

Unsupervised Learning of  
Hierarchical structures in data  
using **Co-occurrence Analysis**

Dr. Shailesh Kumar

# World – a “Hierarchy of Objects”

**“Entities” at one level of the hierarchy combine in certain “contexts” to form entities at the next level**

- **Natural Systems are hierarchical**
  - Atomic Particles → Atoms → Molecules → Crystals → ...
  - Cells → Tissues → Organs → Systems → Individuals → Societies
  - Pixels → Edges → Regions → Objects → Images → Videos
- **Man-made systems are also hierarchical**
  - Letters → Words → Phrases → Sentences → Paragraphs → Documents
  - Tokens → Lines → Functions → Modules → Programs → Software
  - Employees → Teams → Groups → Department → Company → Industry

# The Quest

- **Can we learn this hierarchy from lots of data?**
  - Text, Image, Speech, Genetic, Neural, Retail,...
- **What are the “Building blocks” of such a hierarchy?**
  - Domain agnostic, Level agnostic, Data agnostic...
- **How do we adapt these principles to different domains?**
  - Nature of data (sets, sequences, 2-D sequences), etc.
- **Can we discover how the brain builds hierarchies?**
  - Already there is a lot of evidence around this

# Two building blocks of Learning

FIN\_INDICATOR CHANGED COMPANY EARNINGS

the dow jones industrial average fell  
as ford motor co. released its first  
quarter earnings

STOCK\_EXCHANGE

“NYSE”

“london stock exchange”

“bombay stock exchange”

“hong kong stock exchange”

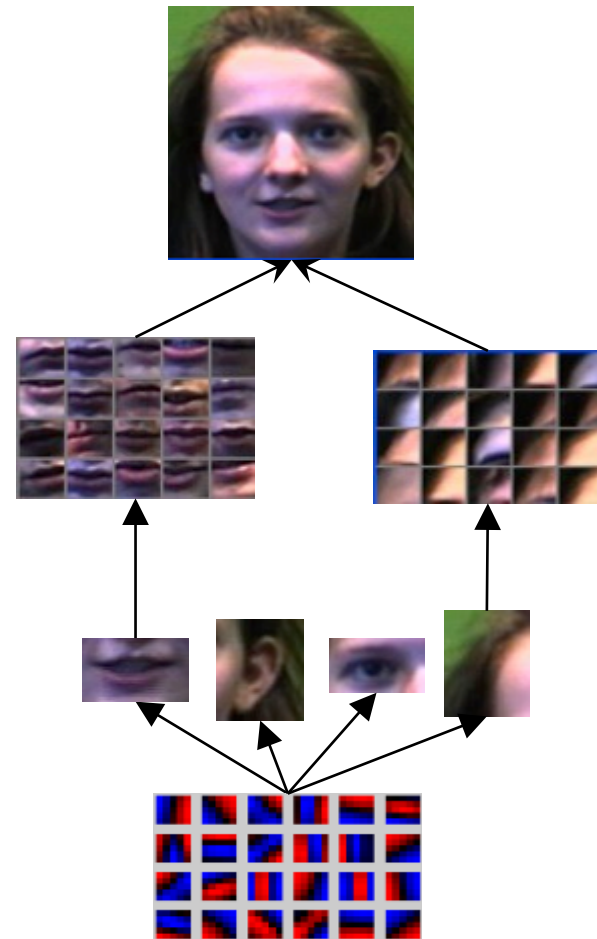
“new york stock exchange”

“new” “york” “stock” “exchange”

Syntactic Composition

Semantic Equivalencing

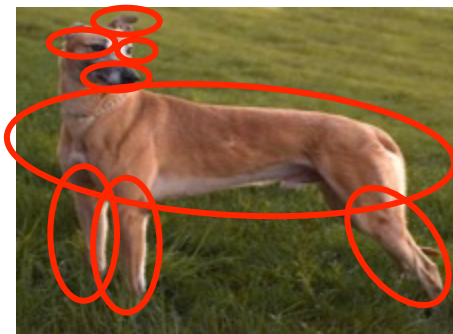
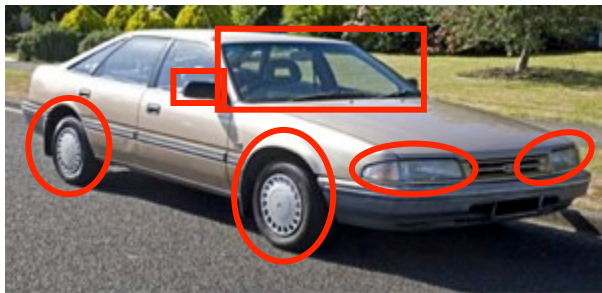
Syntactic Composition





# Two building blocks of learning

- Syntactic Composition/Conjunction – “is a PART of”



- Semantic Equivalence/Disjunction – “is a TYPE of”



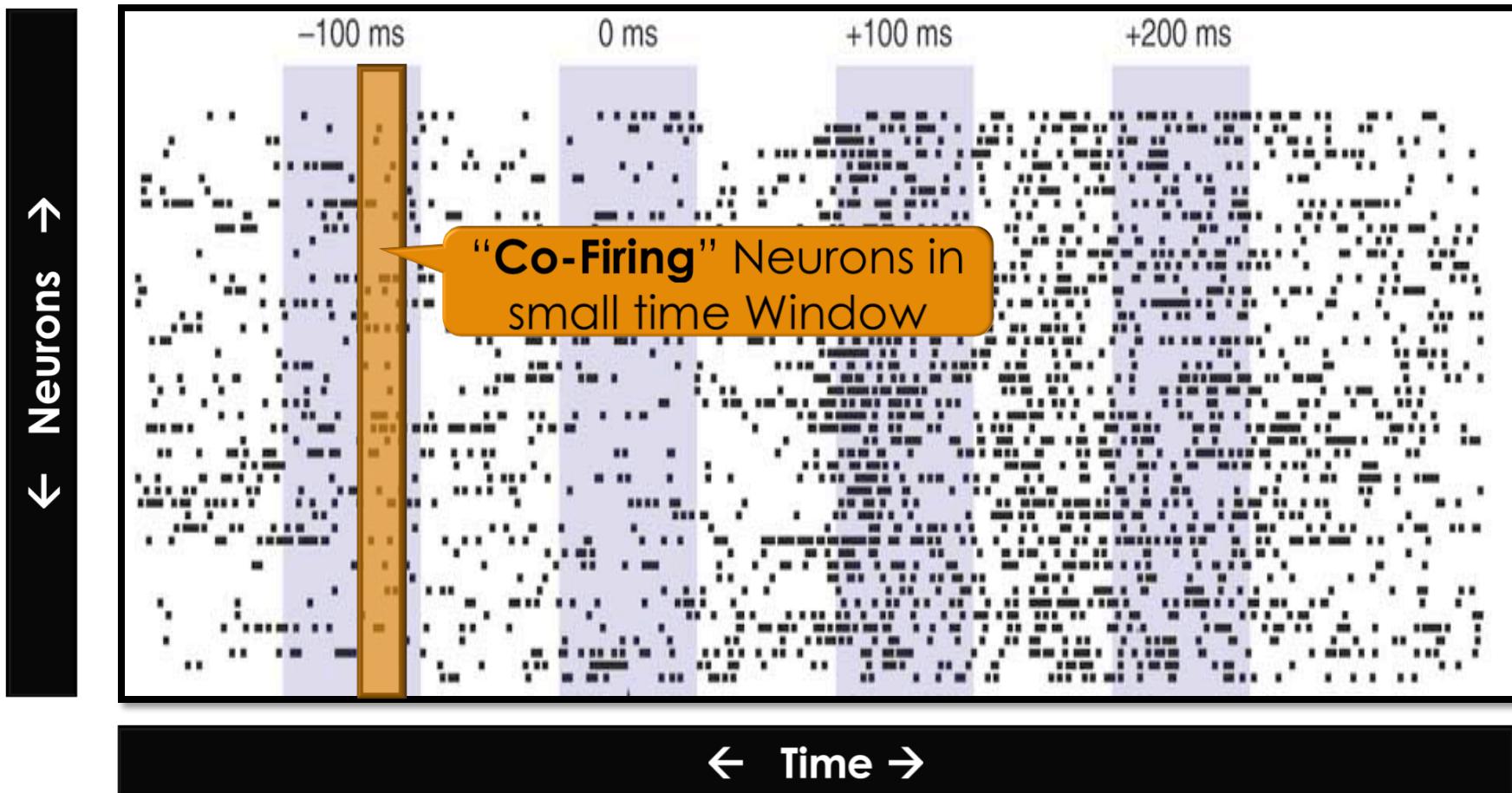
# Co-occurrence Analytics

A scalable, unsupervised, hierarchical framework that

- Analyzes pair-wise relationships among entities
- Co-occurring in various contexts
- To build a Co-occurrence Graph(s) in which
- It discovers coherent higher order structures

