# Requirements Engineering

## (Summer 2021)

## Prof. Nan Niu (nan.niu@uc.edu)

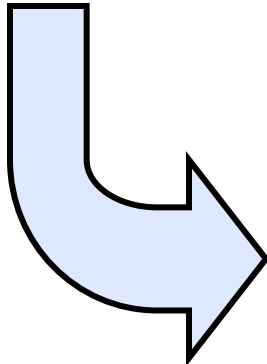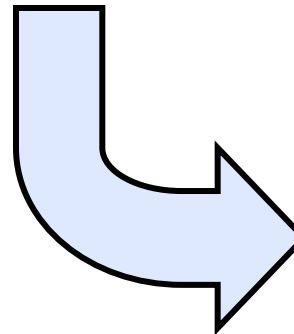## https://github.com/nanniu/RE-Summer2021

# Today's Menu
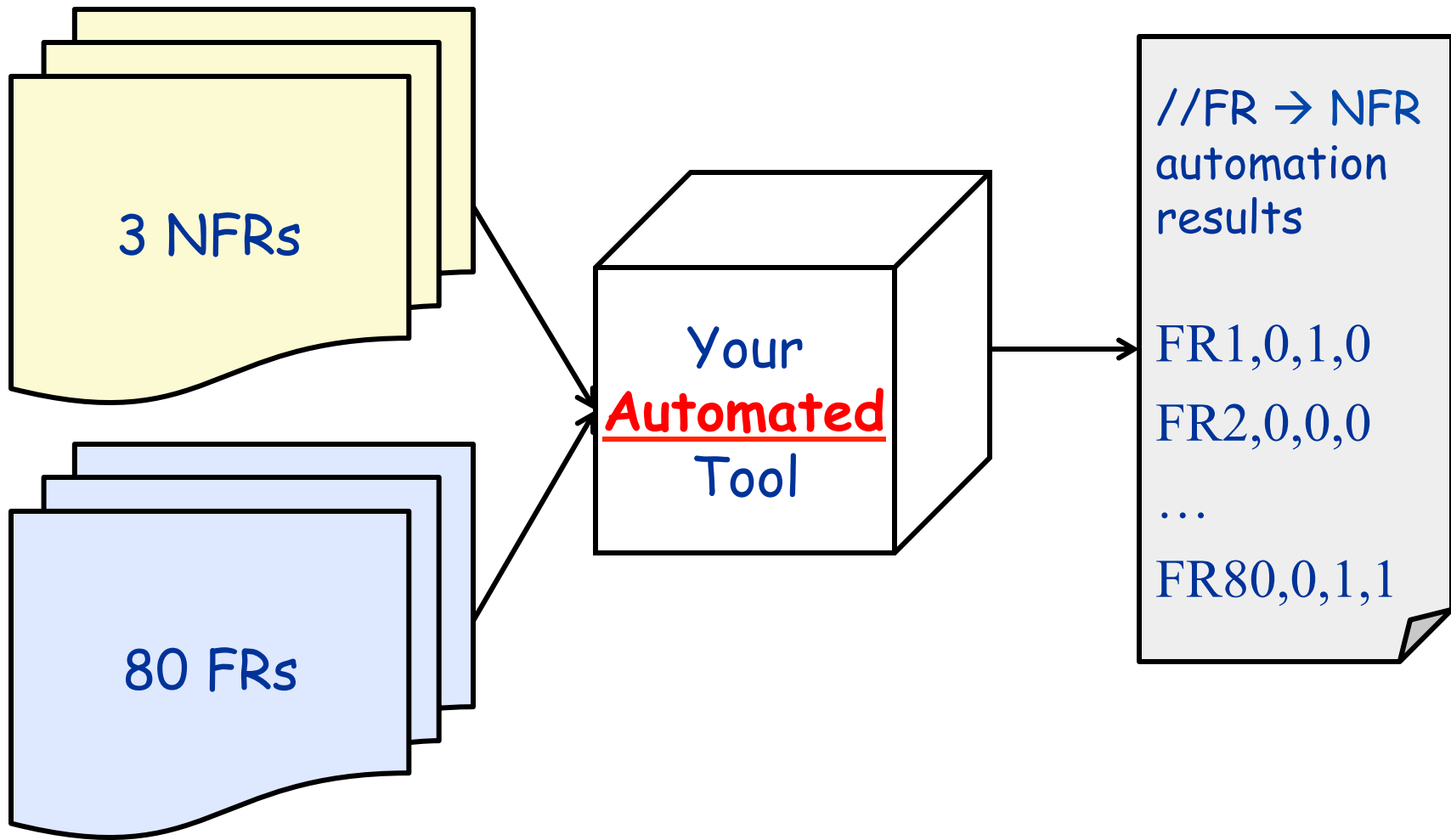
**Thursday (July 22)**
Req.s Traceability
ASN4 Release

**Friday (July 23):**
Unsupervised Learning
(ASN4 Q&A)

**Monday (July 26):**
RE Story
Presentations (ASN2)

# ASN4: A conceptual picture

3 NFRs

80 FRs

Your
**Automated**
Tool

//FR → NFR
automation
results

FR1,0,1,0

FR2,0,0,0

…

FR80,0,1,1

# ML: Will it work for ASN4?



Classical Machine Learning

Task Driven

Data Driven

Supervised Learning
( Pre Categorized Data )
Predications & Predictive Models

Unsupervised Learning
( Unlabelled Data )
Pattern/ Structure Recognition

Classification
( Divide the socks by Color )
Eg. Identity Fraud Detection

Regression
( Divide the Ties by Length )
Eg. Market Forecasting

Clustering
( Divide by Similarity )
Eg. Targeted Marketing

Association
( Identify Sequences )
Eg. Customer Recommendation

*Source:* *Google Images*

