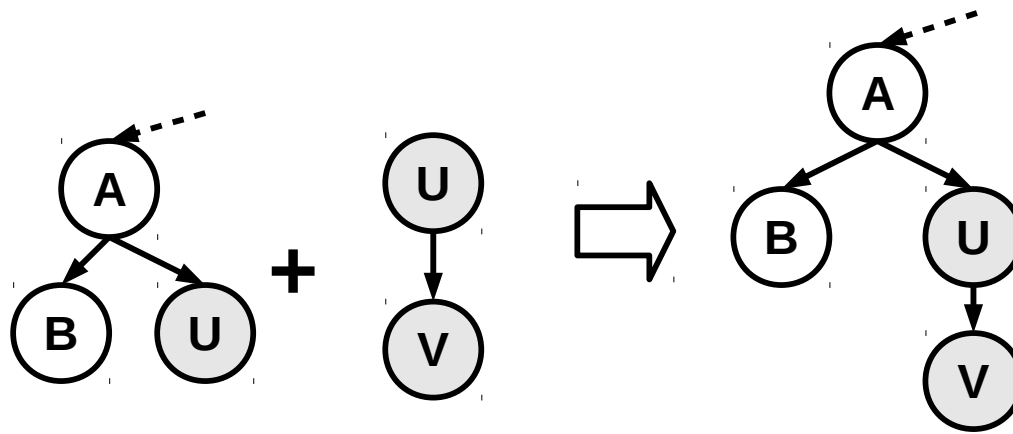
 expand_vertically(u, v);

 backedge_resolution(u, v);

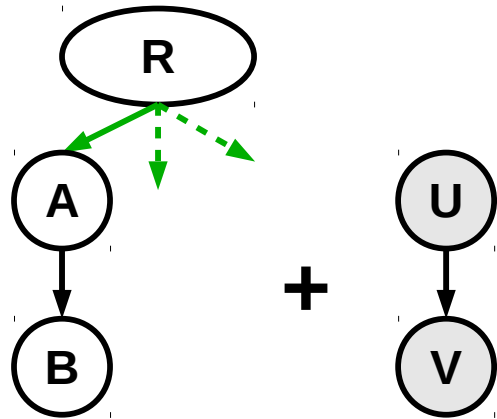
TRSFIRM: ATTACH NEW CHILD



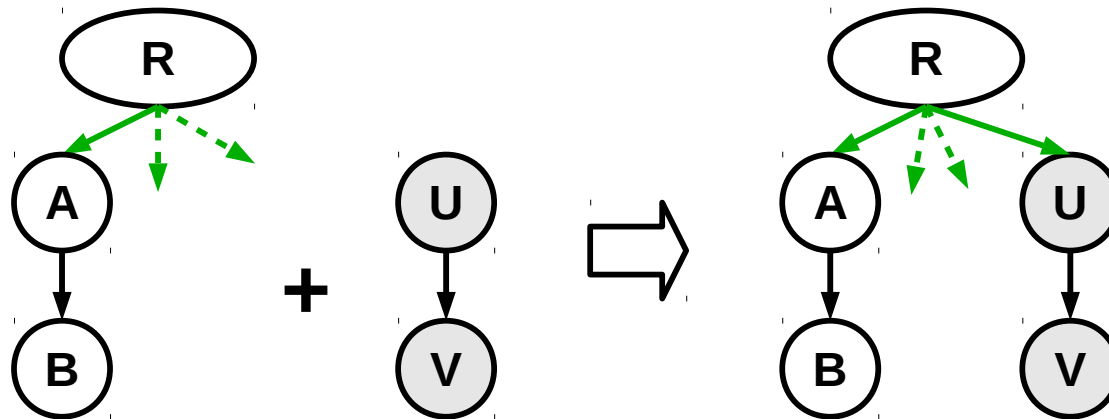
TRSFIRM: ATTACH NEW CHILD



TRSFMR: CREATE NEW TREE

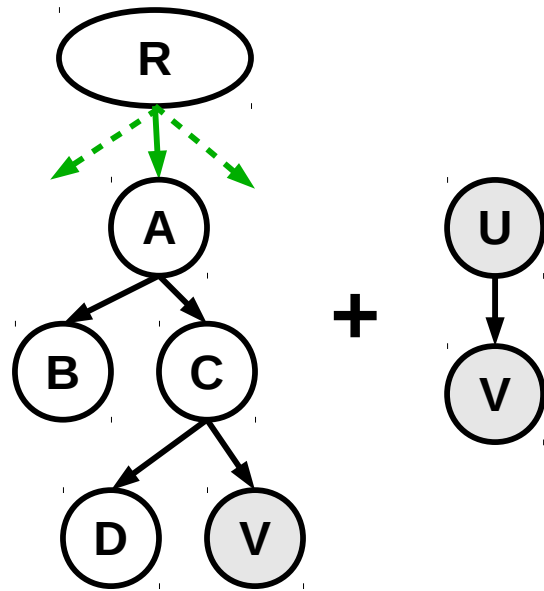


TRSFMR: CREATE NEW TREE

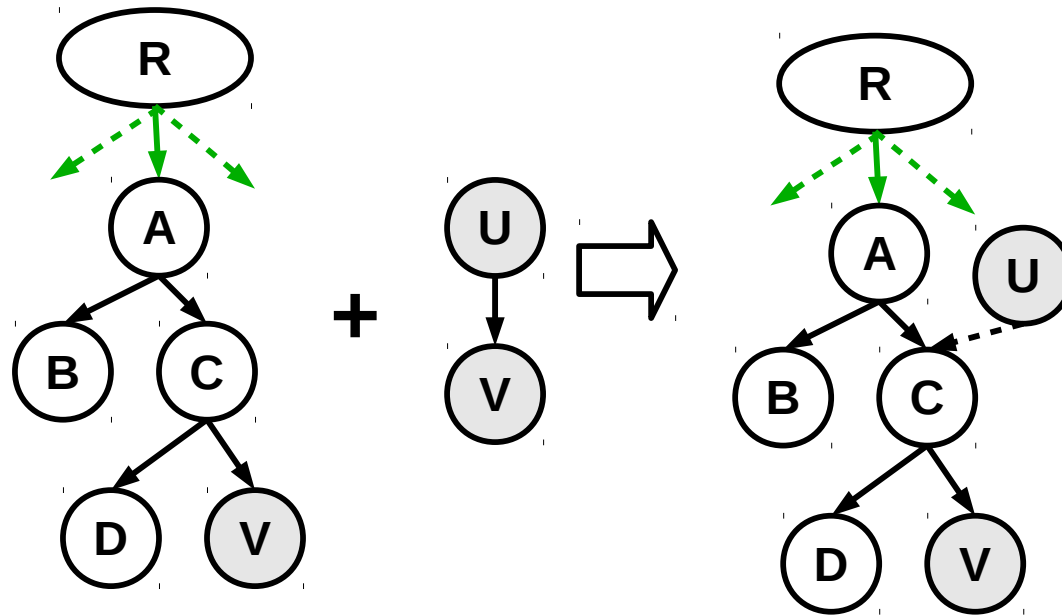


Shortcut: An artificial Parent→Child link that is not an explicit user generated relation, so its weight equals 0 and it is utilized in order to preserve structural continuation.