# Discovering Memory Engrams

# Spiking Neurons – A “Crazy Haystack!”

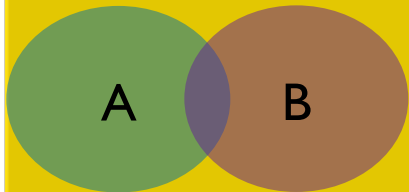


# “Co-Firing” Consistency Graph

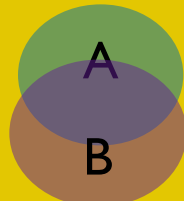
Memory Engrams =  
**Cliques** in the  
Co-Firing Graph

$$\phi_{a,b} = \log \left( \frac{P(a,b)}{P(a)P(b)} \right)$$

**Consistency: Strength**



Low

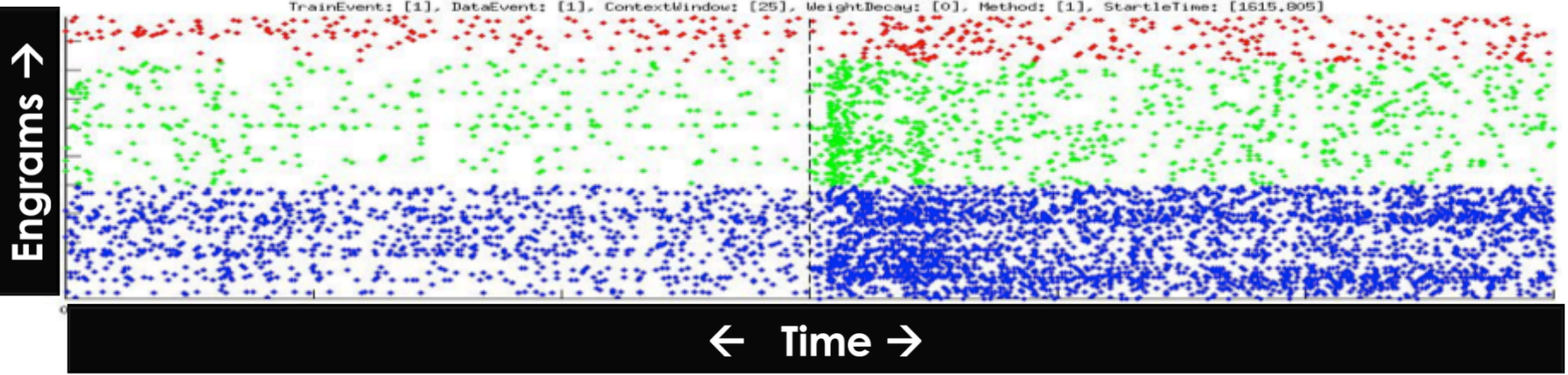


High





# Memory engrams = Co-firing neurons



# Co-occurrence Analytics - steps

- **Context** – Nature of Co-occurrence
  - E.g. Window of time
- **Co-occurrence** – Definition of Co-occurrence
  - E.g. Co-firing
- **Consistency** – Strength of Co-occurrence
  - E.g. Point-wise Mutual Information
- **Coherence** – Tightness of a group of entities
  - E.g. Modularity or Cliqueness score
- **Community** – Locally optimal groups of entities
  - E.g. Memory engrams

# Logical Item-Set Mining



# Why **Market Basket Analysis** failed?

Few buy a complete “logical” item-set in same basket

- ❑ already have other products
- ❑ buy them from another retailer
- ❑ buy them at a different time
- ❑ got them as gifts
- ❑ didn't know they needed it