# Supervised Learning: *not really*

⇨ Run #1: 80 FRs and 3 NFRs

⇨ Run #2: 100 FRs and 3 NFRs (i.e., 20 new/unseen FRs compared to Run #1)

⇨ Run #3: 100 FRs and 4 NFRs (i.e., 1 new/unseen NFR compared to Run #2)

```
1    FR1,0,1,0
2    FR2,0,0,0
3    FR3,0,1,0
4    FR4,0,1,0
5    FR5,0,1,0
6    FR6,0,1,0
7    FR7,0,0,0
     ...
76   FR76,0,1,0
77   FR77,0,1,0
78   FR78,0,1,0
79   FR79,0,1,0
80   FR80,0,1,0
```
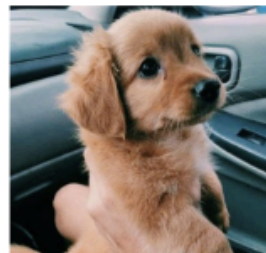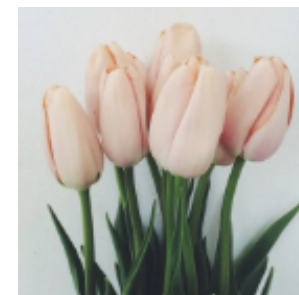
(unseen data)
   FR81: "description of a new FR"
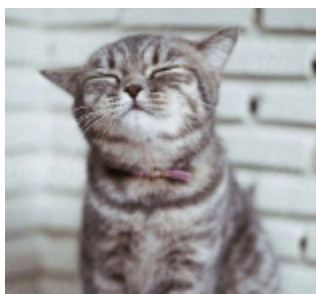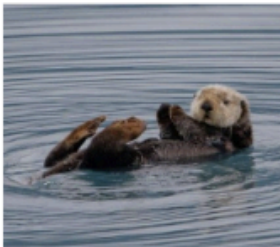
# What to learn?

# Today's Take-Away
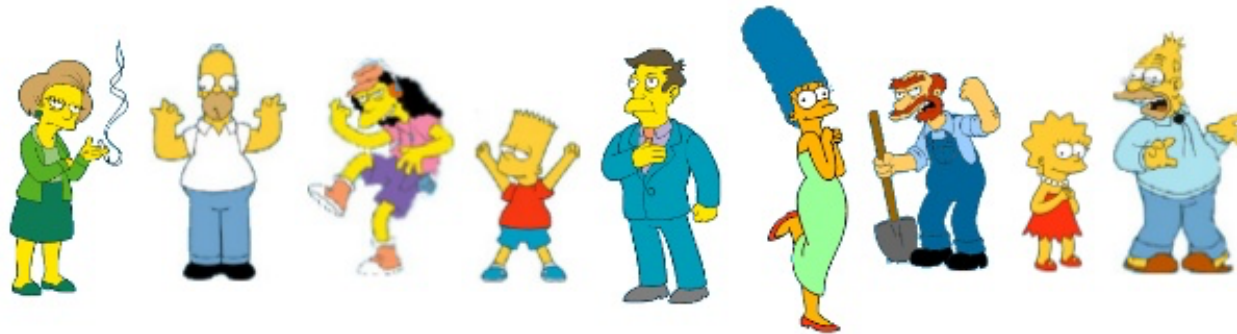
→Unsupervised learning can be used to solve ASN4.

**Unsupervised learning** is a type of machine **learning** algorithm used to draw inferences from datasets consisting of input data without labeled responses. The most common **unsupervised learning** method is cluster analysis, which is used for exploratory data analysis to find hidden patterns or grouping in data.

