



# Large Scale Cross-Document Coreference Resolution

Kambiz Ghoorchian  
Sarunas Girdzijauskas

ghoorian@kth.se  
02/06/2014  
LCN - KTH

## **Problem**

### Cross-Document Coreference Resolution

## **Problem**

Cross-Document Coreference Resolution

## **Current Solution**

Vector Space Modeling  
Limitations

## **Problem**

Cross-Document Coreference Resolution

## **Current Solution**

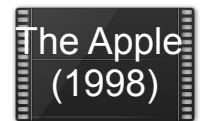
Vector Space Modeling  
Limitations

## **Contribution**

Graph Based Modeling  
Diffusion Based Clustering

# Definition

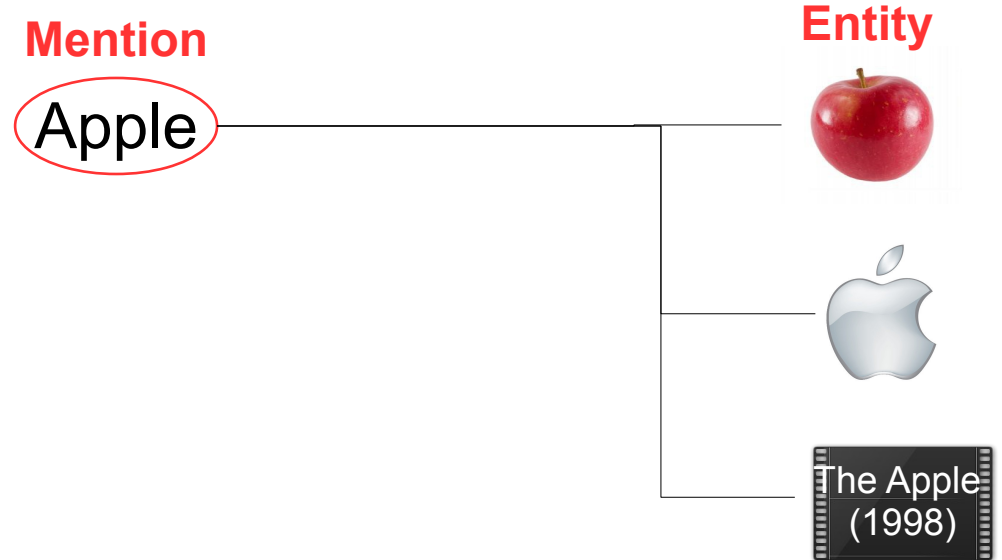
Entity



## Entity

**Distinctive** and **independent** thing in the real world.

# Definition



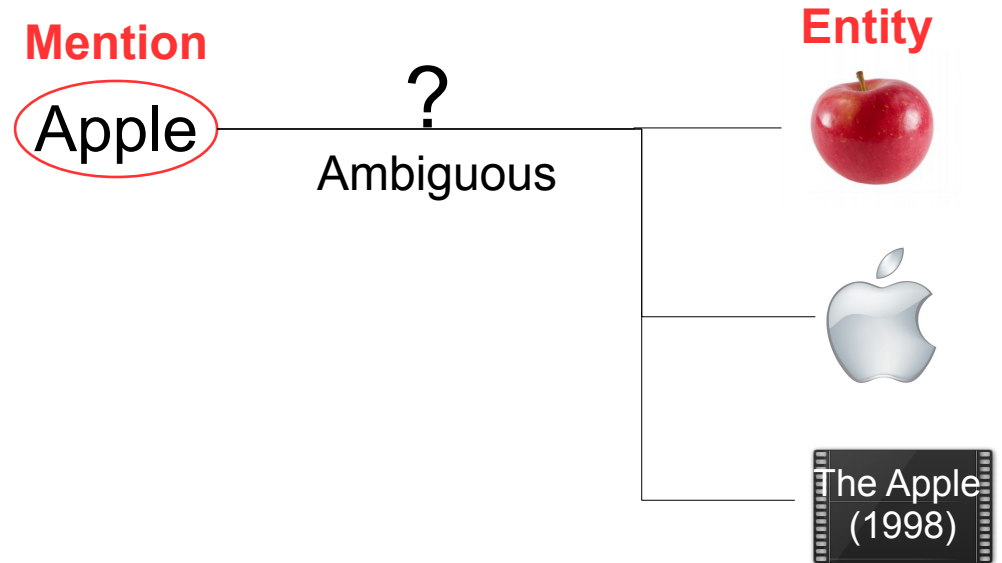
## Entity

**Distinctive** and **independent** thing in the real world.

## Mention

**Linguistic phenomenon** (word or phrase) **refers** to an entity.

# Definition



## Entity

**Distinctive** and **independent** thing in the real world.

## Mention

**Linguistic phenomenon** (word or phrase) **refers** to an entity.

# Problem

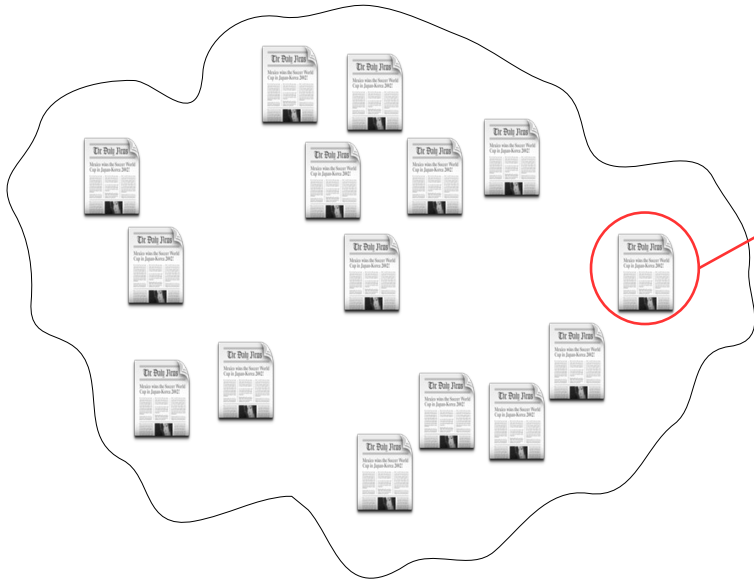
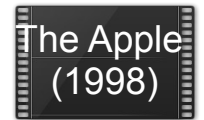
Mention

Apple

?