**Projection** of **latent** customer **intentions**

# The right needle in the wrong haystack?

## Market Basket

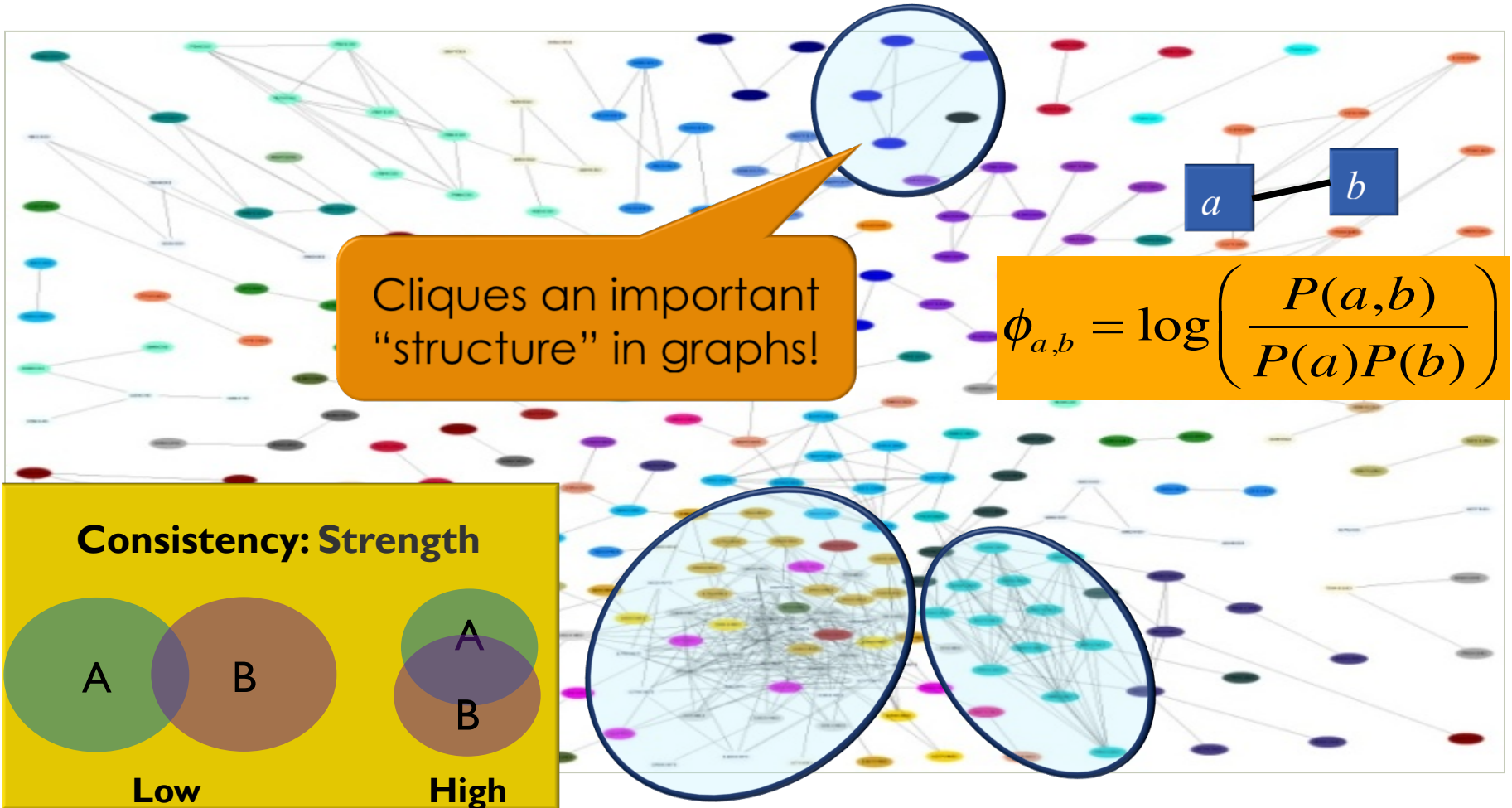


Intention: **Home-office**

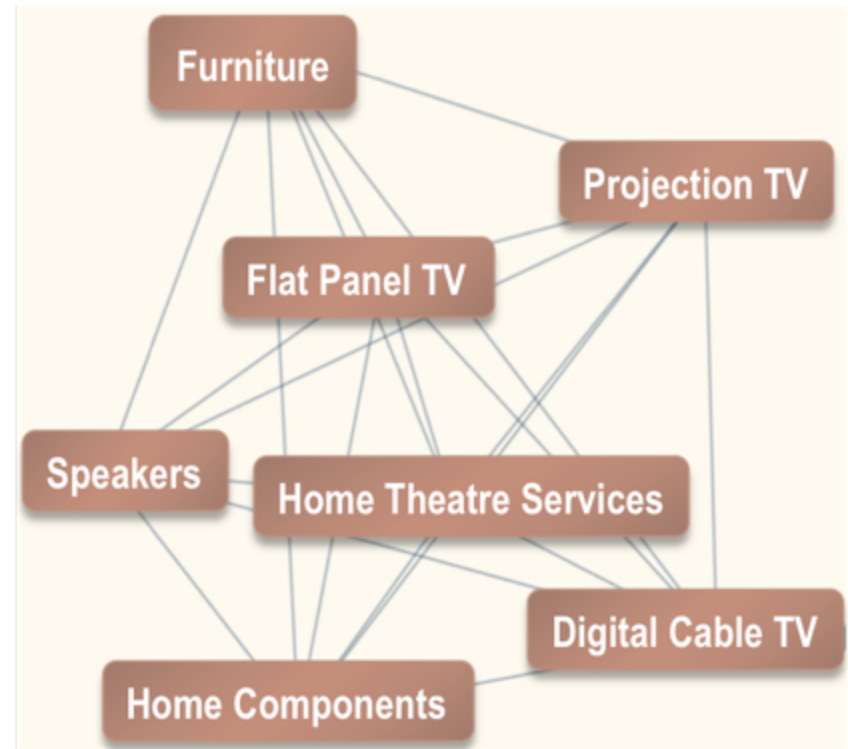
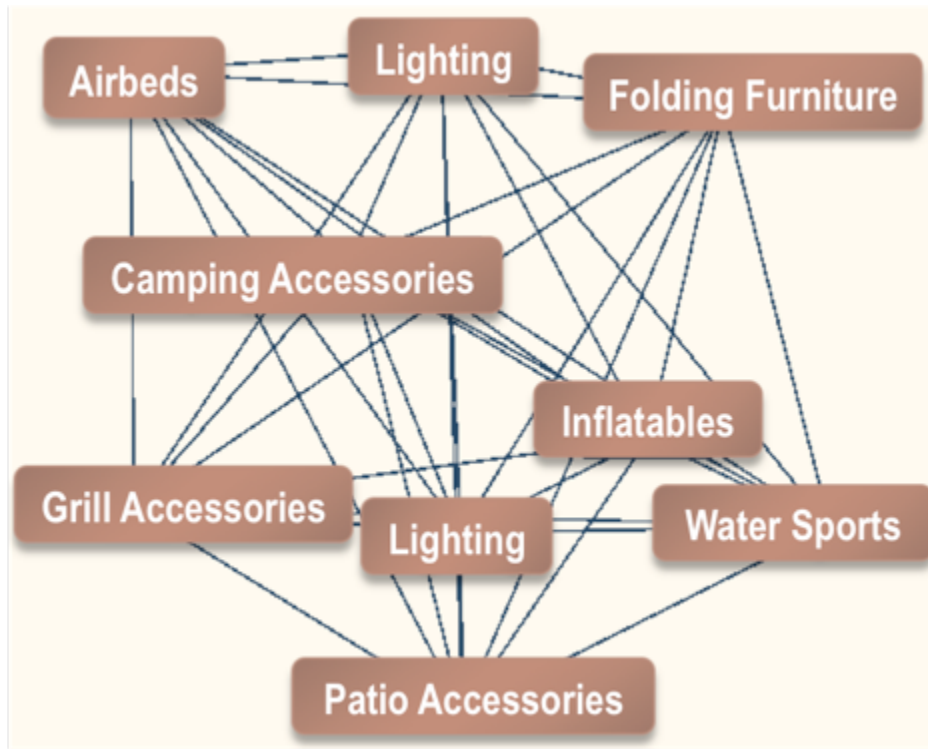
Intention: **Home-tools**

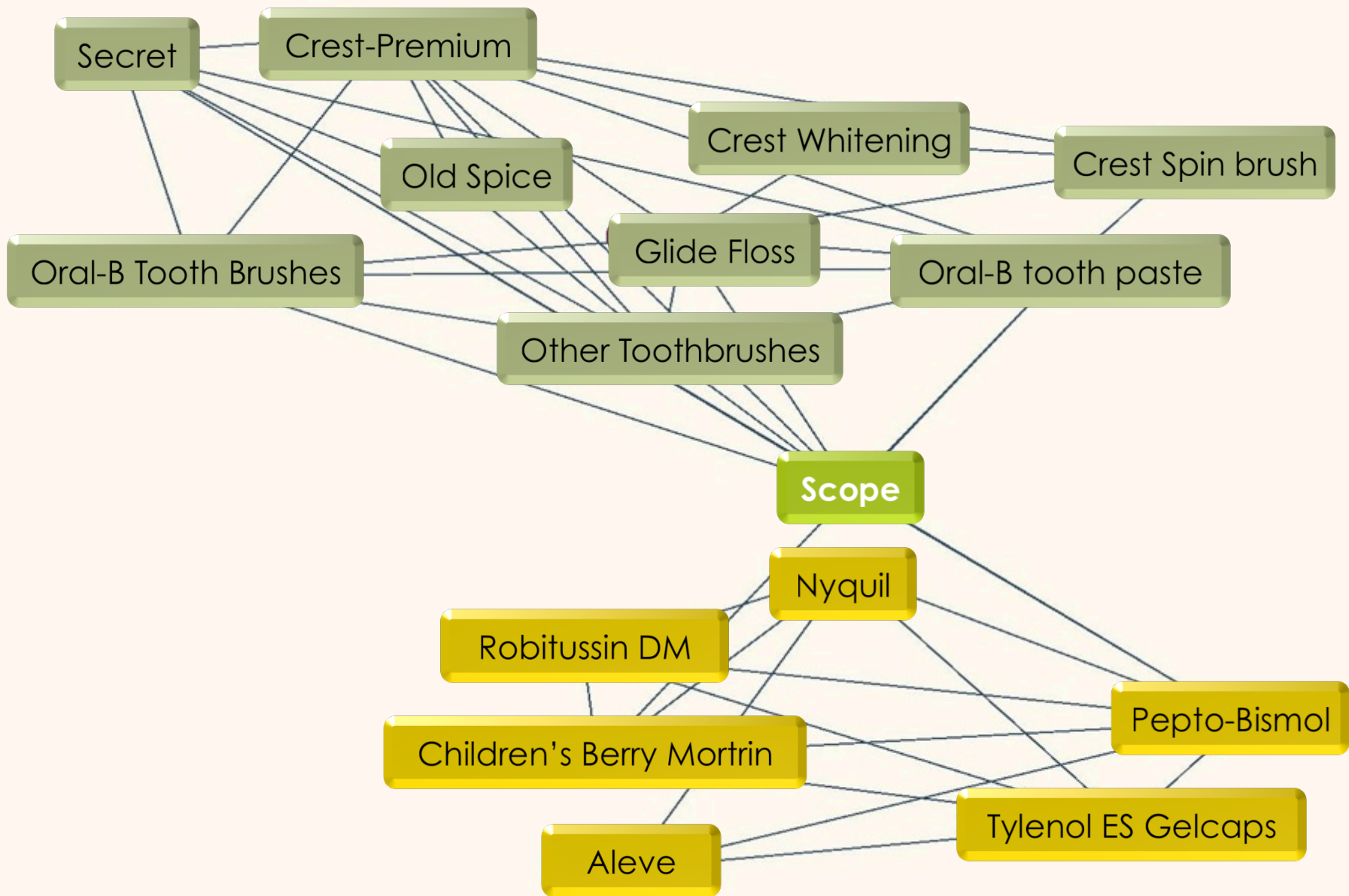
**Mixture** of **Projections** of **latent** customer **intentions**