*Source: Google Search*

# Clustering: Finding ~~the~~ **a** natural grouping of data



What is a natural grouping among these objects?

Clustering is subjective
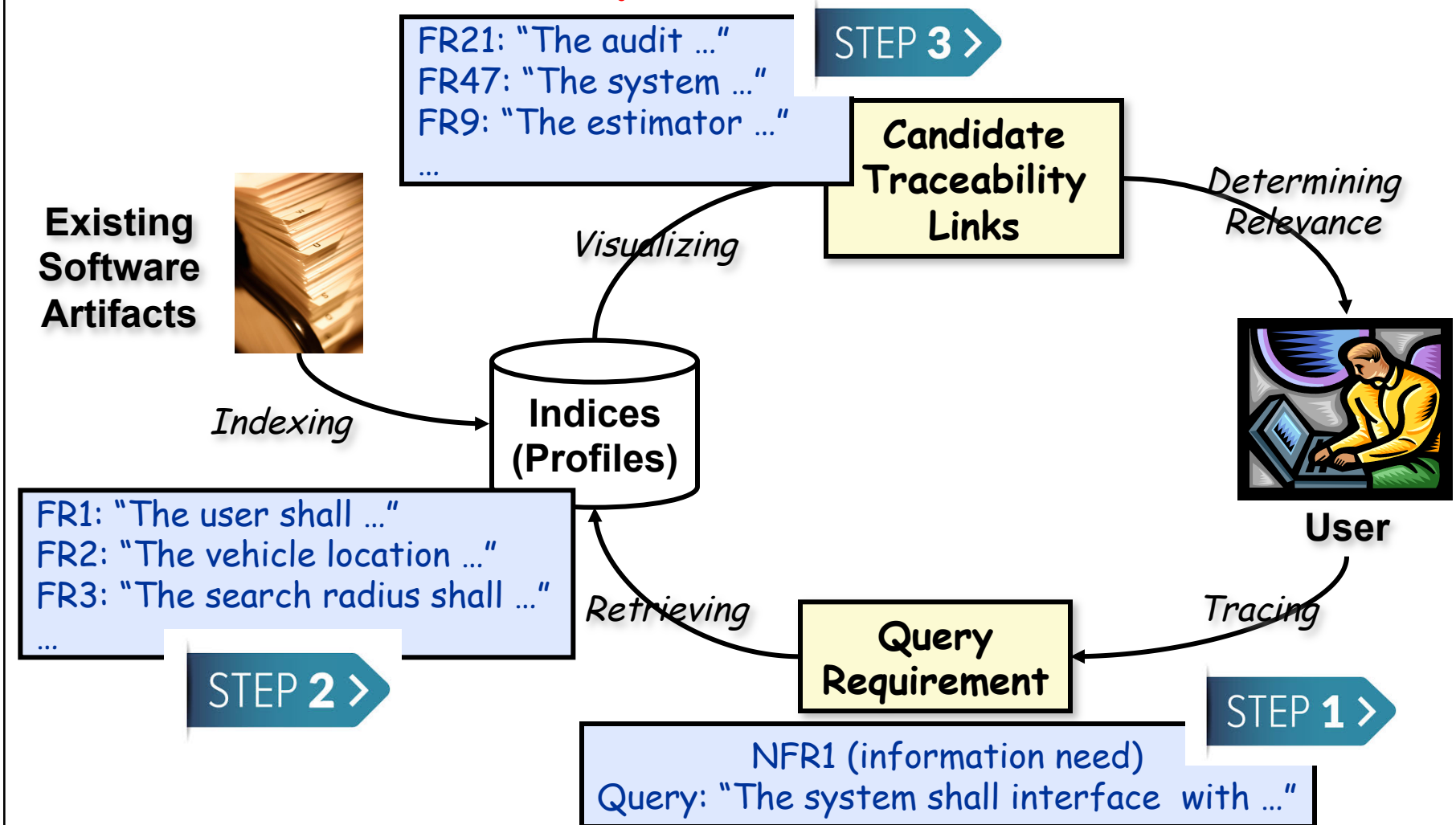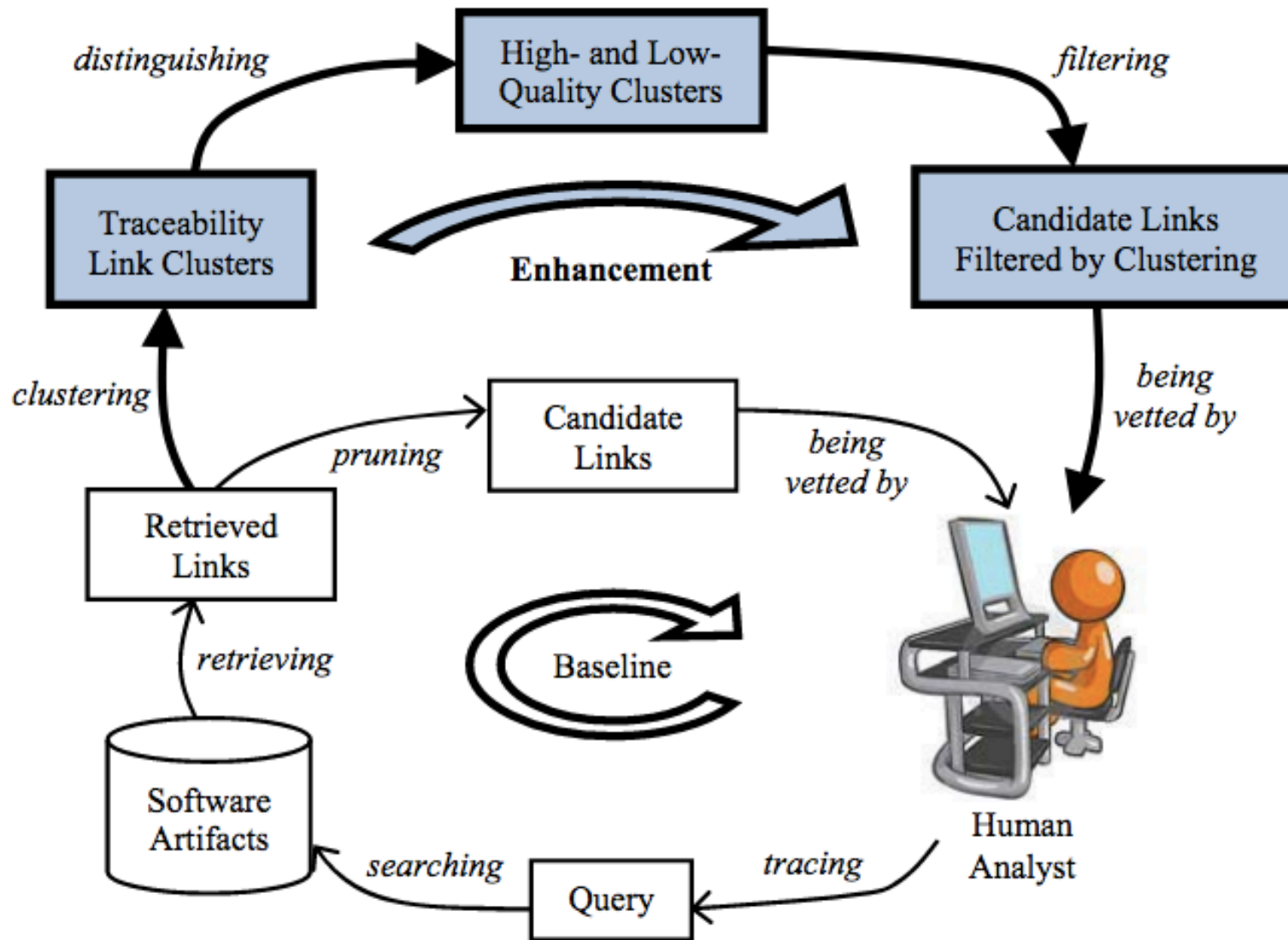
Simpson's Family    School Employees        Females        Males