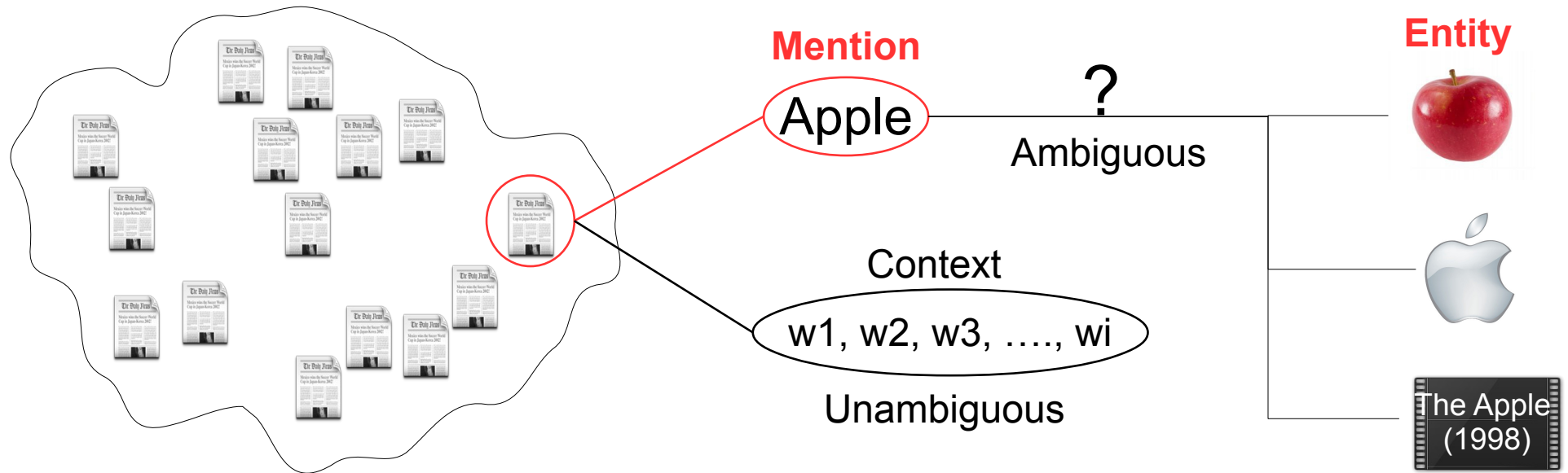
Ambiguous

Entity

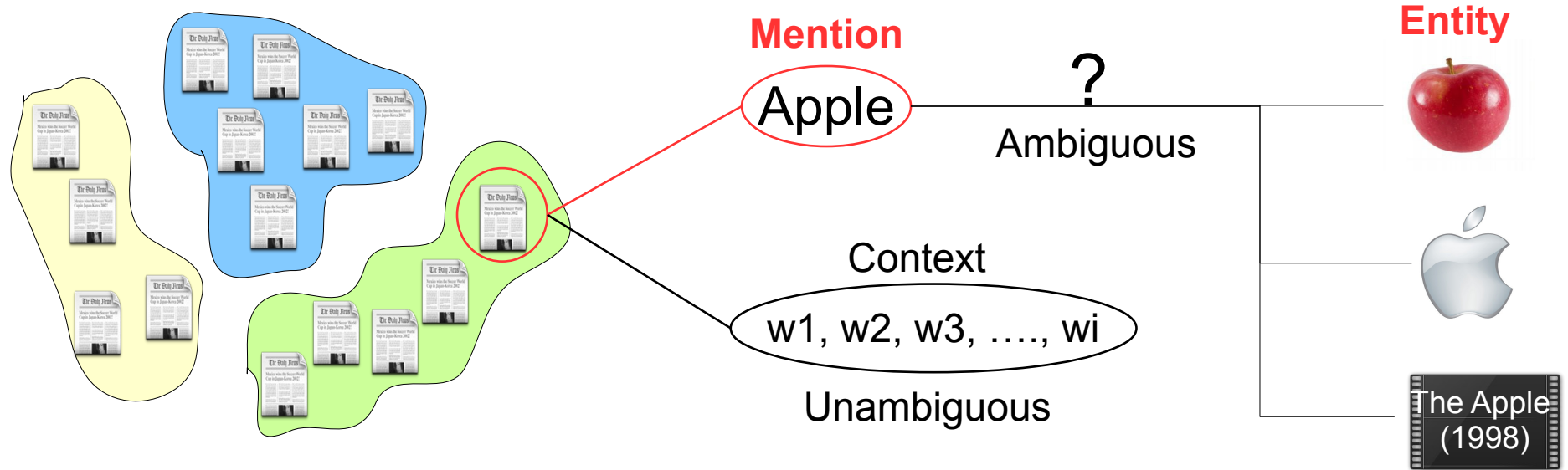




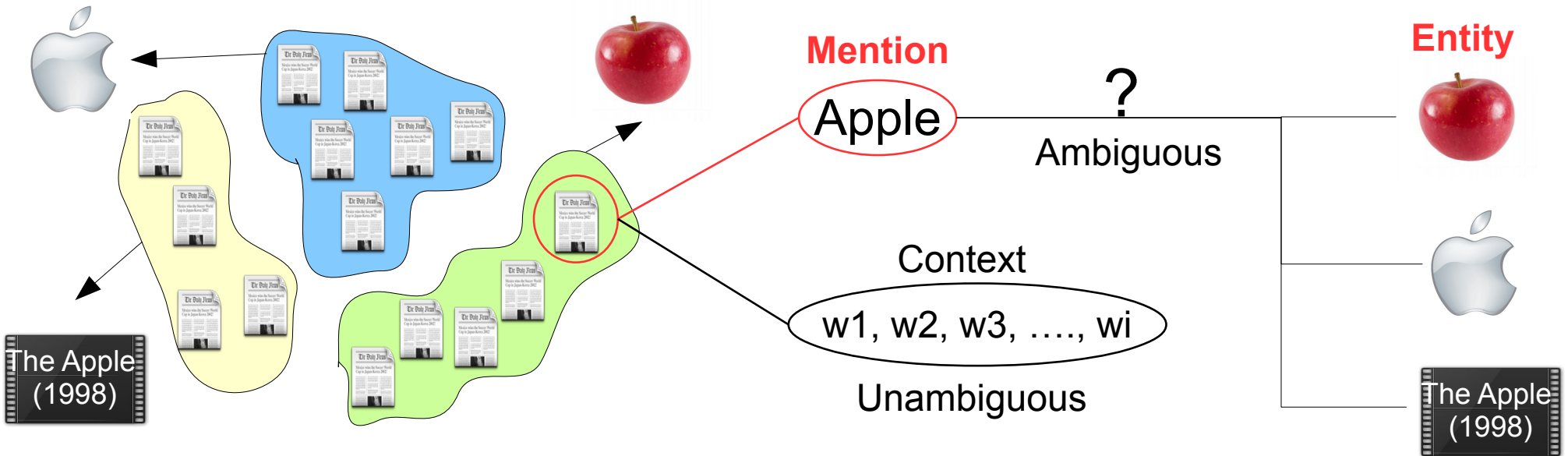
# Problem



# Problem



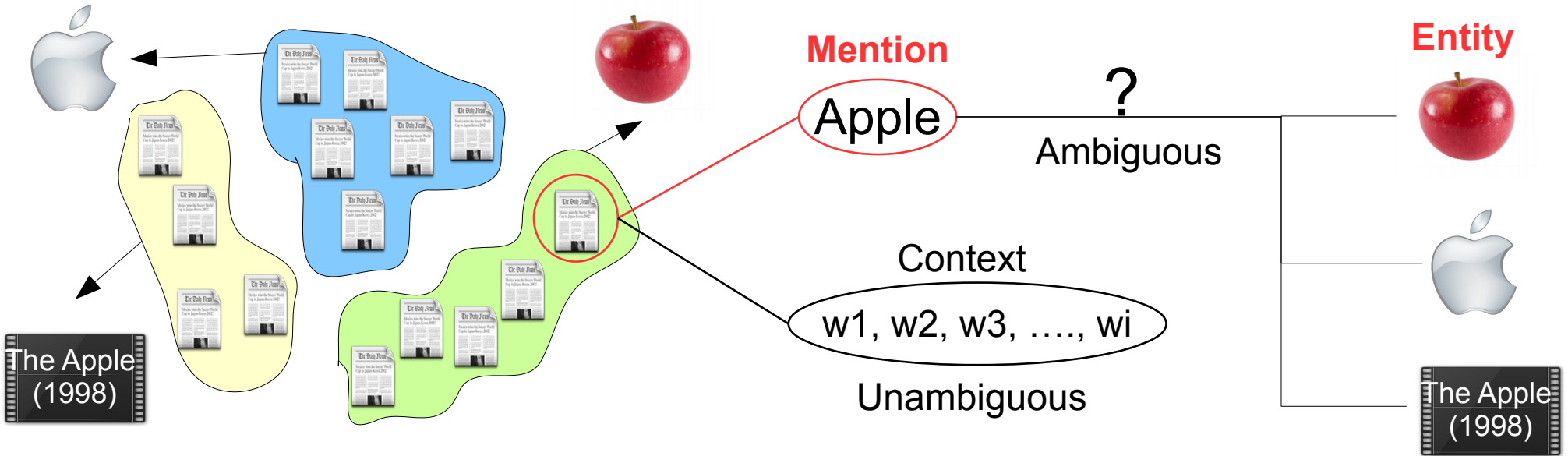
# Problem



## Questions:

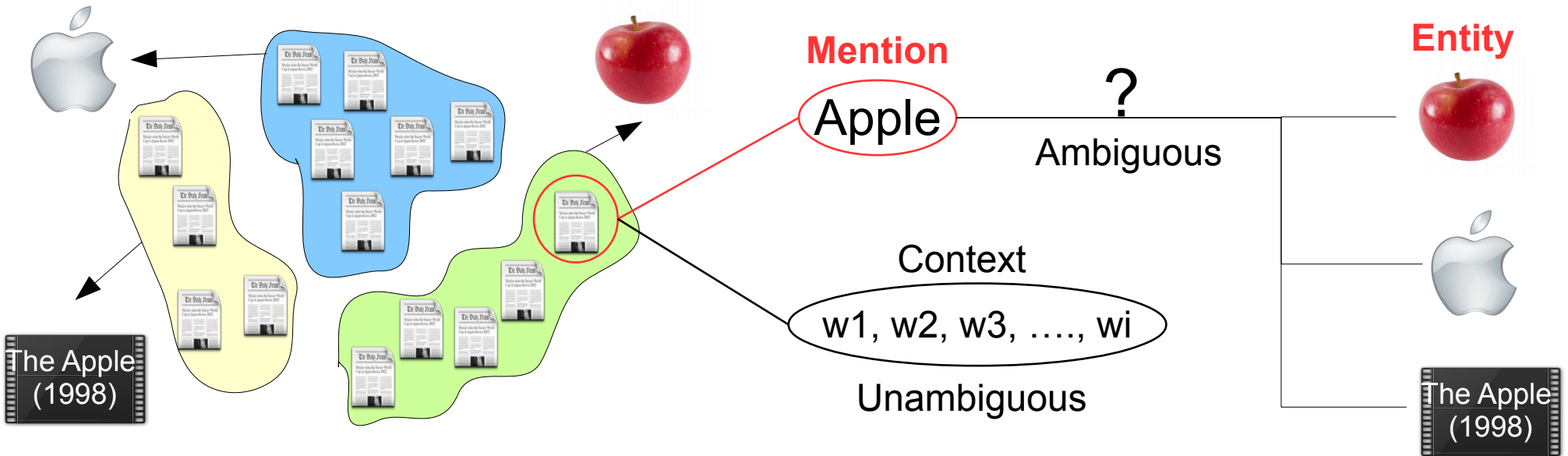
How to group these documents such that all the **Mentions** in each group refer to the same **Entity** in the real world?