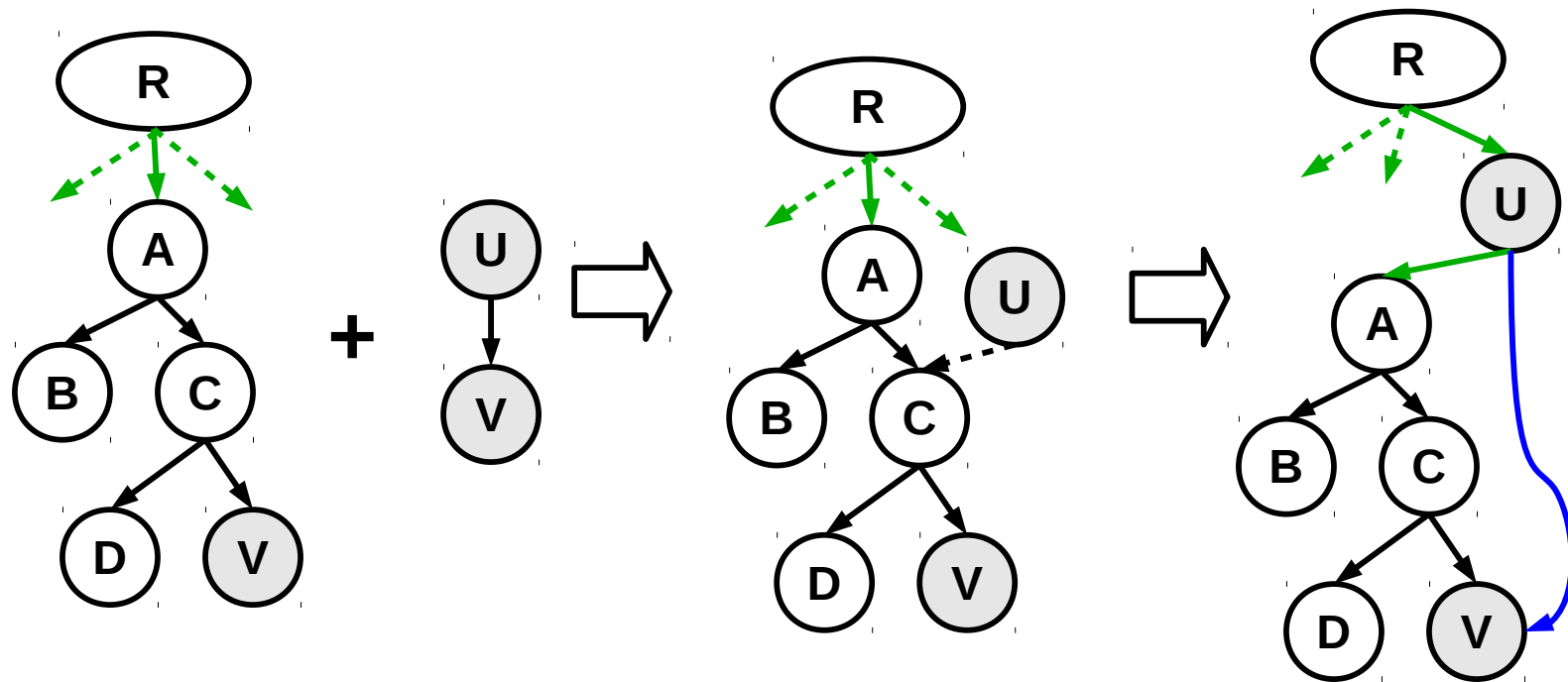
TRSFMR: MERGE



TRSFMR: MERGE

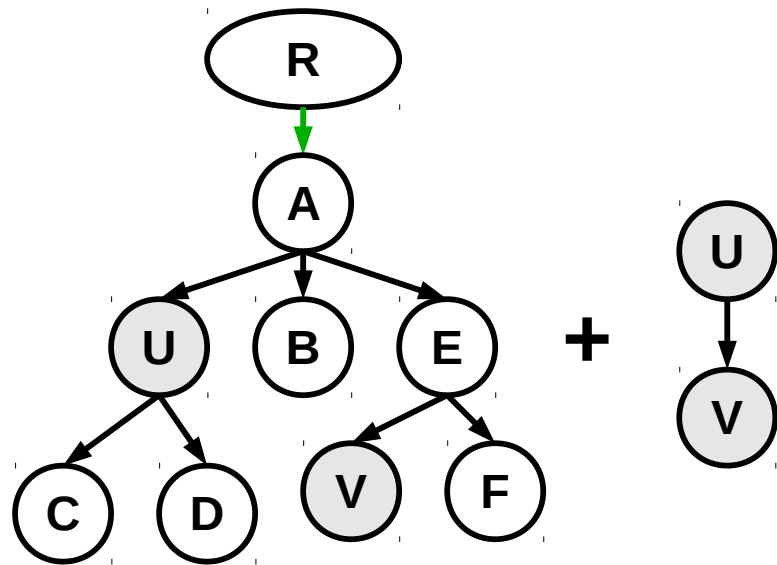


TRSFMR: MERGE

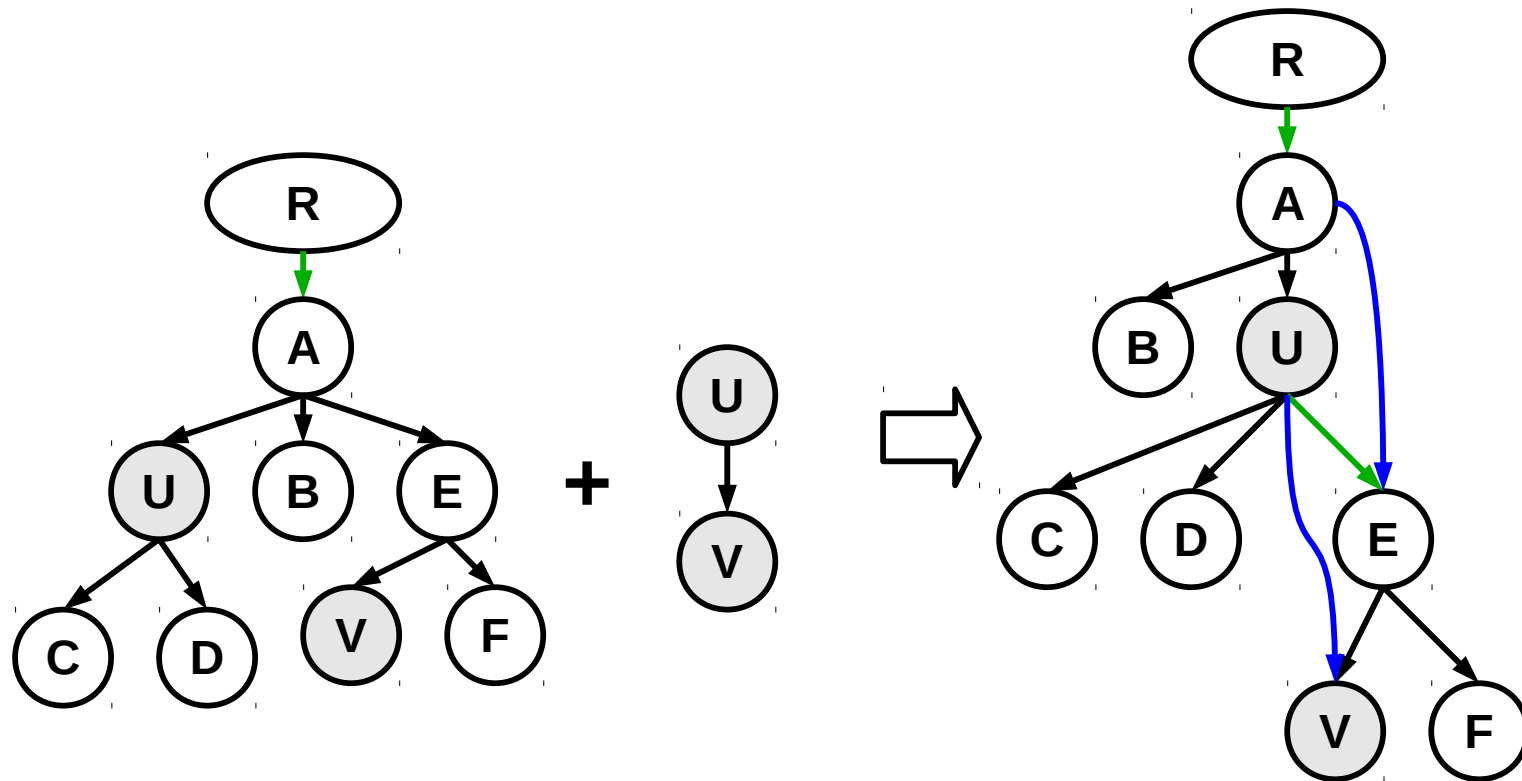


Forward Edge: A latent relation between two nodes. The source node is an ancestor in the taxonomy and the target is a descendant. Forward edges do not refer to Parent → Child links and remain hidden since they violate tree's properties.