# Extract Knowledge, throw away Data



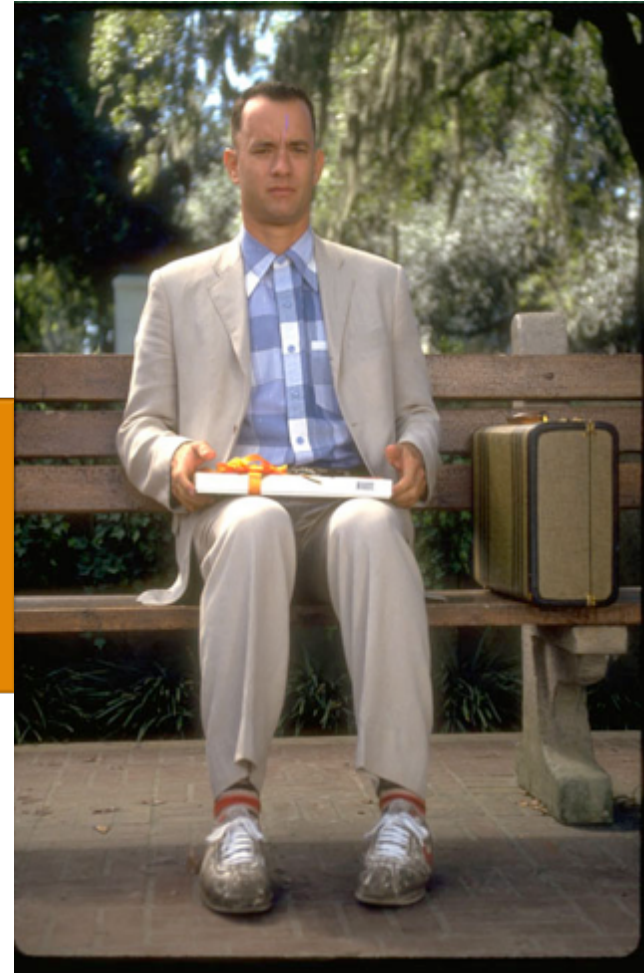
# Reconstruct needles, don't find them!







“**Data Mining** is like a box of chocolates. You never know what (**Insights**) you're gonna get.”!!!



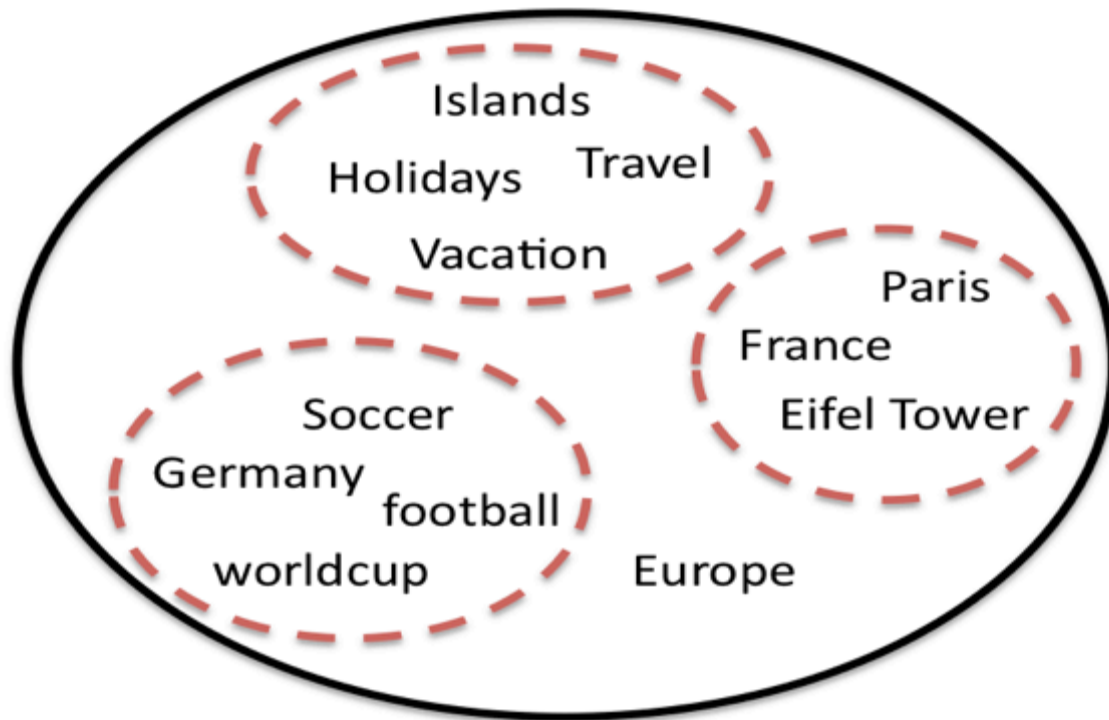
# Discovering Tag Communities

# Community Detection in Tagsets

- **Tagset data...**
  - **Flickr** – tags describing **images**
  - **YouTube** – tags describing **videos**
  - **AdWords** – tags describing **advertisements**
  - **IMDB** – tags describing **movies**
  - **Keywords** – tags describing **scientific publications**
- **Key Challenges**
  - **Noisy Tag-sets** – create robust graphs!
  - **Weighted graphs** – communities in weighted graphs
  - **Overlapping communities** – no clustering!