# Problem



## Clustering

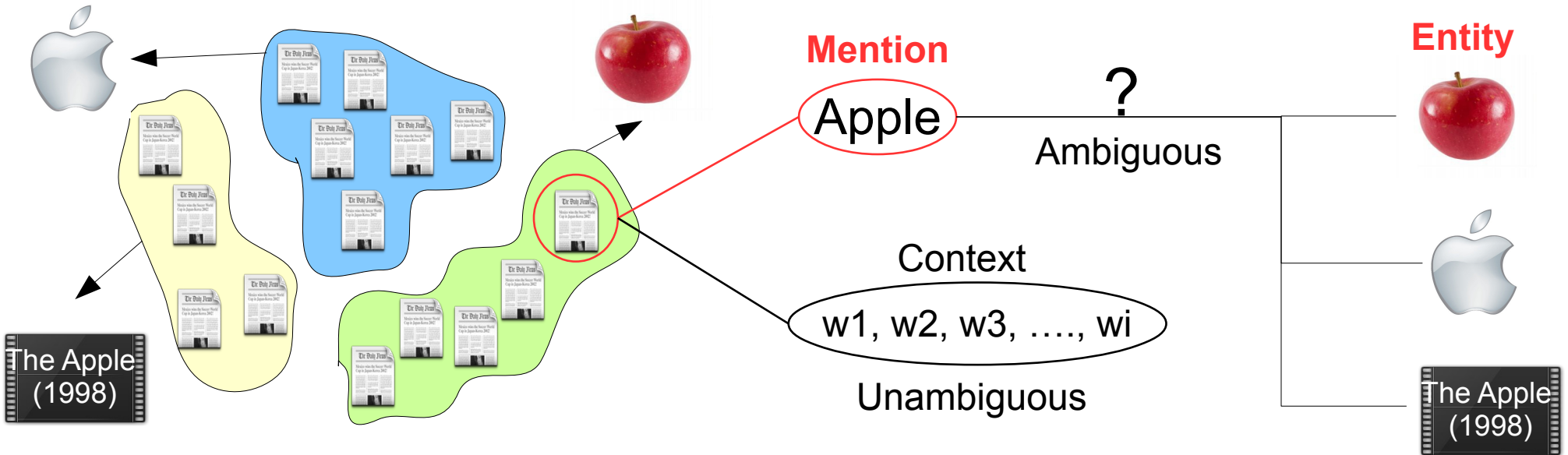
# Problem



## Clustering

Group documents based on their **similarity**  
 $\text{Similarity}(D1, D2) \propto \text{Count}(\text{common context words})$

# Problem

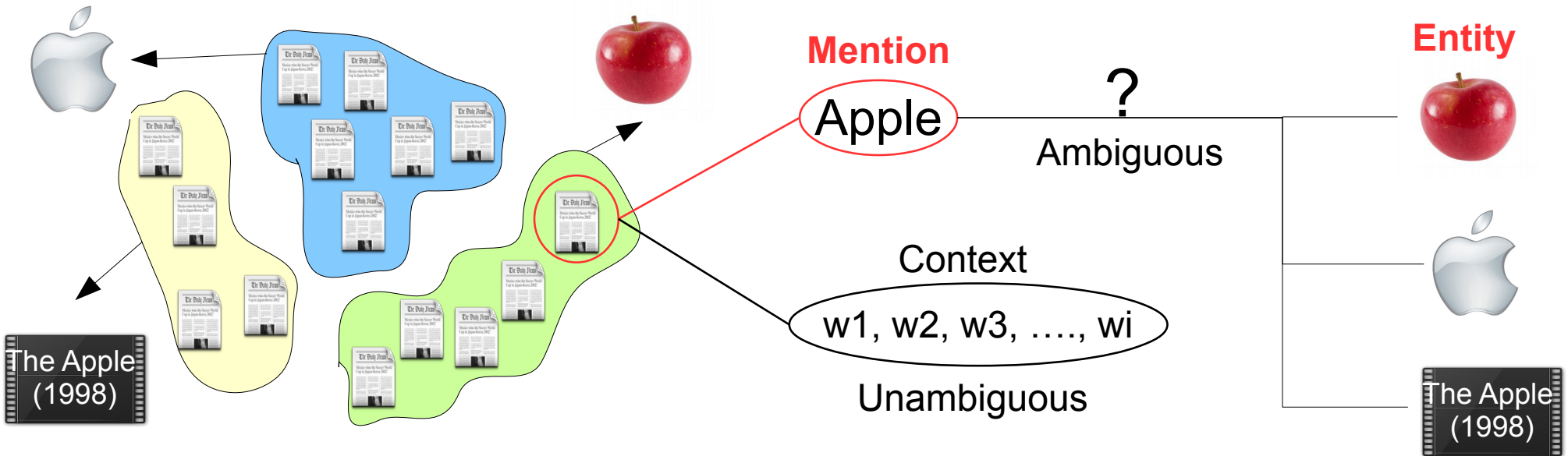


## Clustering

Group documents based on their **similarity**  
 $\text{Similarity}(D1, D2) \propto \text{Count}(\text{common context words})$