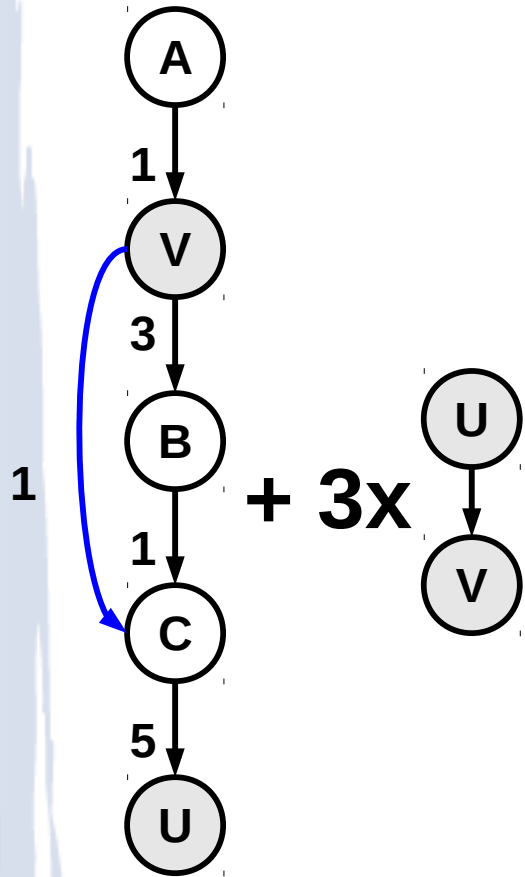
TRSFMR: EXPAND VERTICALLY



TRSFMR: EXPAND VERTICALLY

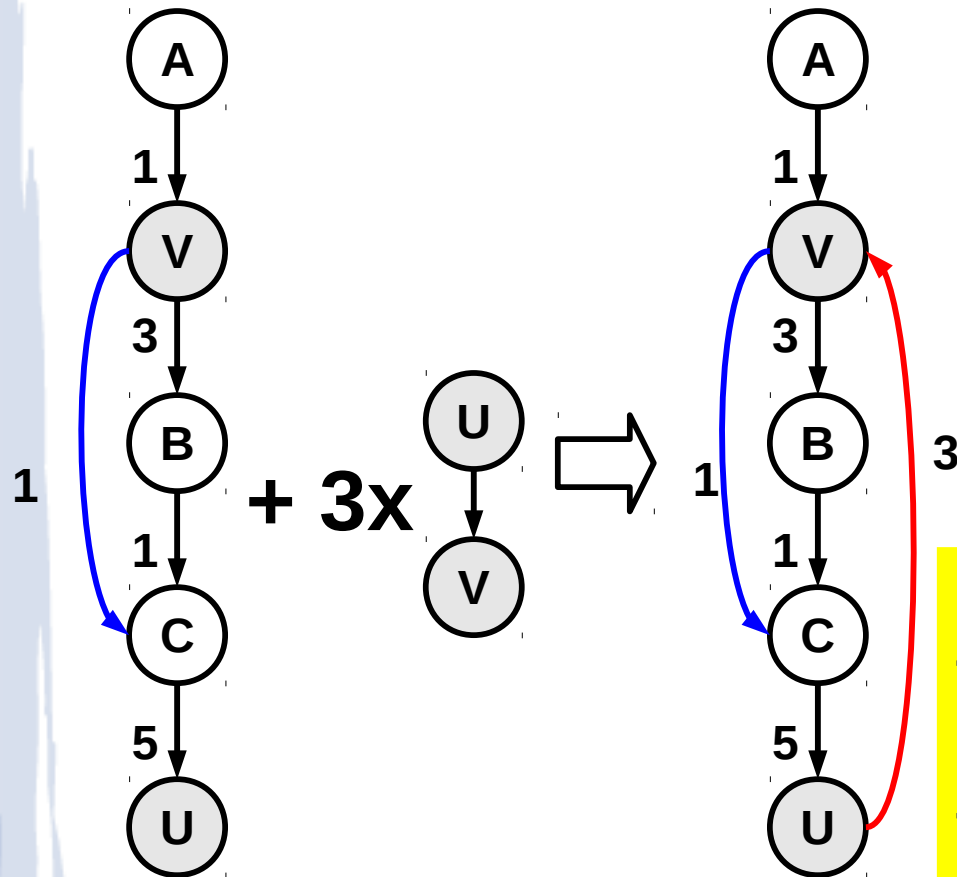


TRSFMR: BACKEDGE RESOLUTION



Initial State S

TRSFMR: BACKEDGE RESOLUTION



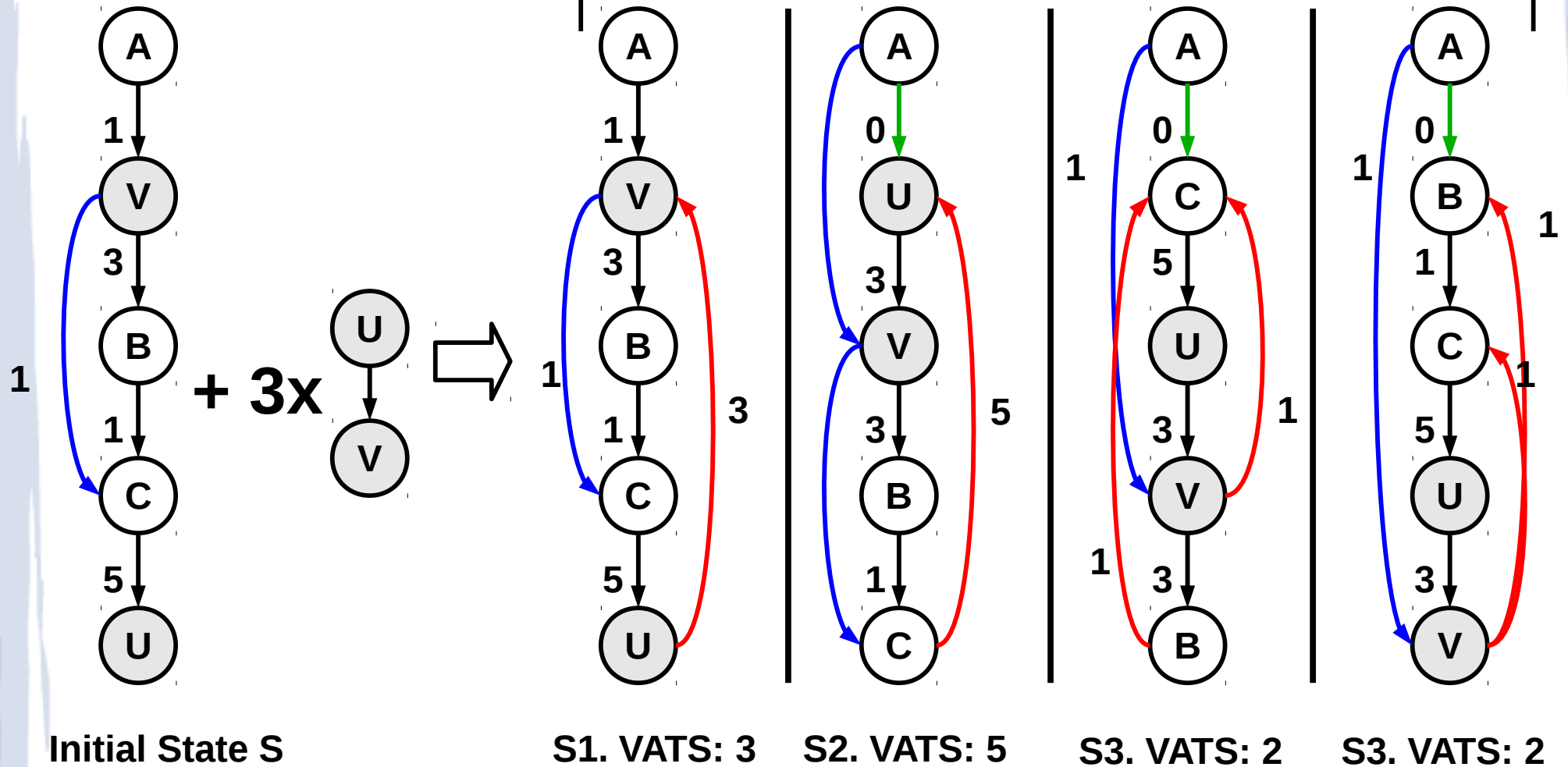
Initial State S

Backedge: A latent relation between two nodes. The source node appears as a descendant in the current taxonomy whereas the target as an ancestor. Backedges remain hidden since they violate the tree's properties.