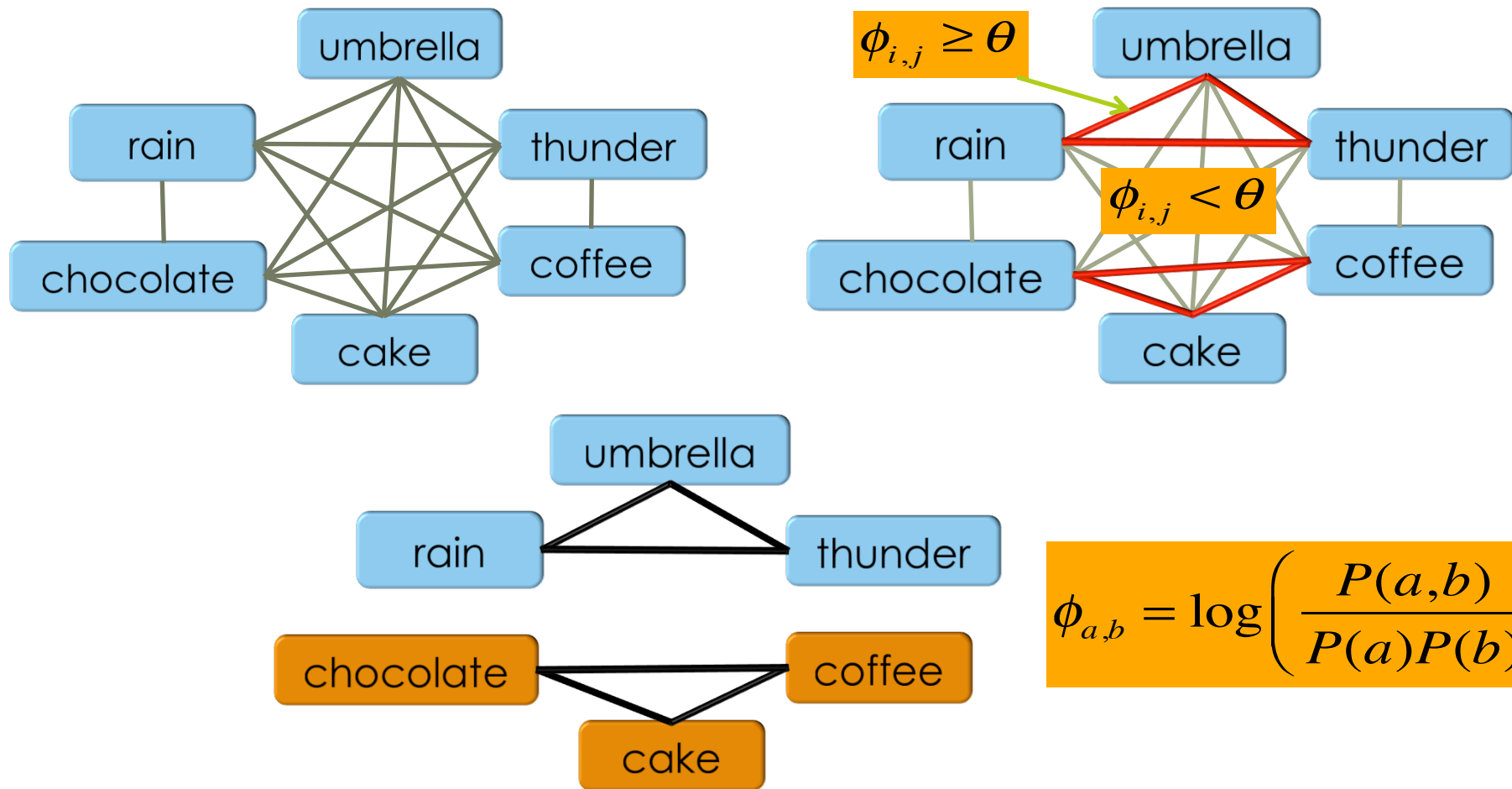


# Tagsets – a “Crazy Haystack”!



**Mixture** of **Projections** of **latent** Concepts

# Creating Robust Co-oc Graph



# De-noising – for better graphs

Co-occurrence of Tags with tag “**wedding**”

Tag	Before Denoising	After Denoising
<b>bride</b>	0.3257	0.5750
<b>reception</b>	0.3720	0.5728
<b>marriage</b>	0.3195	0.5658
<b>cake</b>	0.1699	0.3629
<b>love</b>	0.0148	0.2449
<b>honeymoon</b>	0.0183	0.2262
<i>jason</i>	0.2081	0
<i>chris</i>	0.1461	0