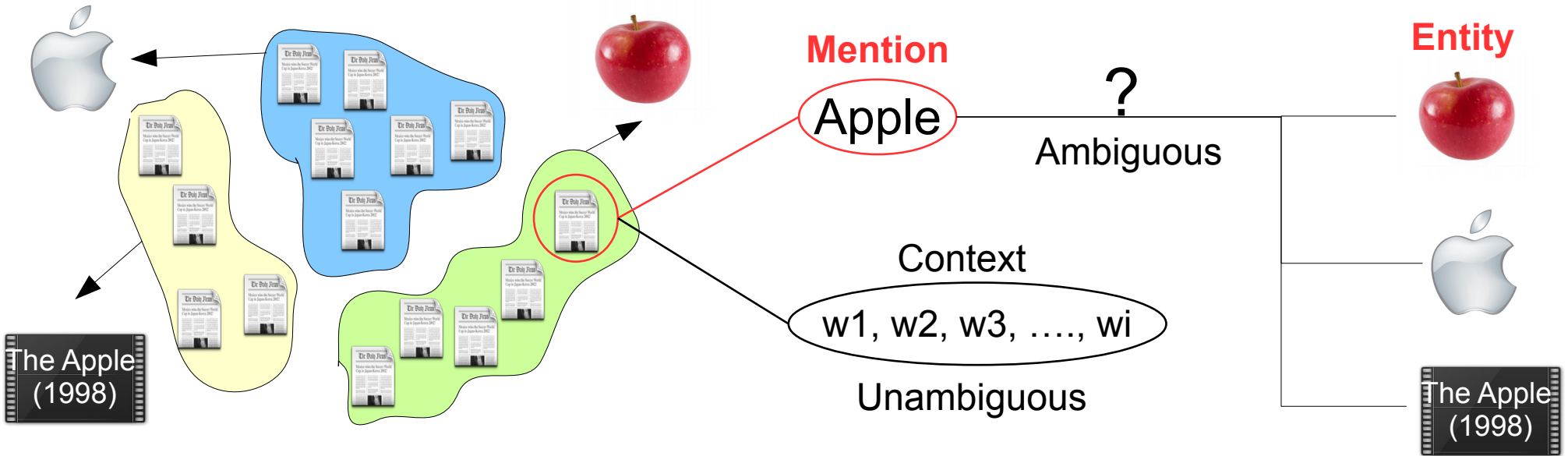
Mathematical **model** of the documents

# Solution – VSM



## Vector Space Model (VSM)

# Solution – VSM

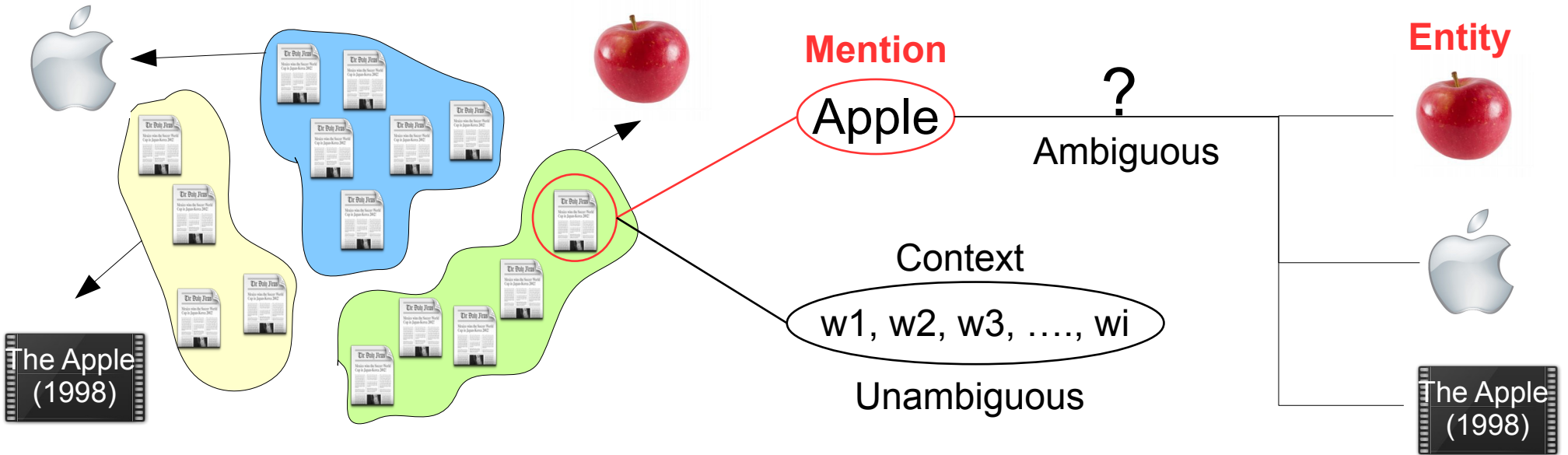


Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$



# Solution – VSM

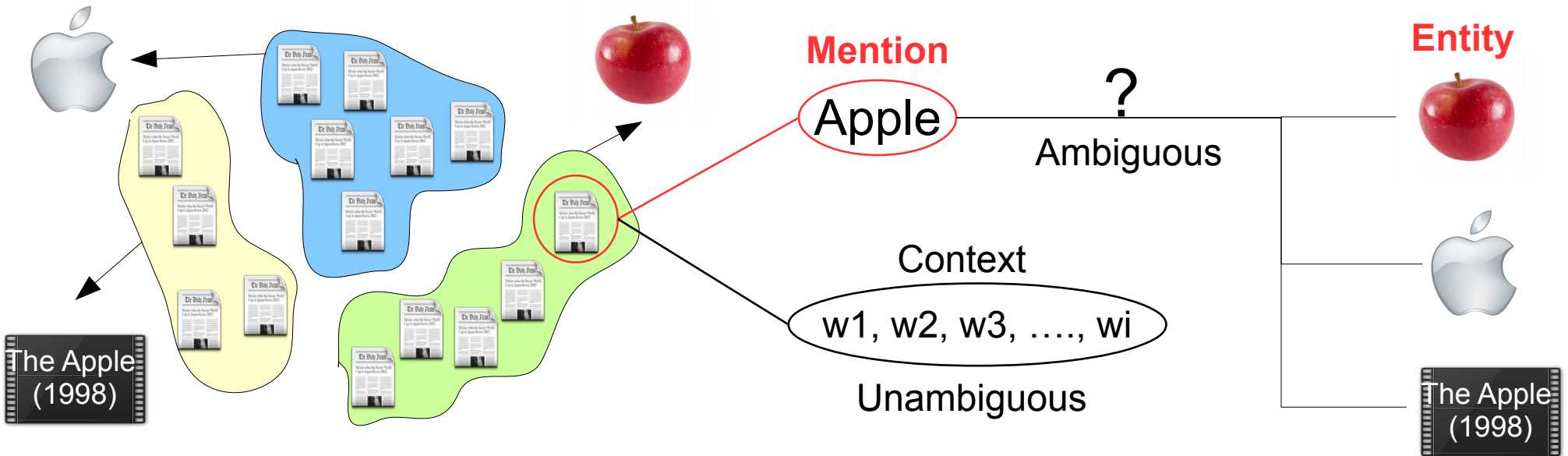


Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \rightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