TRSFMR: BACKEDGE RESOLUTION

VATS: Votes Against This State

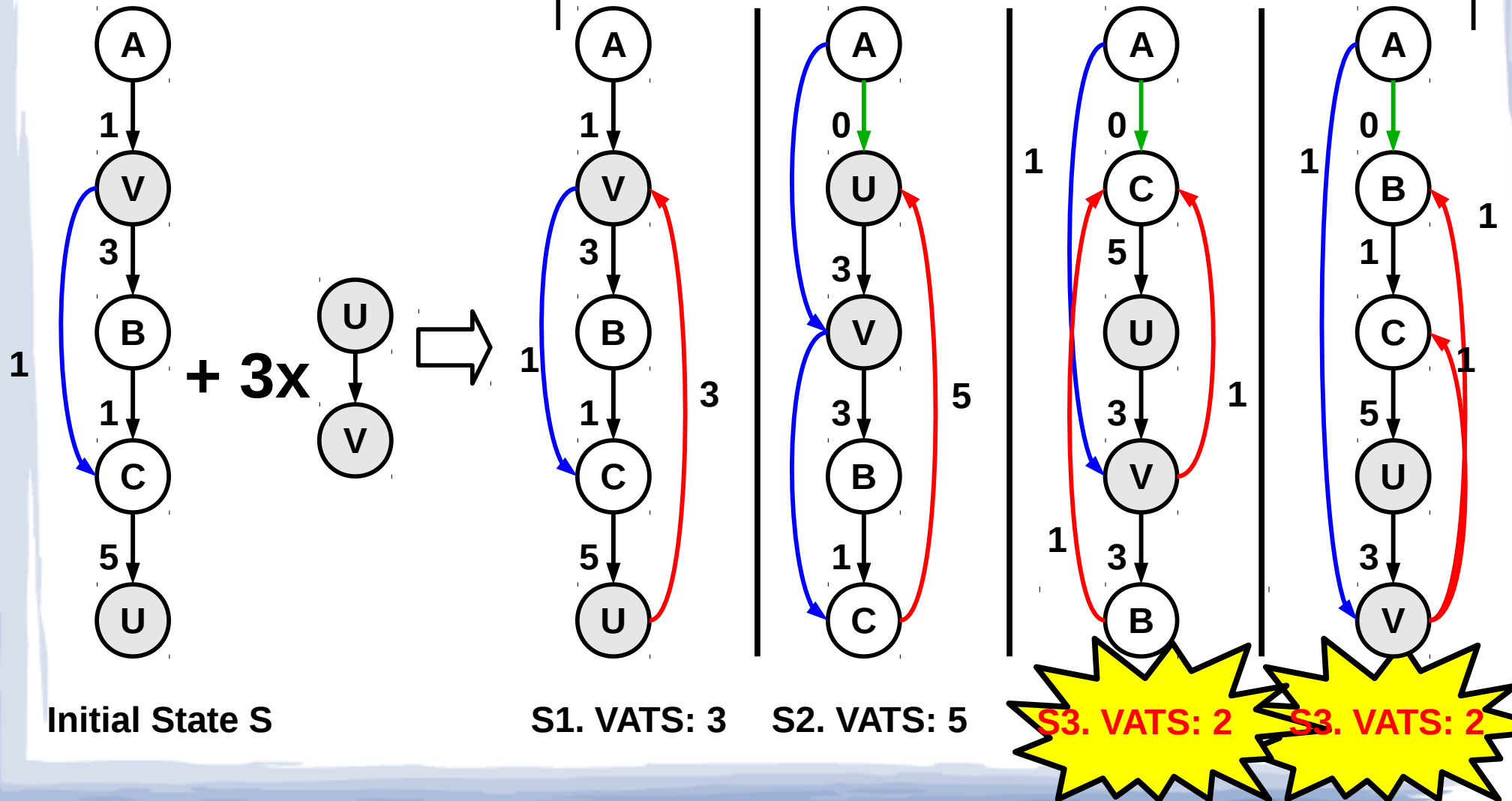
WHAT-IF ANALYSIS



TRSFMR: BACKEDGE RESOLUTION

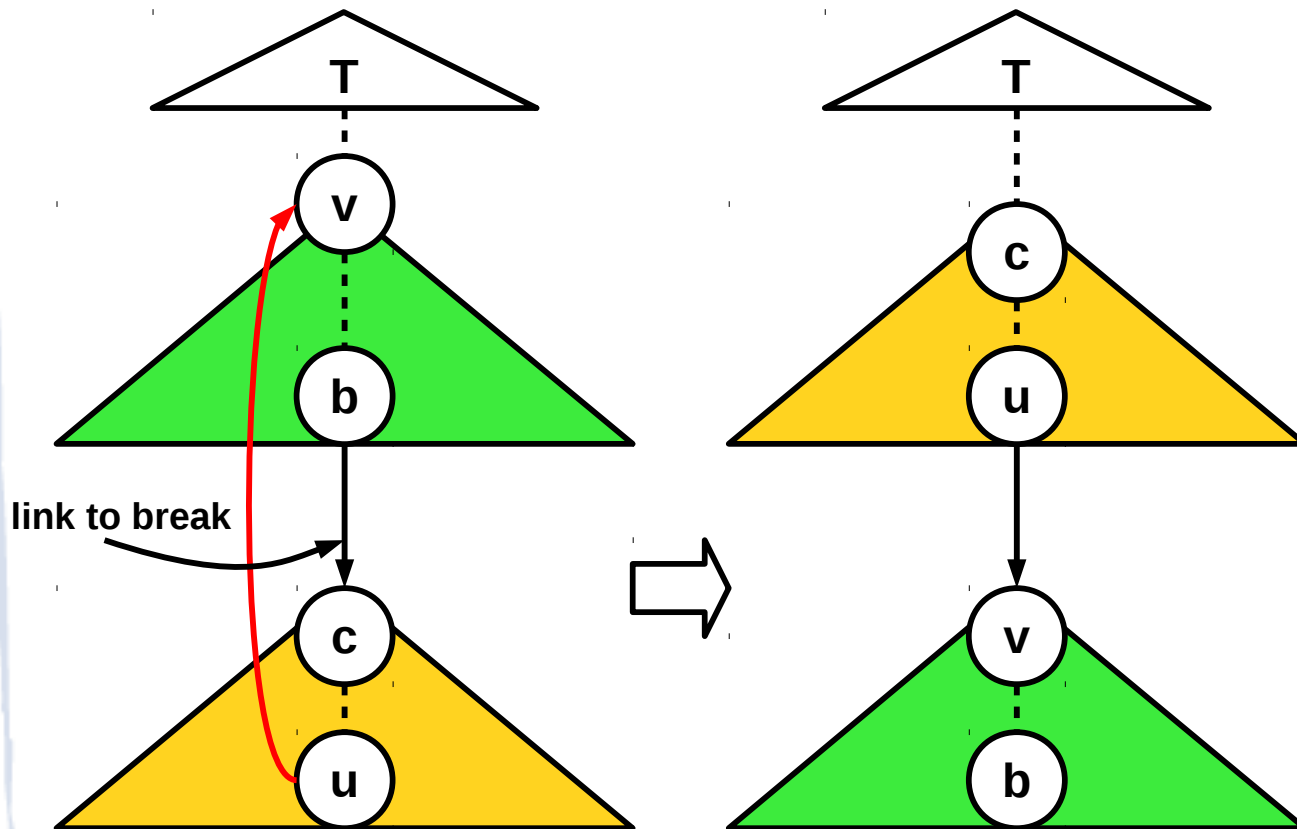
VATS: Votes Against This State

WHAT-IF ANALYSIS



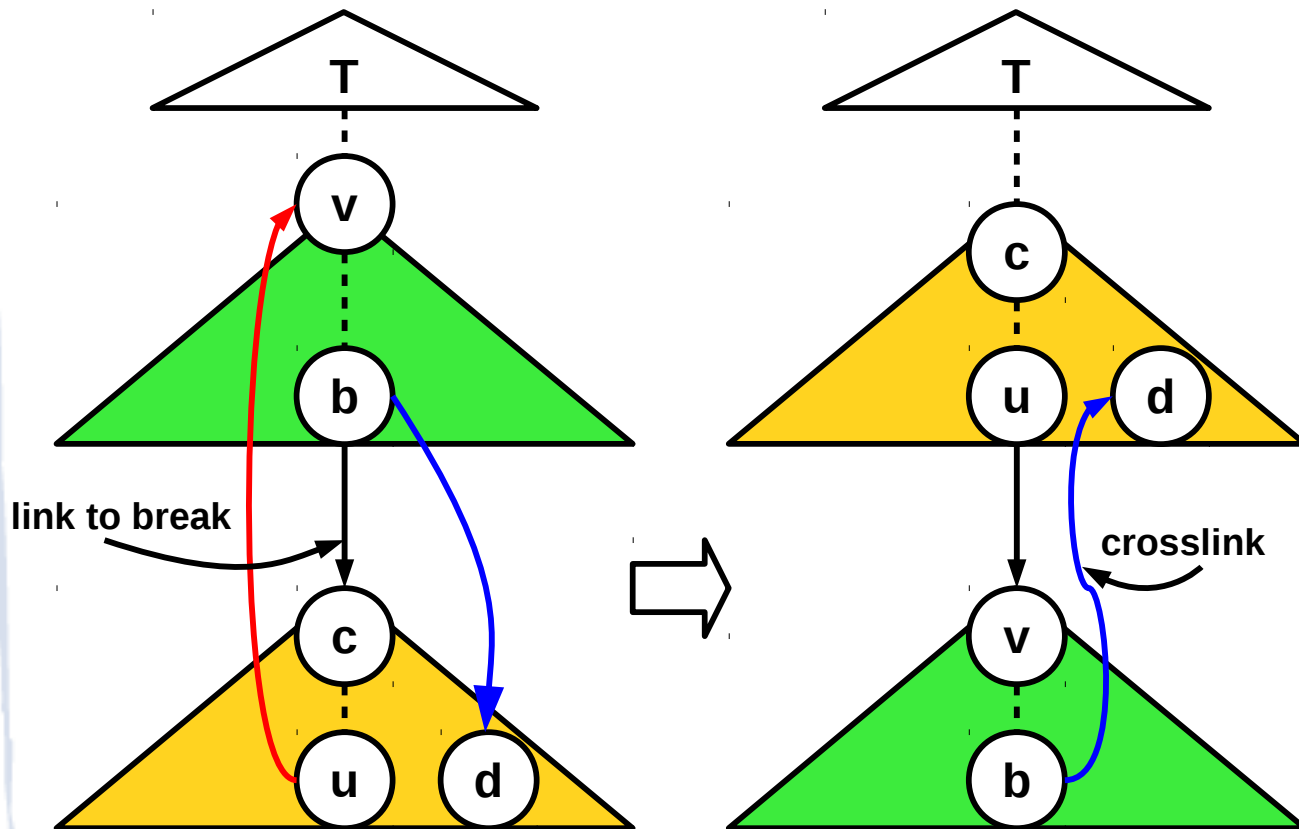