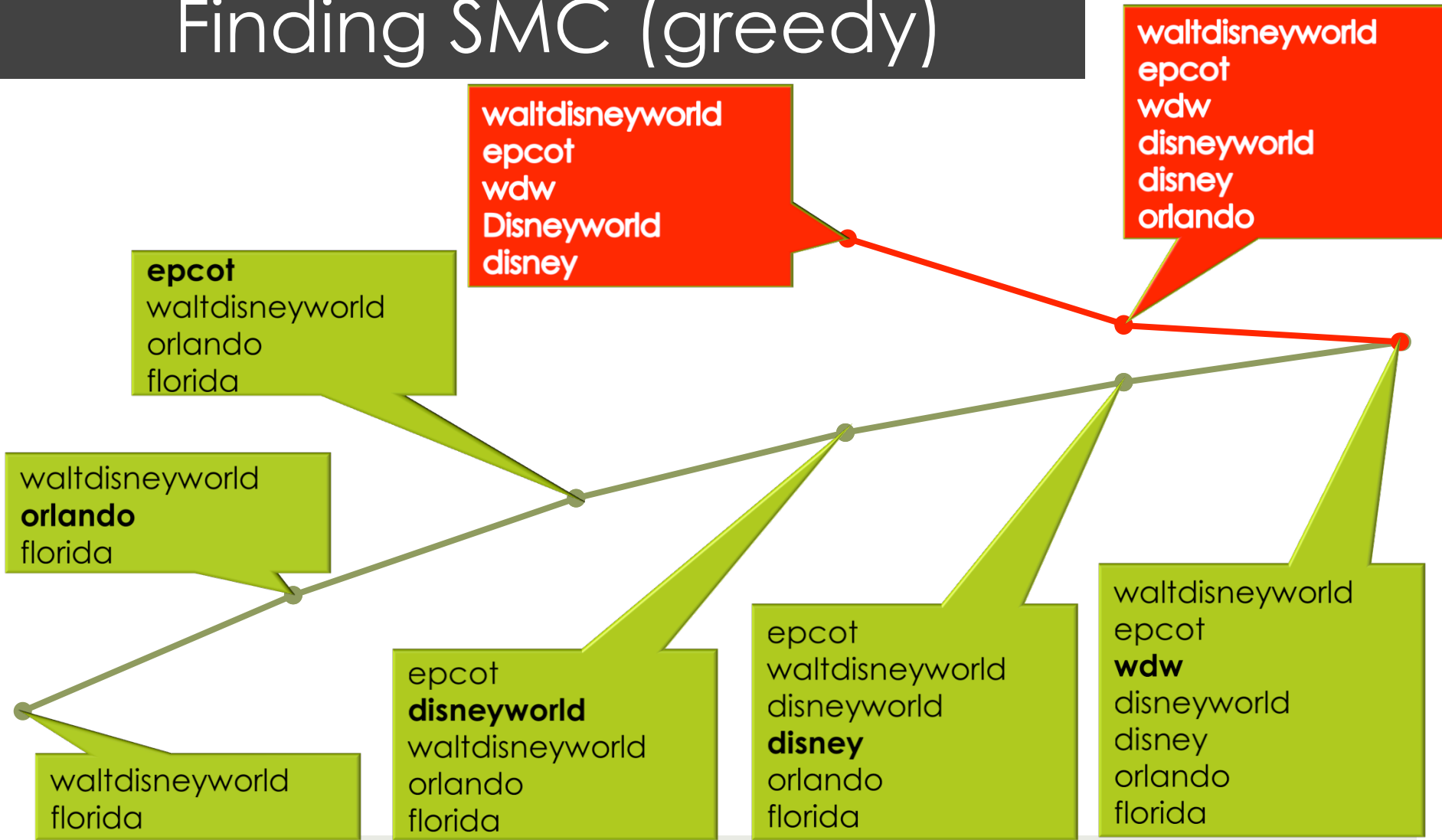
# Consistency + Denoising

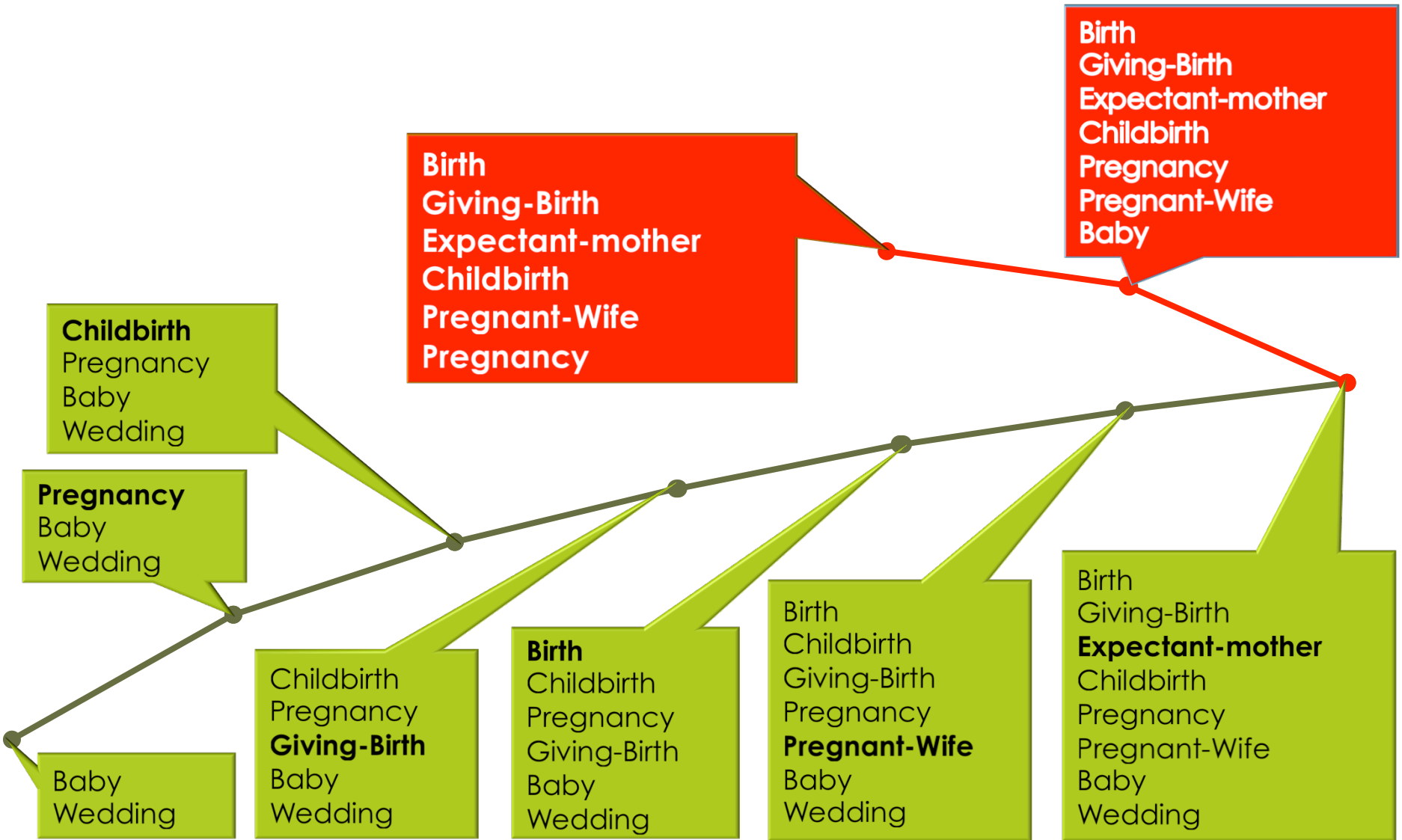
Tag	Most Consistent Tags in IMDB dataset
<b>food</b>	lifestyle, money, restaurant, drinking, cooking
<b>road</b>	truck, motorcycle, car, road-trip, bus
<b>singer</b>	singing, song, dancing, dancer, musician
<b>suicide</b>	suicide-attempt, hanging, depression, mental-illness, drowning
<b>hospital</b>	doctor, nurse, wheelchair, ambulance, car-accident

Tag	Most Consistent Tags in FLICKR dataset
<b>art</b>	painting, gallery, paintings, sculpture, artist
<b>france</b>	paris, french, eiffeltower, tower, europe
<b>island</b>	tropical, islands, newzealand, thailand, sand
<b>animals</b>	zoo, pets, wild, cats, animal
<b>airplane</b>	flying, airshow, fly, military, aviation

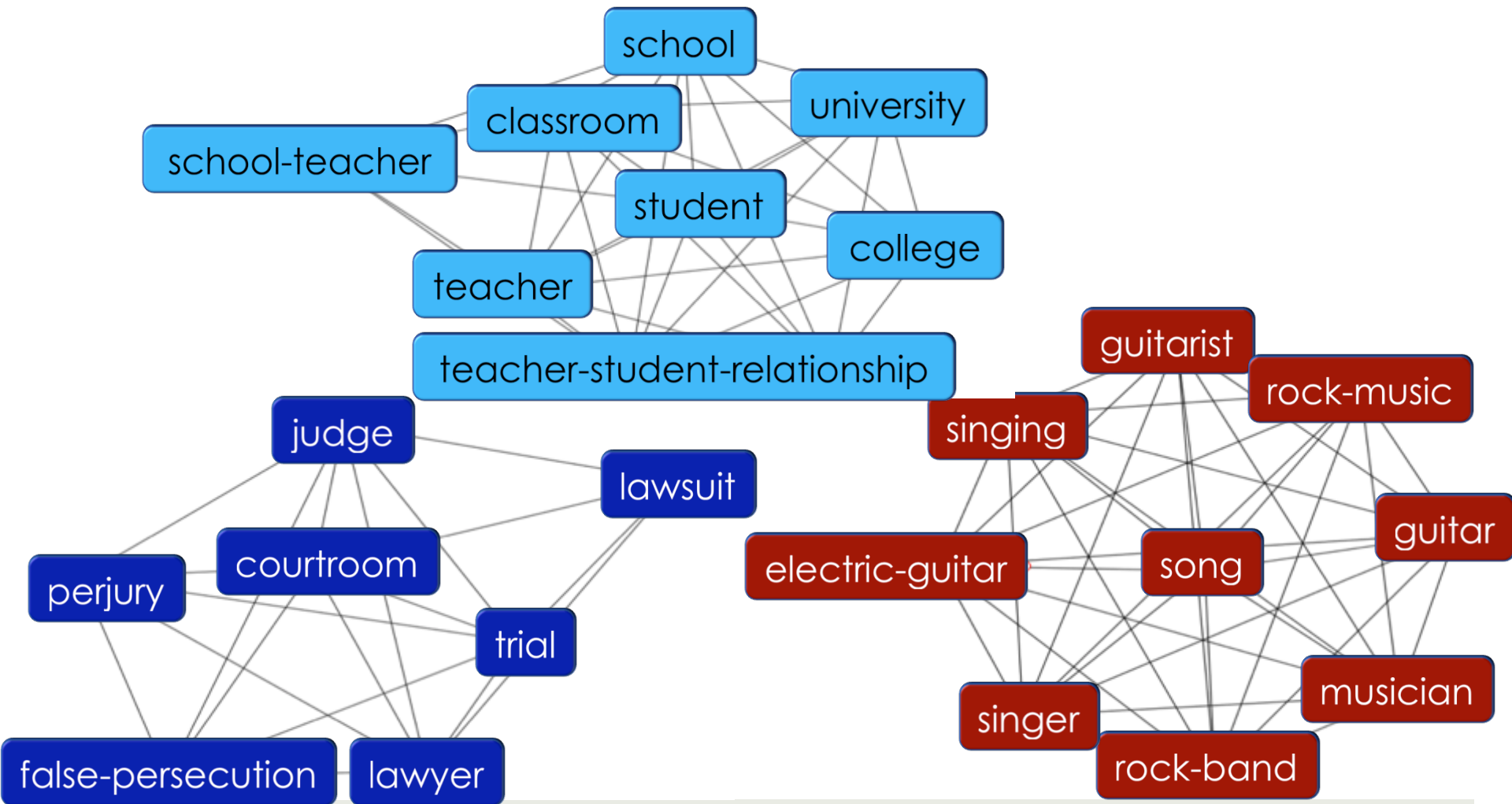


# Finding SMC (greedy)





# Discovered Communities (IMDB)





**“You shall know a word by the company it keeps”**

*– John Rupert Firth, 1957*

# Finding Phrases in Text

# Unsupervised Text “Enrichment”

## ▶ Syntactic Tokenization

- ▶ He distributes **Time** magazine in **new york**
- ▶ Today **new york times** reported **new** rise in crime