# Solution – VSM



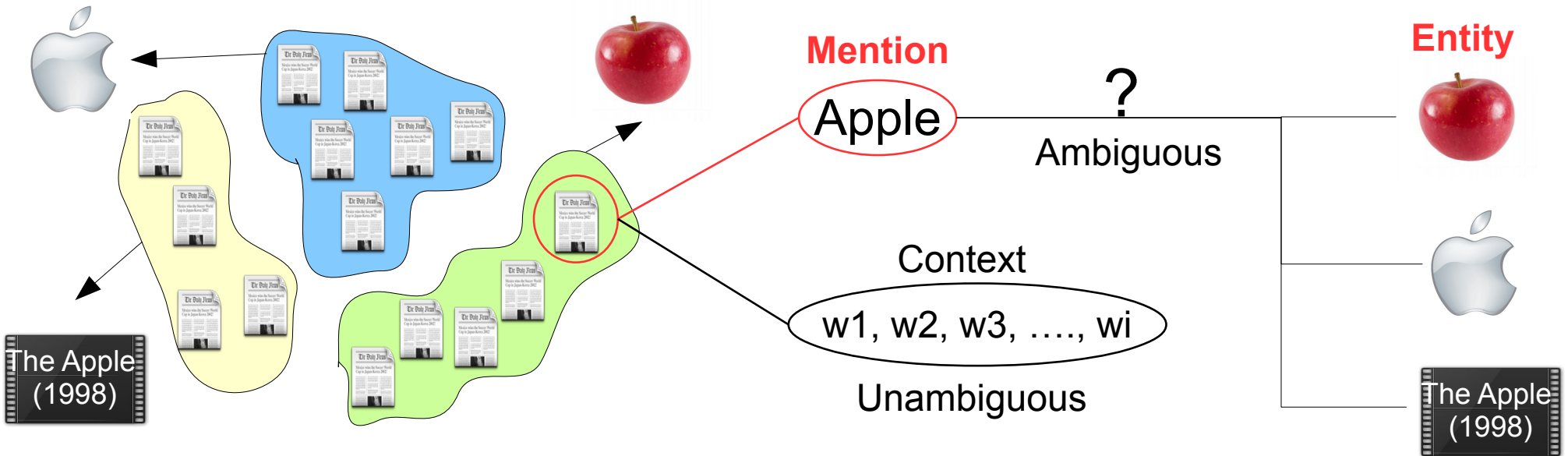
Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \rightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

$$D_2 \rightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

# Solution – VSM



Unique Context Words

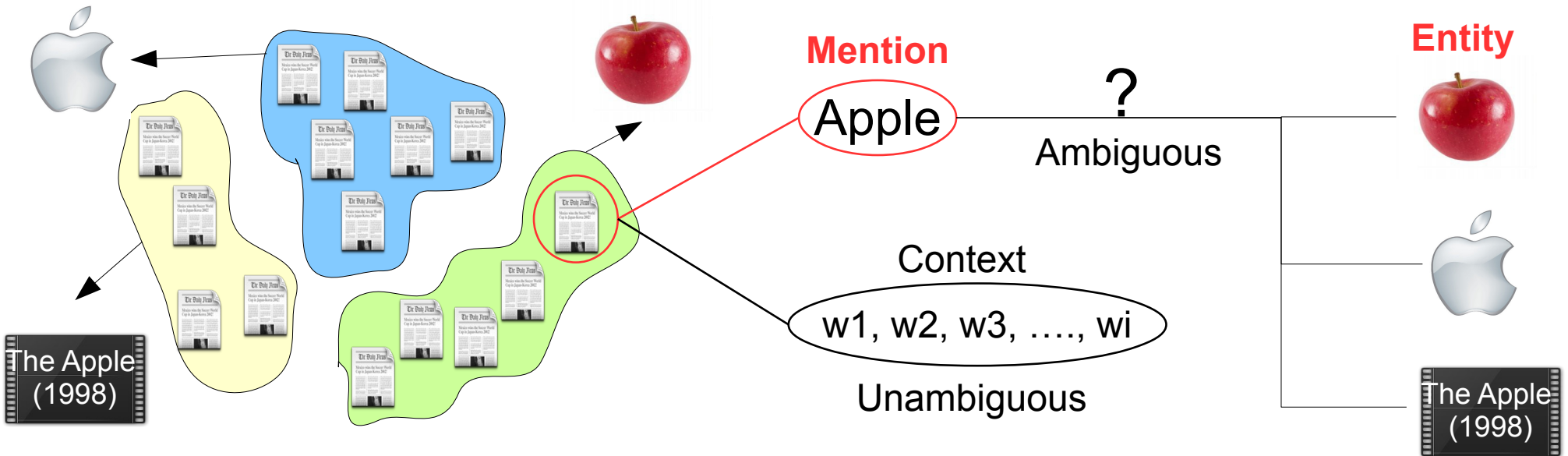
$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \longrightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

$$D_2 \longrightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

$$D_n \longrightarrow \vec{V}_n = \{1, 1, 1, \dots, 0\}$$

# Solution – VSM



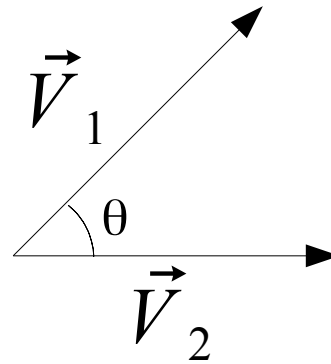
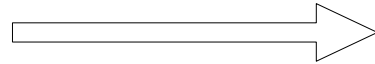
Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \longrightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

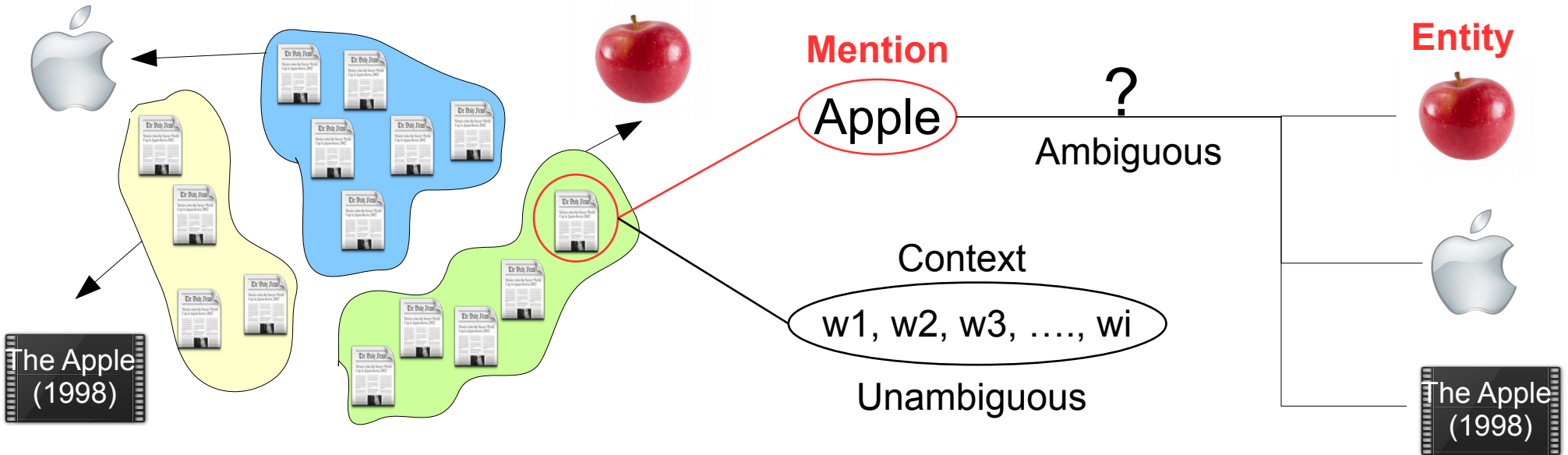
$$D_2 \longrightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