TRSFRM: BACKEDGE RESOLUTION

The Big Picture



TRSFMR: BACKEDGE RESOLUTION

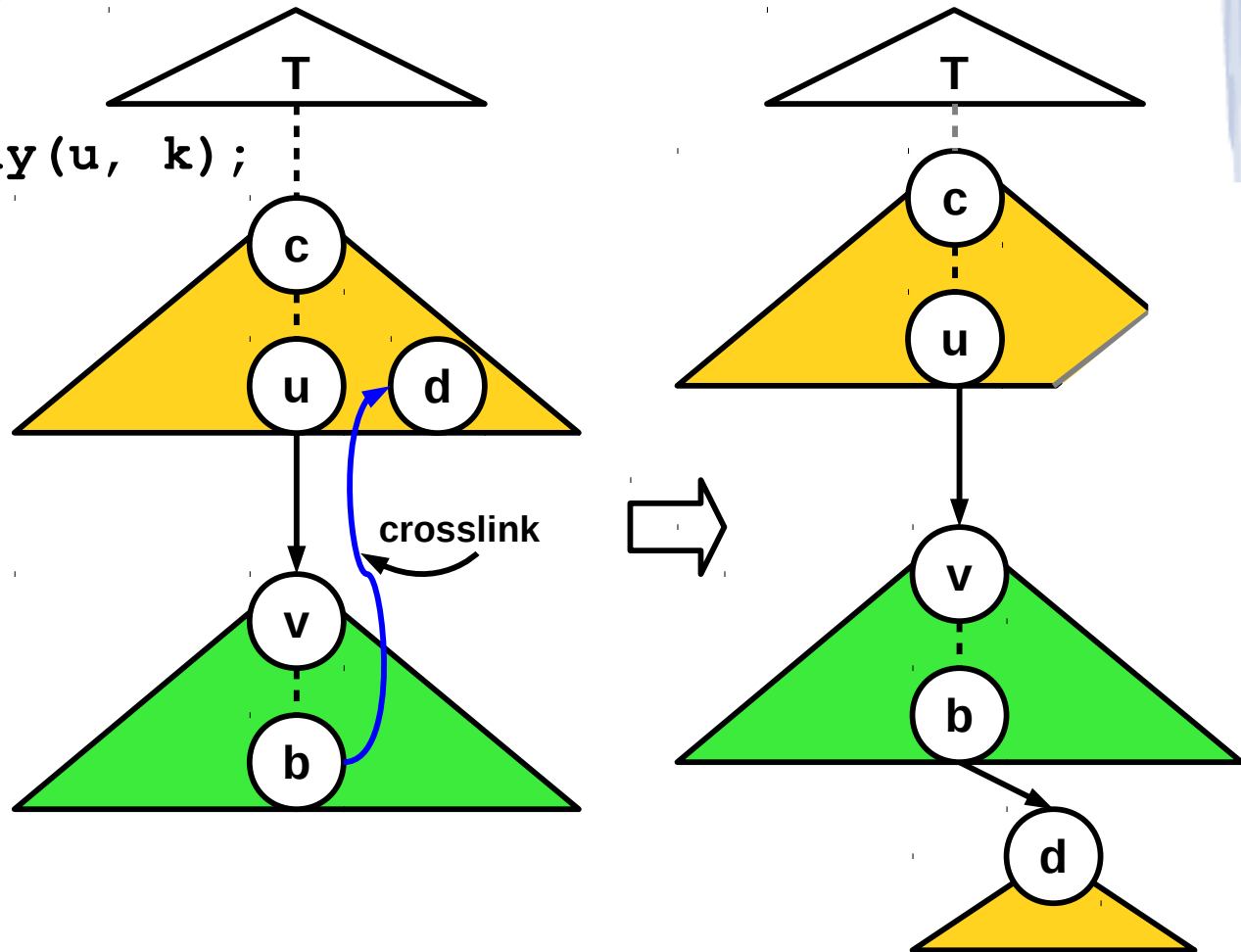
Follow ups



TRSFMR: BACKEDGE RESOLUTION

Follow ups

```
for all nodes u ∈ {v...b}
  for all u->fwd_edge(k)
    if (u == LCA(u, k))
      expand_vertically(u, k);
```