# IR-Based ASN4 Solution
# (fully automatic)

FR21: "The audit ..."
FR47: "The system ..."
FR9: "The estimator ..."
...

**STEP 3 >**

**Candidate Traceability Links**

*Determining Relevance*

**Existing Software Artifacts**

*Visualizing*

*Indexing*

**Indices (Profiles)**

FR1: "The user shall ..."
FR2: "The vehicle location ..."
FR3: "The search radius shall ..."
...

**STEP 2 >**

**User**

*Retrieving*

*Tracing*

**Query Requirement**

**STEP 1 >**

NFR1 (information need)
Query: "The system shall interface with ..."

# Clustering as Enhancement

# Cluster Hypothesis

→**In IR, *cluster hypothesis* suggests, "relevant documents tend to cluster near other relevant documents and farther away from irrelevant ones".**

　　↳Applying *cluster hypothesis* in automated requirements tracing, it suggests, "correct links tend to be more similar to each other than to incorrect links".

　　　　*Source: https://homepages.uc.edu/~niunn/papers/RE12.pdf*　　　11

# Key questions answered

**Does cluster hypothesis hold in automated requirements traceability, and if so, how to best exploit it?**
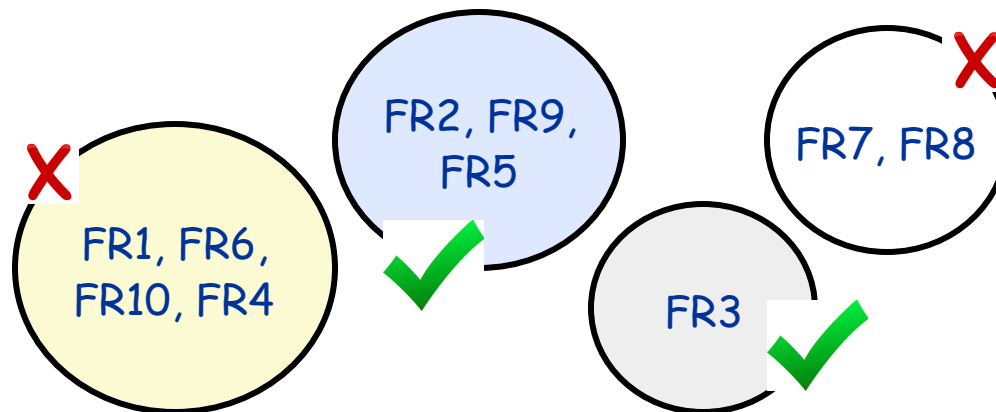
> ## Enhancing Candidate Link Generation for Requirements Tracing: The Cluster Hypothesis Revisited
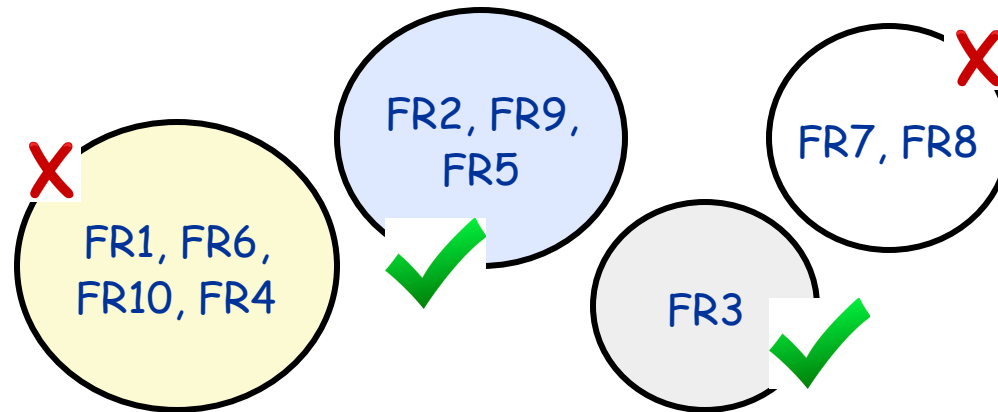>
> Nan Niu     Anas Mahmoud

# Conceptually

**For a particular trace query (e.g., an NFR), identify candidate traceability links (e.g., via Jaccard index). Then, cluster the candidate traceability links:**



**each cluster (as opposed to each link / FR) will be judged "correct / traceable" or "incorrect / not traceable".**

**Today's poll question:** *"In addition to which clustering algorithm to use, <u>how many</u> other decisions are required?"*