$$D_n \longrightarrow \vec{V}_n = \{1, 1, 1, \dots, 0\}$$



$$\text{Similarity}(D_1, D_2) \propto \cos(\theta)$$

# Solution – VSM



Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

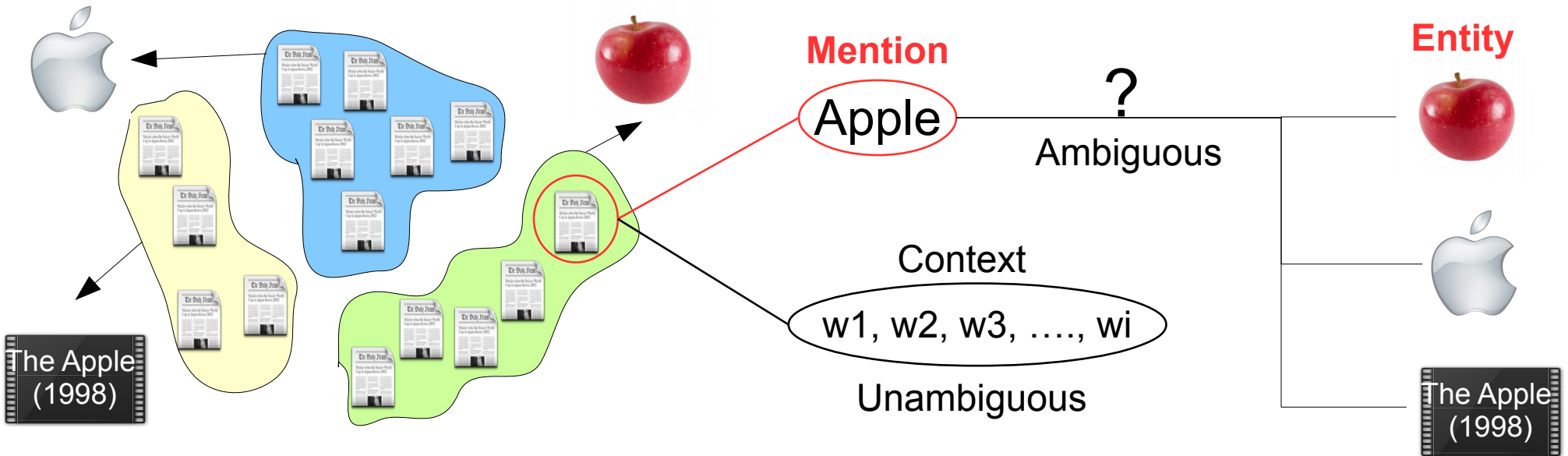
$$D_1 \rightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

$$D_2 \rightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

$$D_n \rightarrow \vec{V}_n = \{1, 1, 1, \dots, 0\}$$

## Limitations

# Solution – VSM



Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \longrightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

$$D_2 \longrightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

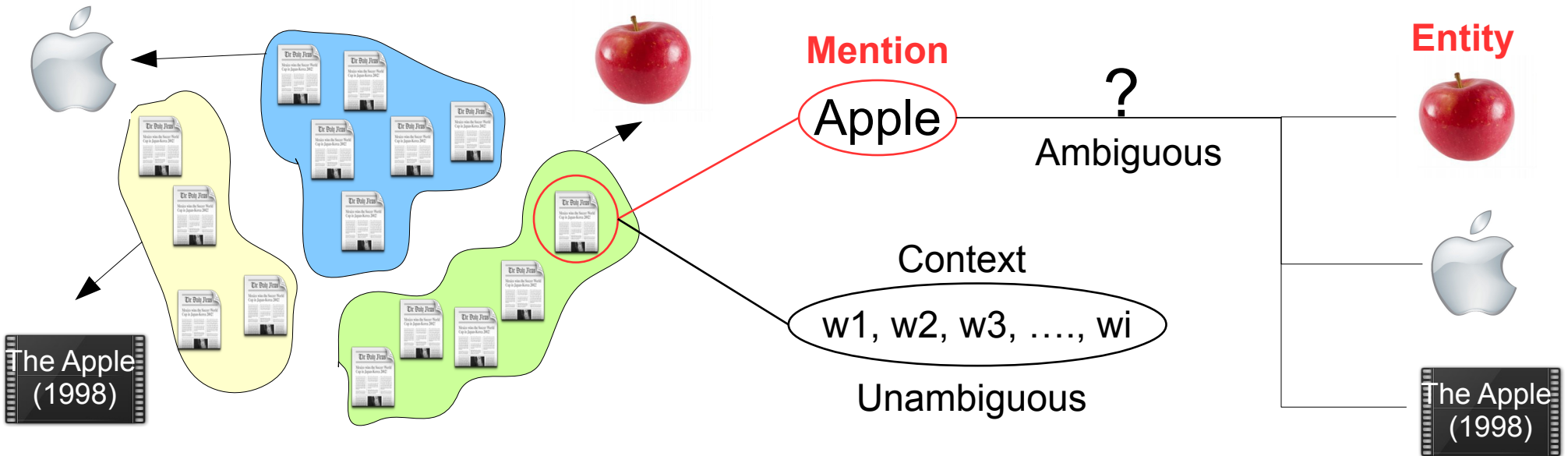
$$D_n \longrightarrow \vec{V}_n = \{1, 1, 1, \dots, 0\}$$

## Limitations

**Number** of Comparisons

Not every two documents need to be compared

# Solution – VSM



Unique Context Words

$$\{w_1, w_2, w_3, \dots, w_m\}$$

$$D_1 \longrightarrow \vec{V}_1 = \{1, 0, 1, \dots, 0\}$$

$$D_2 \longrightarrow \vec{V}_2 = \{0, 0, 1, \dots, 1\}$$

$$D_n \longrightarrow \vec{V}_n = \{1, 1, 1, \dots, 0\}$$

## Limitations

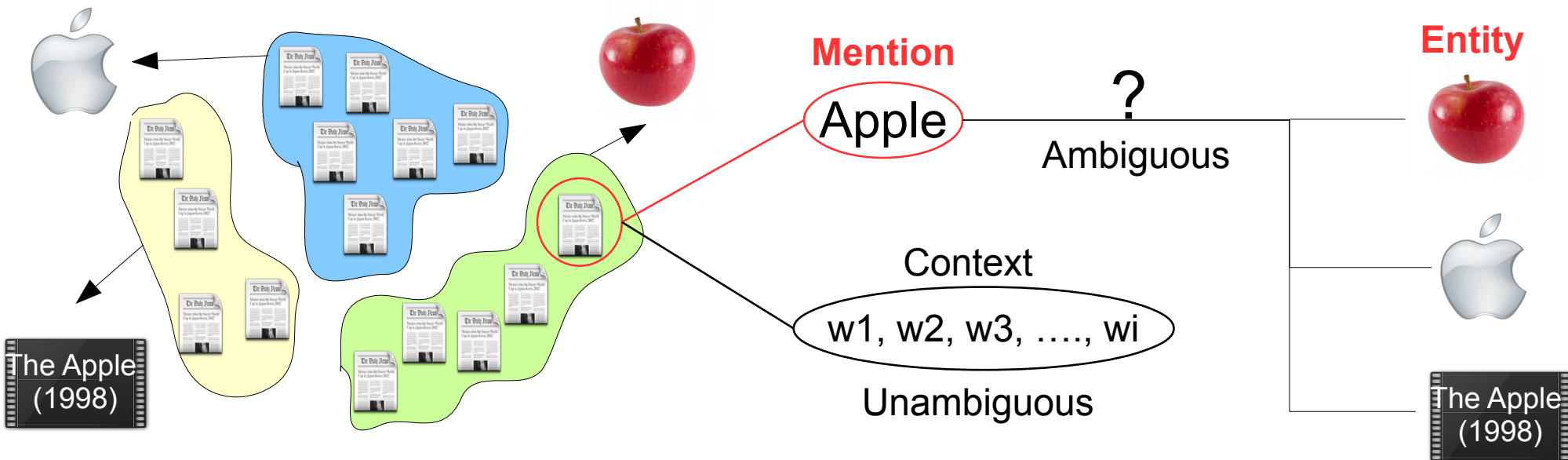
**Number** of Comparisons

Not every two documents need to be compared

**Size** of each Comparison

Total number of unique context words (n)