ROADMAP

- Introduction
- Model & Theory
- Algorithm
- **Experimentation**
- Conclusions

DATASET & METRICS

Dataset: ACM Computing Classification System

1473 Nodes

3824 (Ancestor → Descendant) Relations

Metrics:

To → Original Tree

T → Tree derived from Crowd

Ro → Set of Parent → Child relations in To

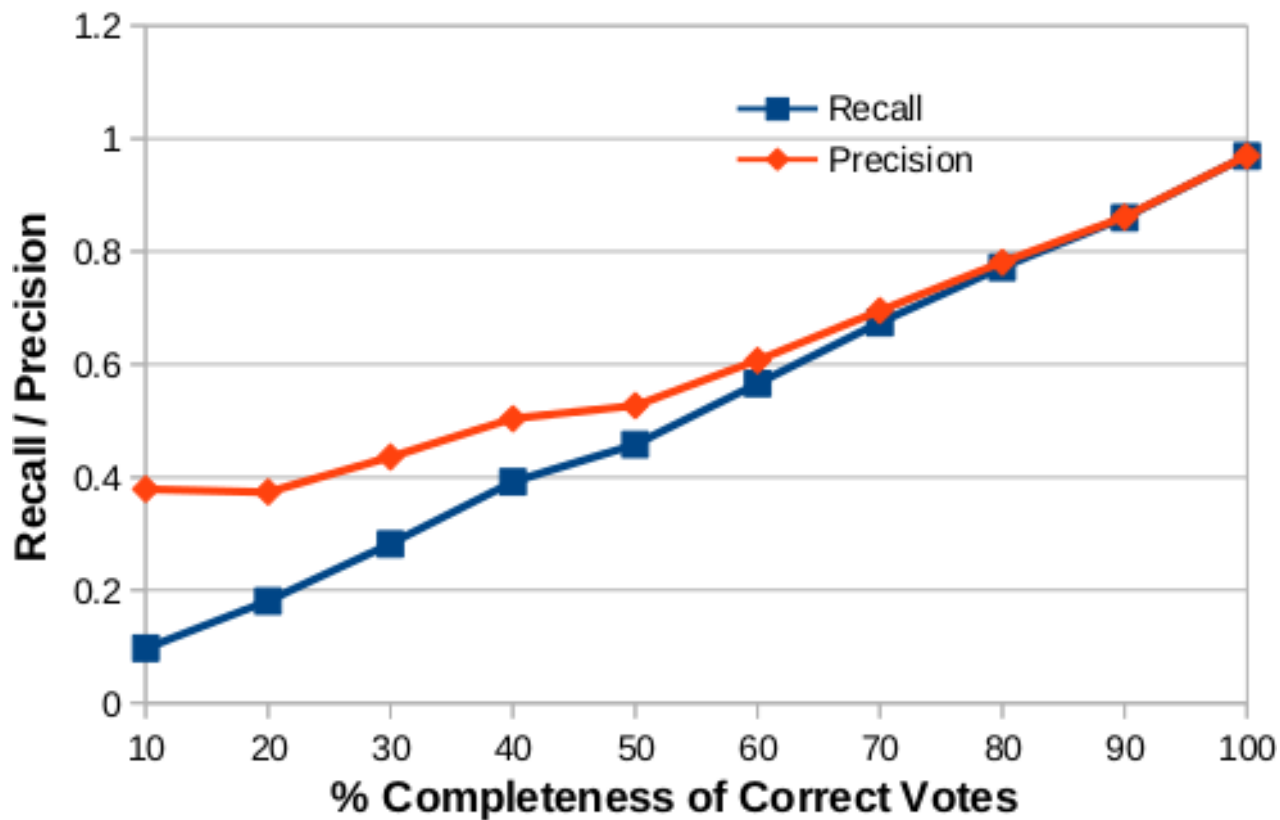
R → Set of Parent → Child relations in T