## ▶ Semantic Tokenization / Disambiguation

- ▶ I was **right** to avoid a **suit** against **apple**
- ▶ Man in red **suit** on my **right** was drinking **apple** juice

## ▶ Conceptual Tokenization

- ▶ **filed** a **suit** **charging** **orange** of **illegal** **behavior**
- ▶ **submitted** a **case** **accusing** **apple** of **unauthorized** **conduct**

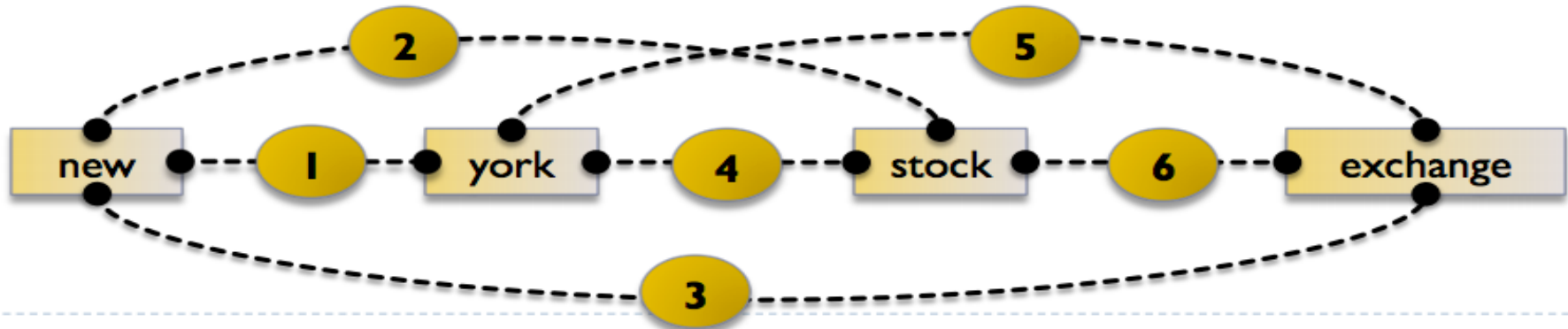
## ▶ Contextual Weighting

- ▶ rain, thunder, umbrella, lightening, **chocolate**
- ▶ kids, birthday, candies, **chocolate**, cake, candles

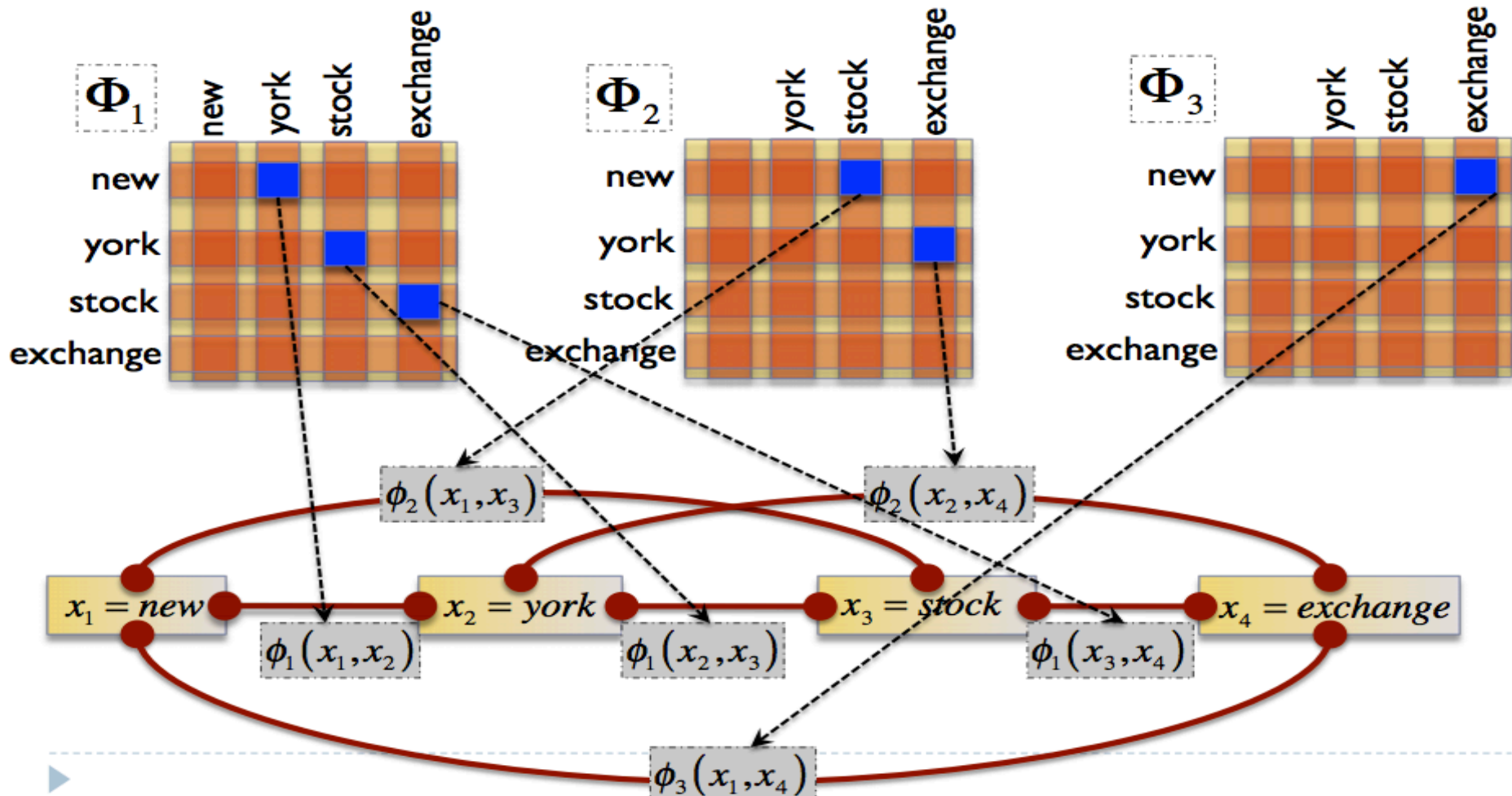
# Positional Bigrams $O(nV^2)$

**Positional Bi-grams** = Co-occurrence of a pair of **words** in a **specific relative position** w.r.t. each other.

1. Word **new** occurs **one position before** word **york**
2. Word **new** occurs **two positions before** word **stock**
3. Word **new** occurs **three positions before** word **exchange**
4. Word **york** occurs **one position before** word **stock**
5. Word **york** occurs **two positions before** word **exchange**
6. Word **stock** occurs **one position before** word **exchange**