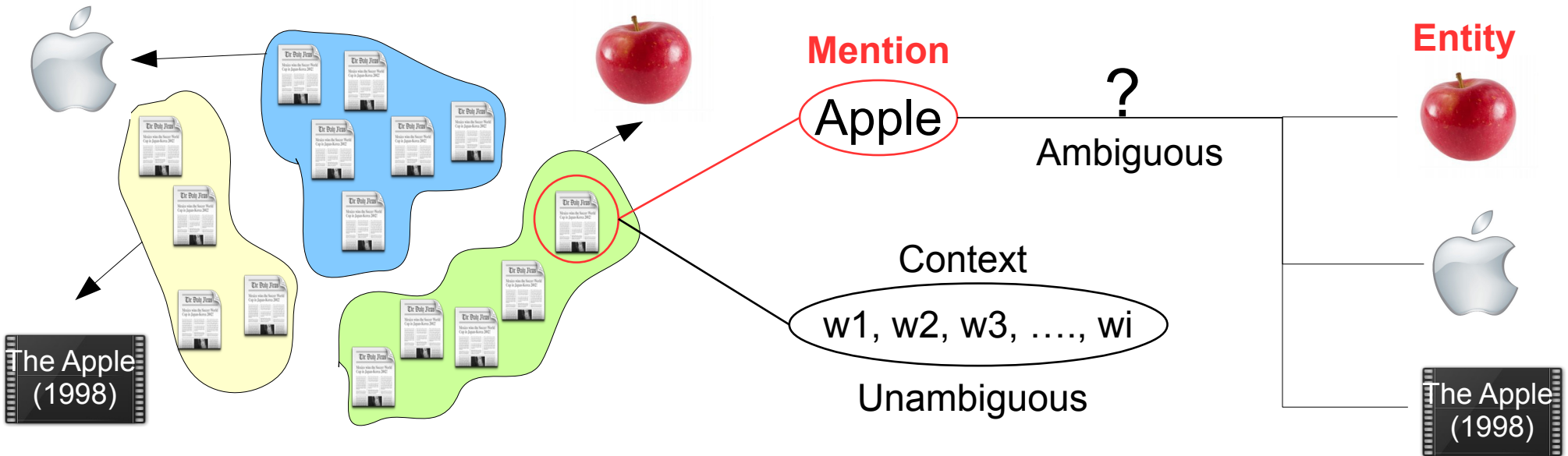
# Graph Based Modeling



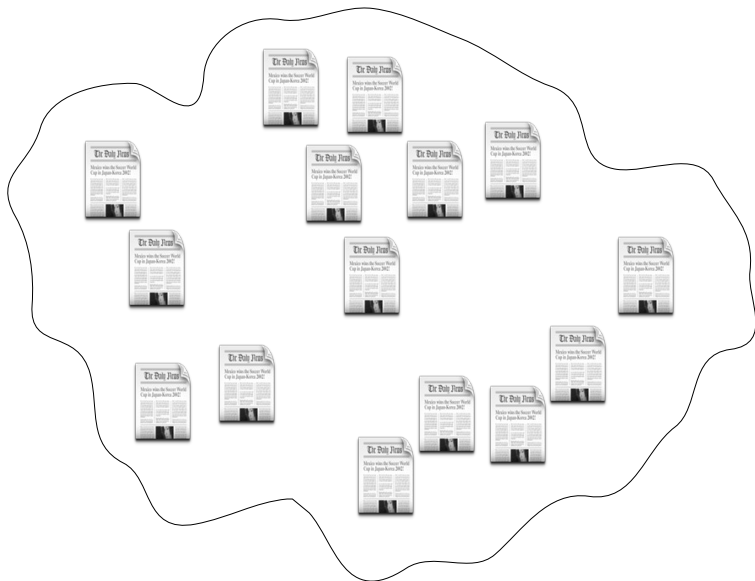
## Graph Based Modeling



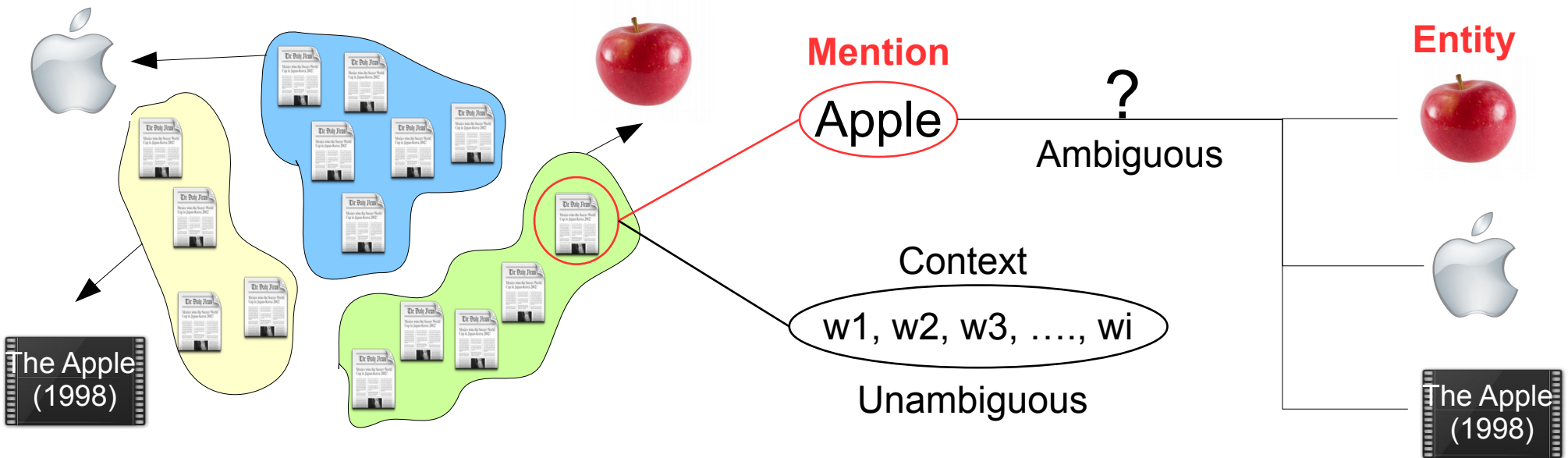
# Graph Based Modeling



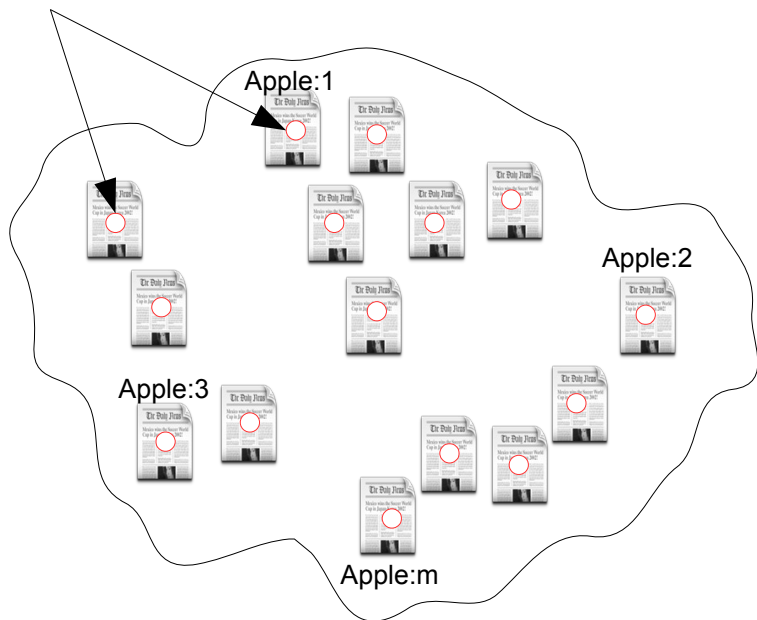
## Graph based modeling



# Graph Based Modeling



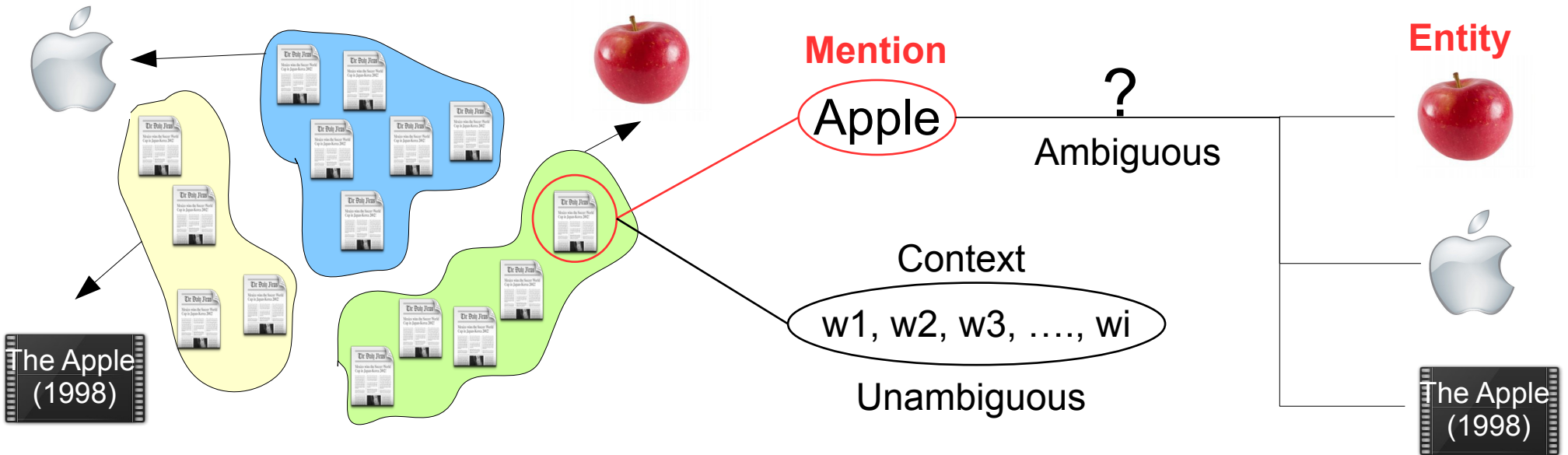
Mention



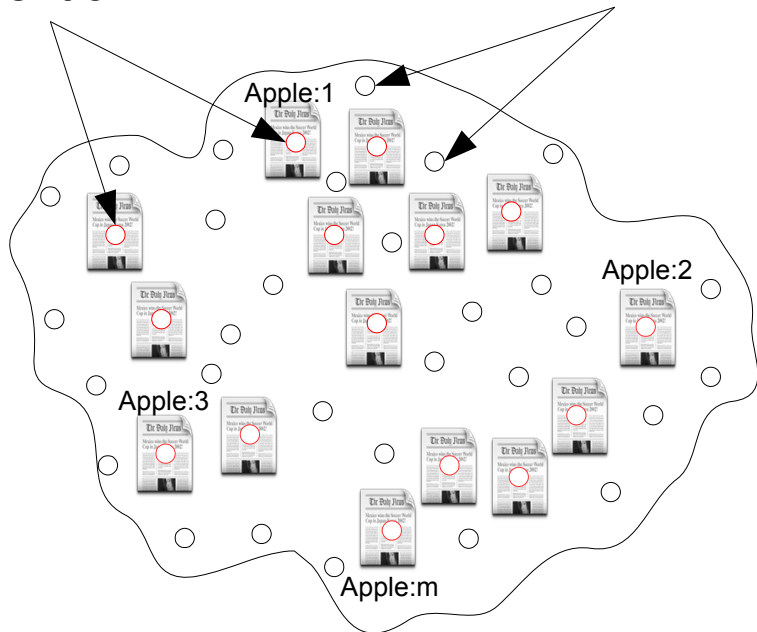