$$Recall = \frac{|R \cap Ro|}{|Ro|} \quad Precision = \frac{|R \cap Ro|}{|R|}$$

$$FScore = \frac{2 \times Recall \times Precision}{Recall + Precision}$$

SYNTHETIC EXPERIMENT - 1

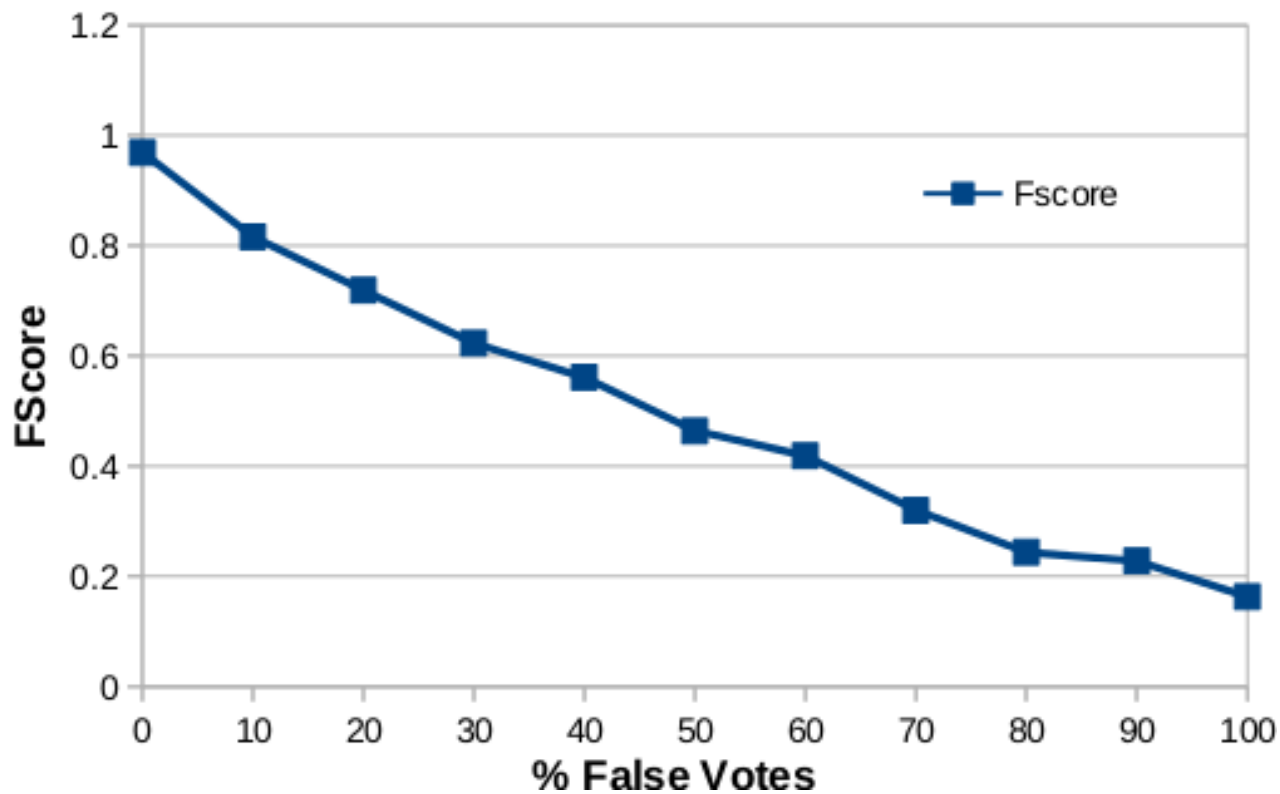
Setup: On every iteration reduce by 10% the number of Correct Votes



SYNTHETIC EXPERIMENT - 2

Setup: 100% Correct Votes

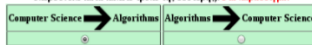
On every iteration increase by 10% the number of False Votes



CROWDSOURCING EXPERIMENT

- 245 Nodes
- 620 Ancestor → Descendant Relations (+ inverted counterparts)
- Preprocessing: Remove misleading concepts ('Misc', 'General')
- HIT: 25 pairs of relations presented in groups of 5
- 3-day Period

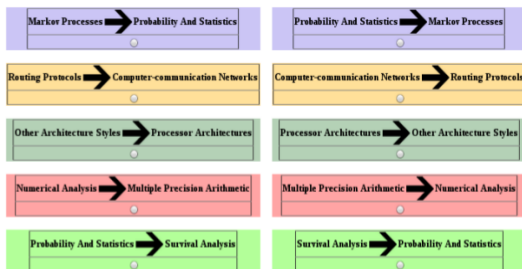
Ζητούμενο του πειράματος η σήμανση της σωστής σχέσης μεταξύ δύο εννοιών. Προσβλέπουμε να επιλέξετε κάθε φορά αυτή στην οποία το ενδιαφέρον θέλετε έχει φασά από τη γενικότερη προς την ειδικότερη ή πρώτα είναι υπαρκτό και δεύτερο τρόπο της δεύτερης. Για παράδειγμα:



Εδώ φαίνεται πως η έννοια Computer Science είναι γενικότερη της έννοιας Algorithms.

Αρκεί να τις σχέσεις είναι διαφορετικές. Χρησιμοποιήστε τη λογική σας ή το κοινό για να δείτε ποια είναι σωστή. Μπορείτε επίσης να παραλείψετε το συγκεκριμένο ζεύγος. Για πρόβλημα στείλτε mail στο karampini@cc.upatras.gr εντυπωσιακά το # με @!

Set: 1 of 5



OK

EXIT

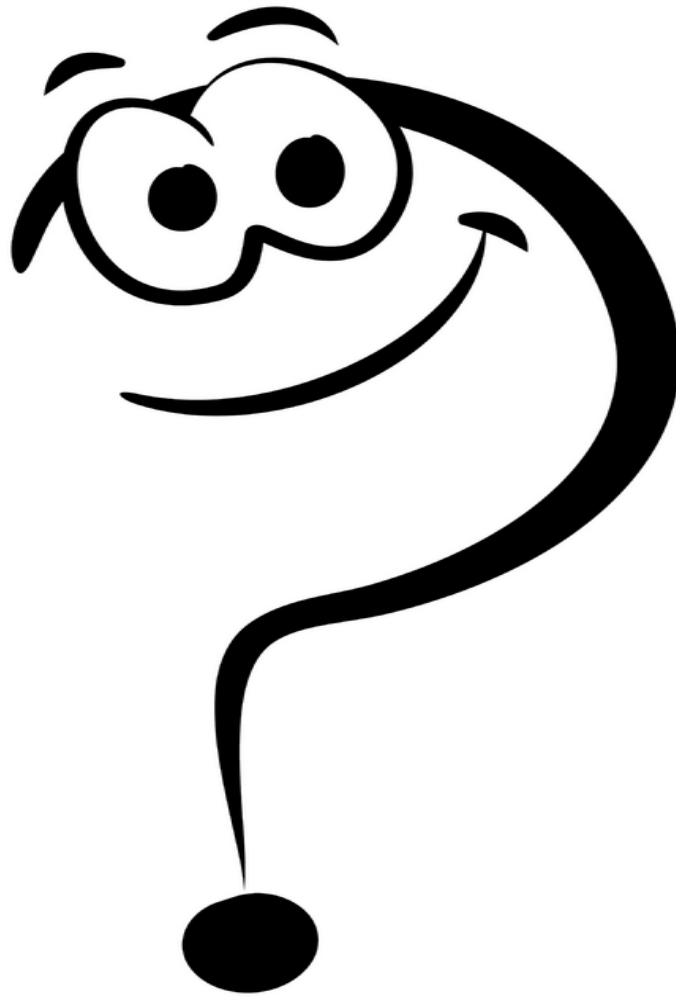
Users	102
Total Votes	2155
Correct Votes	1501
False Votes	435
Undefined	229
Rel. without Votes	24
Rel. without Correct Votes	33
Rel. without False Votes	302
Max Votes on a Rel.	10
FScore	0.486
Predicted FScore	0.519

ROADMAP

- Introduction
- Model & Theory
- Algorithm
- Experimentation
- **Conclusions**

CONCLUSIONS

- Novel idea for a Taxonomy development in Crowdsourcing environments
- Challenges & Model
- Problem definition (NP-Hardness)
- Algorithm
 - Online
 - Copes with lack of input & user conflicts
- Synthetic & Real world crowdsourcing experiments



Thank you crowd...