# Conceptually, to achieve _full automation_



FR2, FR9, FR5 ✓

FR7, FR8 ✗

FR1, FR6, FR10, FR4 ✗

FR3 ✓

**(0) Which clustering algorithm to use?**

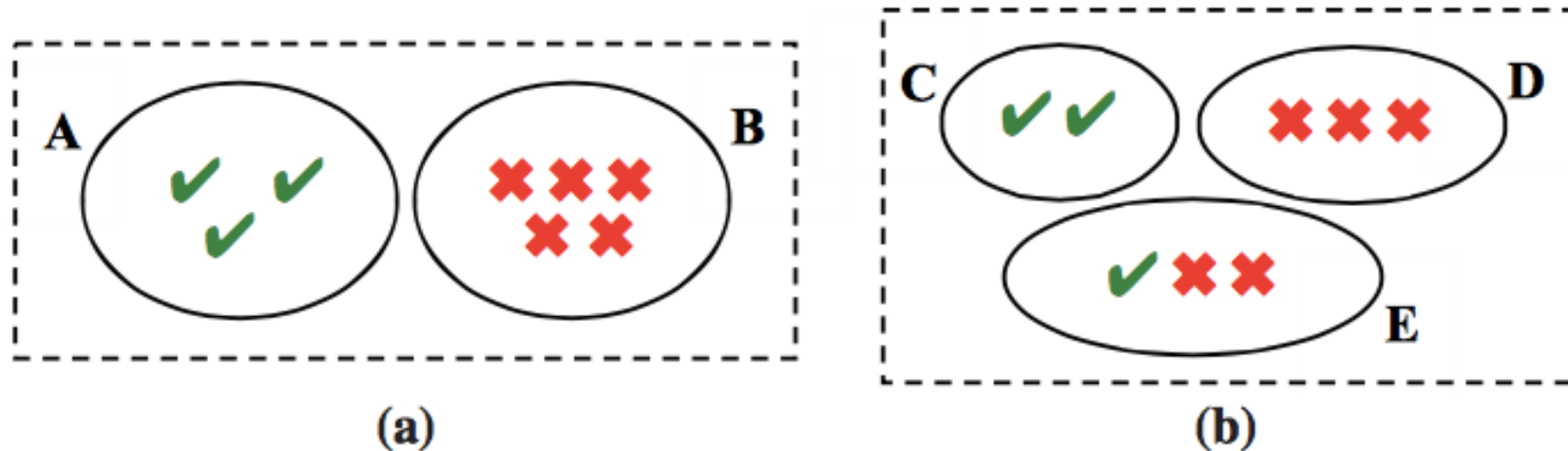**(1) How many clusters ($k$) to produce?**

**(2) How to rank the resulting clusters?**

**(3) After ranking, how many "good" clusters to keep and how many "bad" clusters to remove?**

# Results

- The cluster hypothesis holds in traceability.

- Single-link (SL), at the $k=8$ clustering granularity, represents a good candidate mechanism for fulfilling the potential suggested by the cluster hypothesis.

- The quality of clusters can be adequately inferred by their maximum similarity (MAX) to the trace query, and the 3 lowest-quality clusters contain such a high density of false positives that discarding them significantly improves the overall quality of the candidate link generation.

*Source: https://homepages.uc.edu/~niunn/papers/RE12.pdf*

# Evaluating Clustering Results
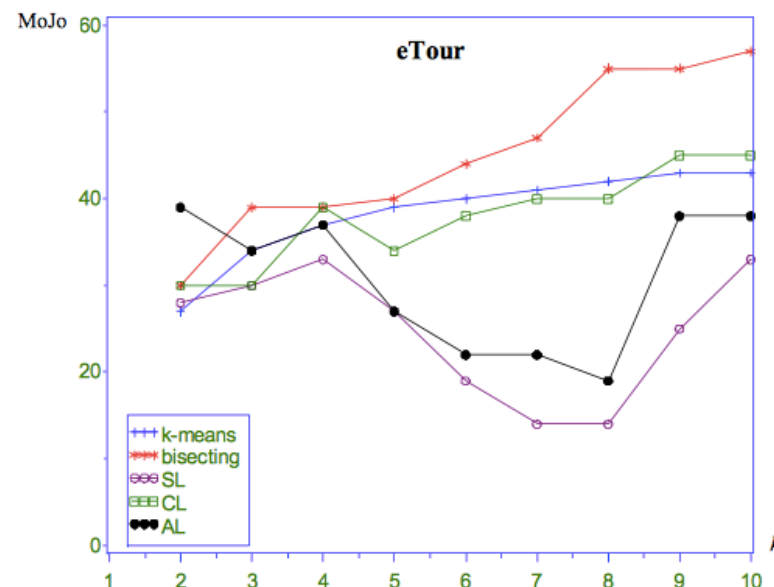


(a)                     (b)