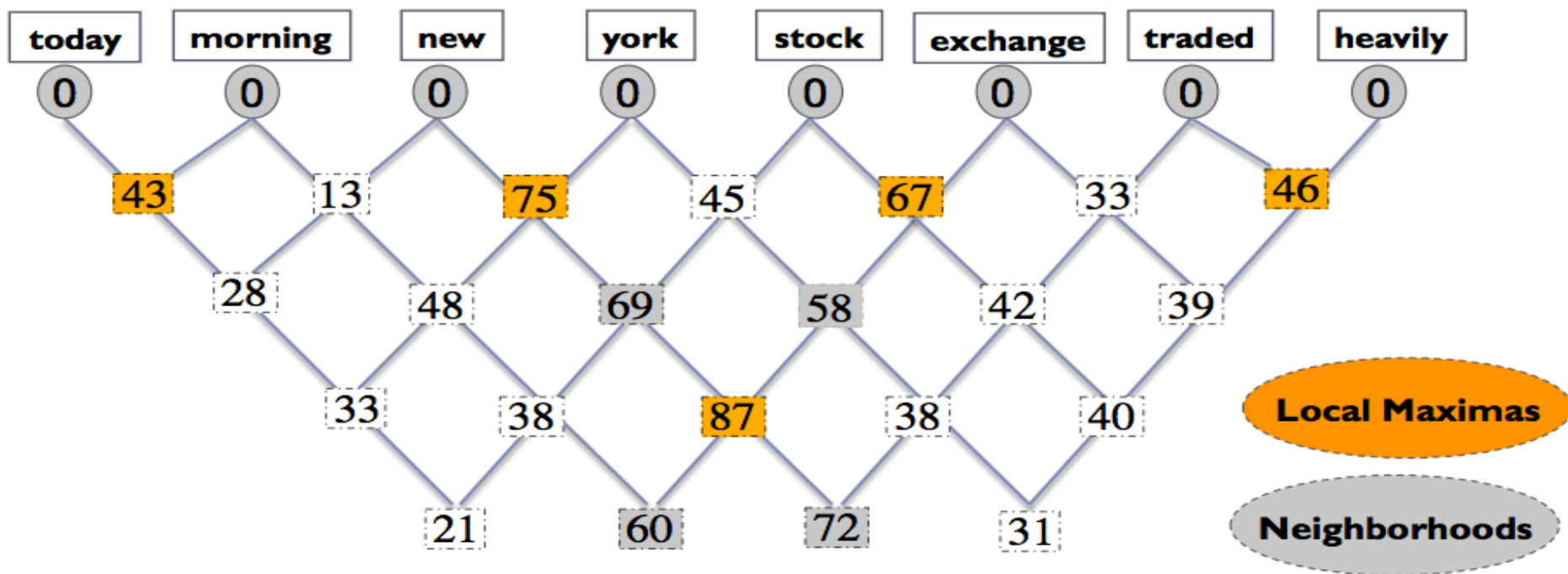


# Using Positional Bigrams



# Phrase: Soft Maximal Sequence

Local maxima in “Sequence coherence lattice”



# Phrases in Wall Street Journal data

## High Cohesiveness Low Frequency

russian president boris yeltsin  
prime minister viktor chernomyrdin  
prime minister ryutaro hashimoto  
palestinian leader yasser arafat  
serbian president slobodan milosevic  
syrian president hafez al-assad  
british foreign secretary douglas hurd  
iraqi deputy prime minister tariq aziz  
secretary general kofi annan  
senate majority leader bob dole  
lieutenant general raoul cedras  
nba commissioner david stern  
bulls coach phil jackson  
federal reserve policy makers

## Medium Cohesiveness Medium Frequency

new york stock exchange  
dow jones industrial average  
nasdaq composite index  
new york mercantile exchange  
automated teller machines  
international monetary fund  
u . n . officials  
u . n . peacekeepers  
u . n . spokesman  
an international tribunal  
intensive care unit  
highway traffic safety administration  
mitsubishi heavy industries ltd  
rio de janeiro  
freshly ground black pepper

## Low Cohesiveness High Frequency

he would not elaborate  
sometime next year  
few weeks before  
nearly two years ago  
made no comment  
gain market share  
tons of cocaine  
had no chance  
made no comment  
materials contained herein  
people have been killed  
as soon as possible  
dollar rose as high as  
threatened to pull out  
forced to give up  
as far as possible  
as early as tomorrow  
kilometers ( <num> miles ) north  
kilos ( <num> pounds )

# And... It's language agnostic

世界贸易组织	World Trade Organization	10.504958
总统克林顿	President Clinton	9.011669
亚洲金融危机	Asian Financial Crisis (1997)	9.702676
巴塞罗那奥运	Barcelona Olympics (1992)	10.333340
总统布什	President Bush	7.339422
俄罗斯总统叶利钦	Russian President Yeltsin	11.698954
糖尿病	diabetes	11.583491
艾滋病毒	Human immunodeficiency virus	14.140515
络绎不绝	in endless stream	8.637969
紧锣密鼓	an intense publicity campaign	11.915520
兴致勃勃	be highly interested in	11.415492

“You shall know a **PIXEL** by the company it keeps”

Recognizing “objects” in Images

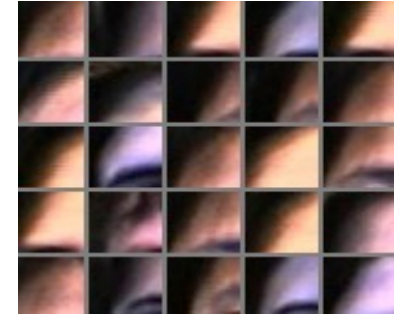


# Images – the “Craziest Haystack!”

- Image Understanding (the way humans do) is very difficult
- Images have far too many degrees of freedom!
  - Translation, Scale, Rotation, Illumination, Color, Shadows, ...
- Even for Human brain its hard!
  - 30-40% of brain processes vision!
- Most computer vision systems:
  - Use Low level features
  - Complex models to compensate
- SIFT features →
  - Still don't capture semantics



# Phrases and Concepts in Images



# Content Based Image Retrieval



# Content Based Image Retrieval





# Content Based Image Retrieval



# Content Based Image Retrieval

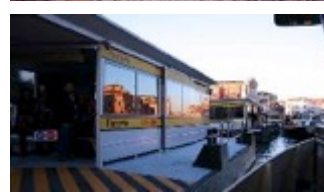




# Content Based Image Retrieval

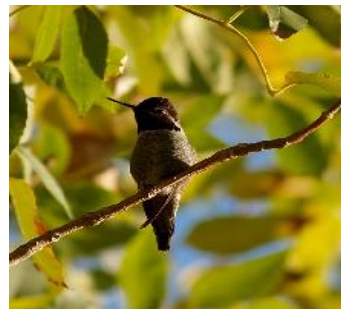
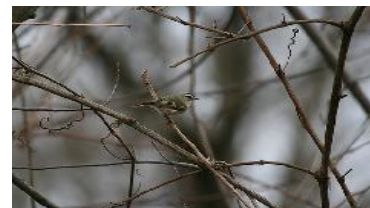
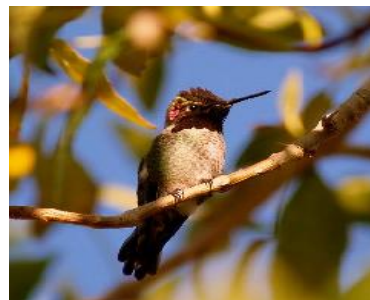


# Content Based Image Retrieval

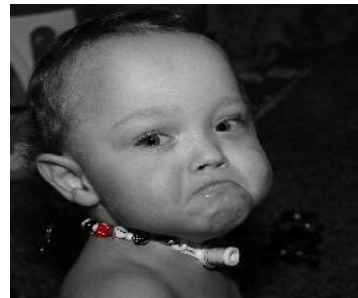
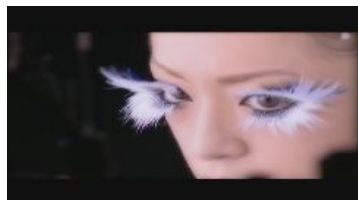




# Content Based Image Retrieval



# Content Based Image Retrieval





# Content Based Image Retrieval



# Content Based Image Retrieval



# Conclusions

- **World is a hierarchy of objects**
  - Raw data is the lowest level objects
- **Discovering the “grammar of the data”**
  - Syntactic Composition
  - Semantic Equivalencing
- **Co-occurrence Analysis Framework**
  - Unsupervised Hierarchical Structure Discovery
  - Language and Domain agnostic