**Suppose (a) is the "answer set", then the quality of (b) can be measured by the MoJo distance (i.e., the number of <u>Mo</u>ves and <u>Jo</u>ins to transform (b) to (a)); here MoJo distance = 2 (i.e., move the correct link from E to C, and then join D and the revised E together).**
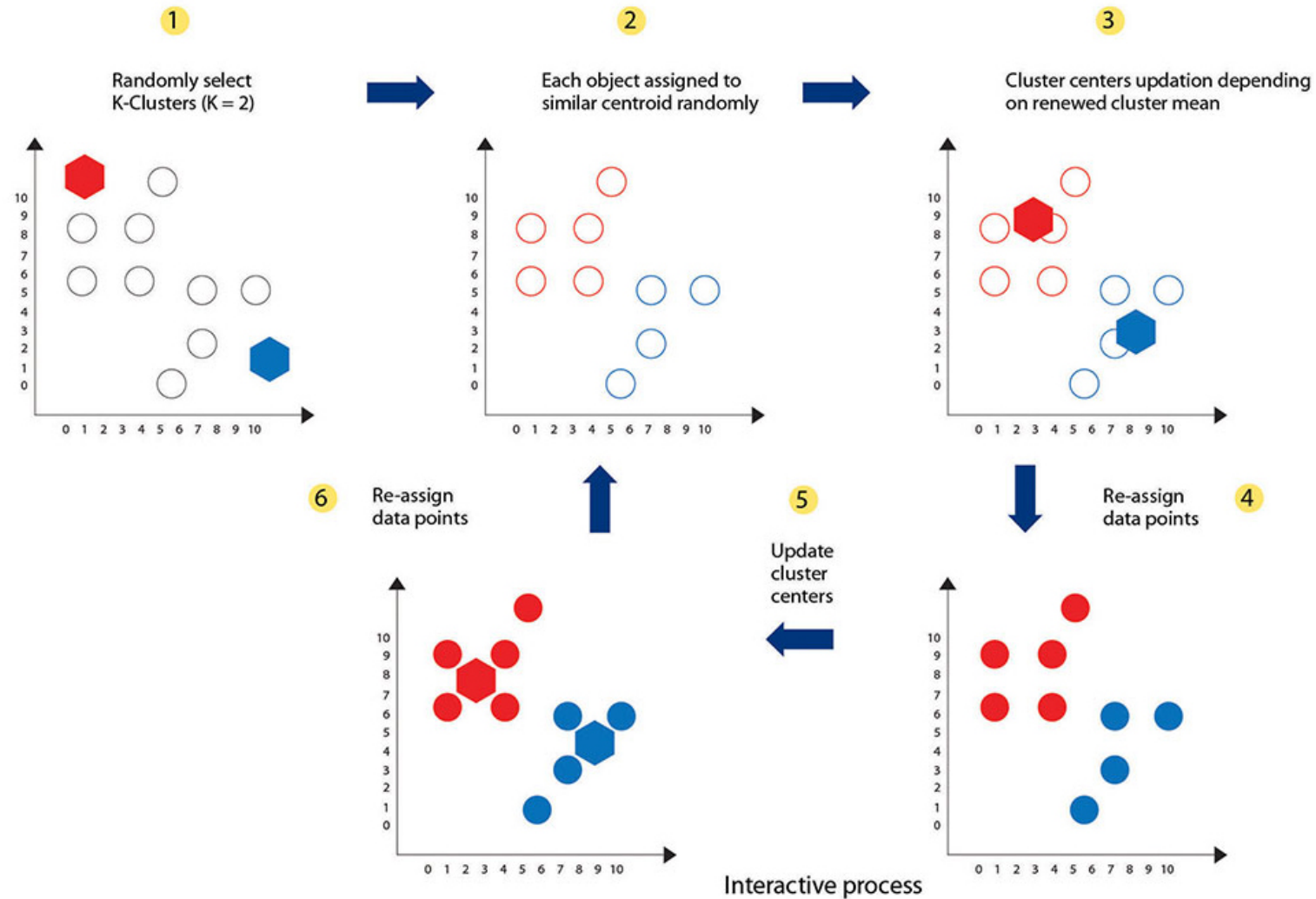
**The smaller the MoJo distance, the better the clustering.**

# Clustering Algorithms & *k*

→ *k*-means: centroid-based

→ bi-secting: top-down

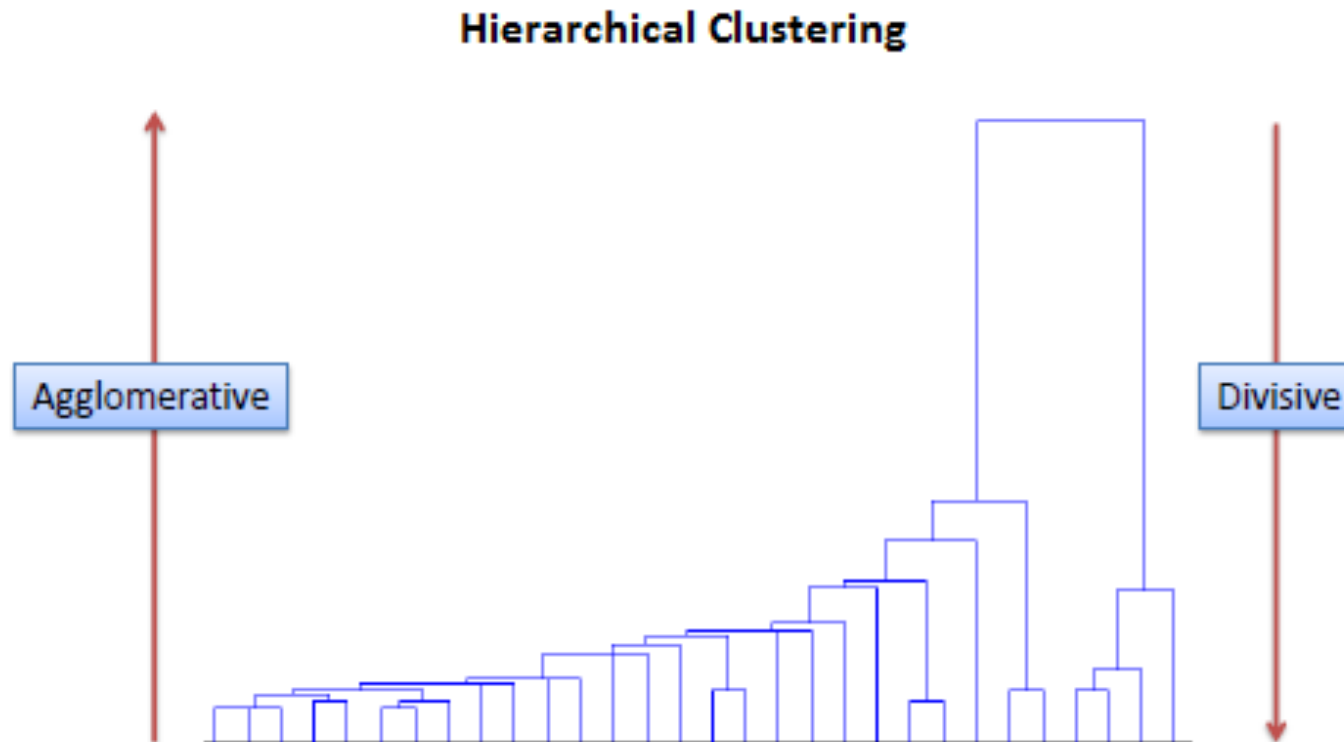→ single-linkage (SL), complete-linkage (CL), average-linkage (AL): bottom-up
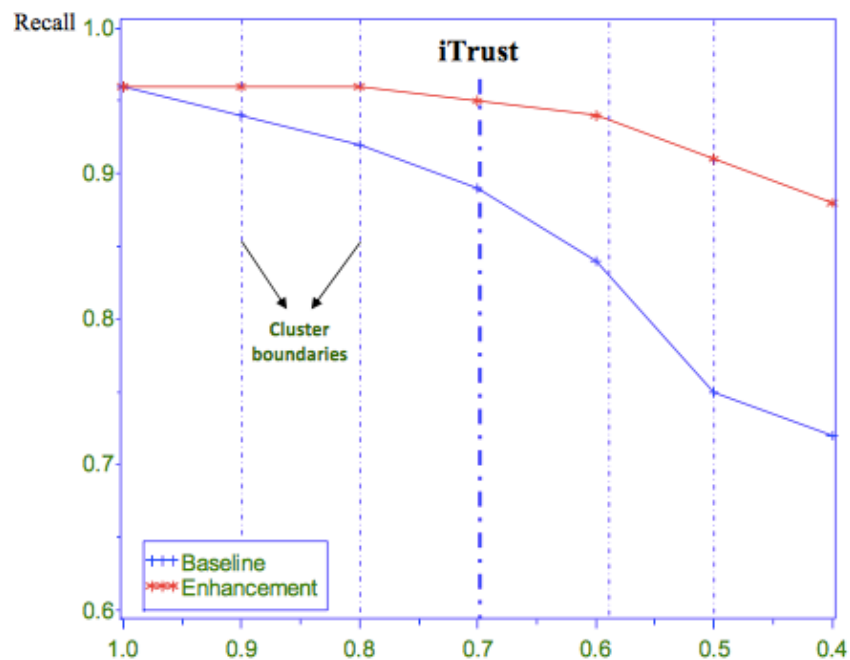
# k-means (e.g., k=2)

# Hierarchical Clustering

# How much better?

# ASN4 grader program

# RE Story (ASN2) Schedule

→ Check 'ASN2-Schedule' on the course website

→ We'll start at 9am on Monday (July 26)

→ Each presenter has 5-10 minutes

→ All the students are required to attend all the presentations

| | |
|---|---|
| 1. Jiachang | 10. Muyu |
| 2. Weijiang | 11. Luyao |
| 3. Bo Z. | 12. Yuqi |
| 4. Bo L. | 13. Xiling |
| 5. Jinzhi | 14. Xiaye |
| 6. Hongrong | 15. Chenxi |
| 7. Xiaoyu | 16. Zimao |
| 8. Miaoyu | 17. Zichun |
| 9. Shuang | 18. Pengxiang |