## Graph based modeling

- Assign a unique **vertex** to each **Mention**

# Graph Based Modeling



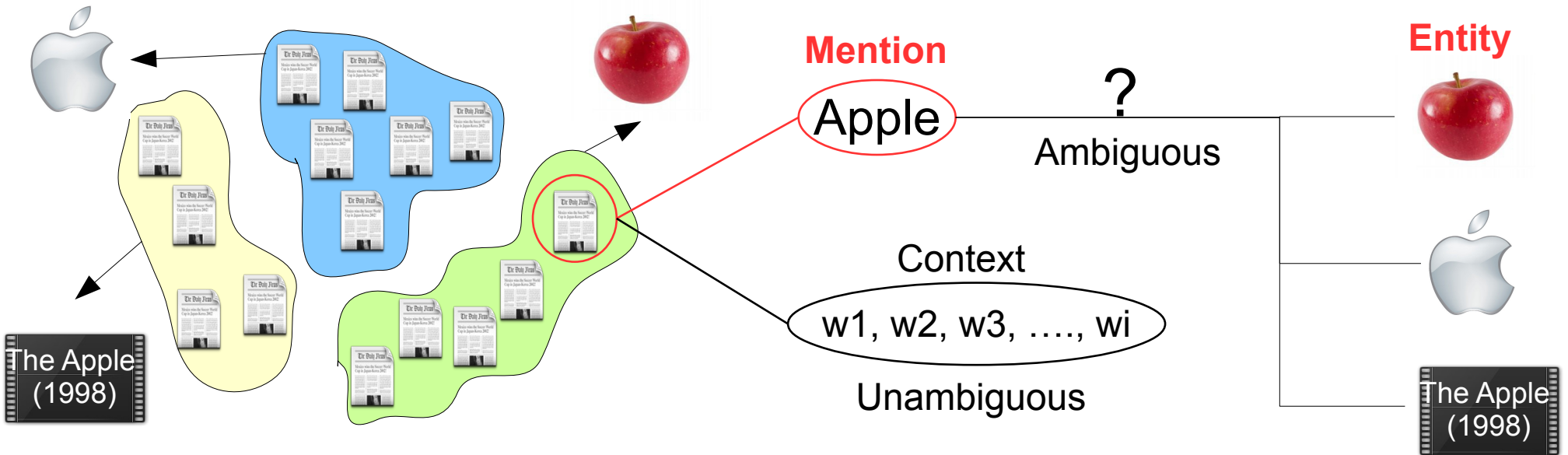
Mention Context



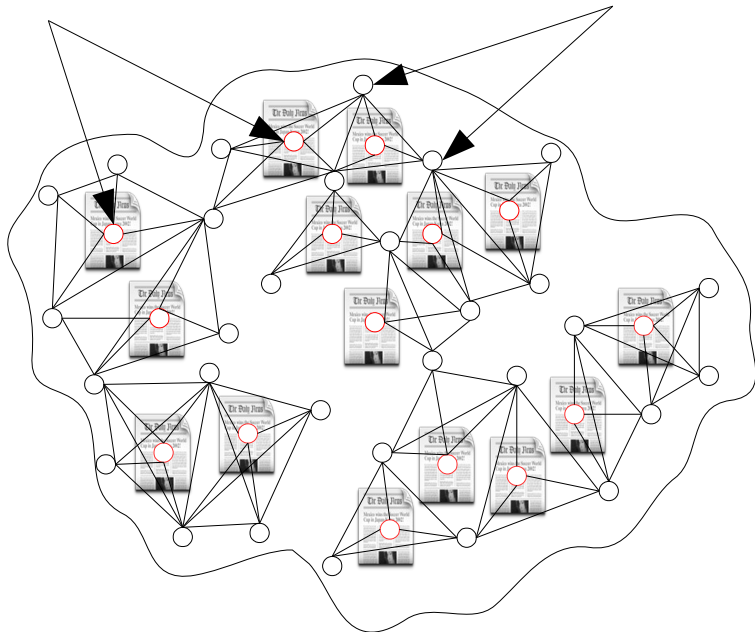
## Graph based modeling

- Assign a unique **vertex** to each **Mention**
- Assign a unique **vertex** to each **context word**

# Graph Based Modeling



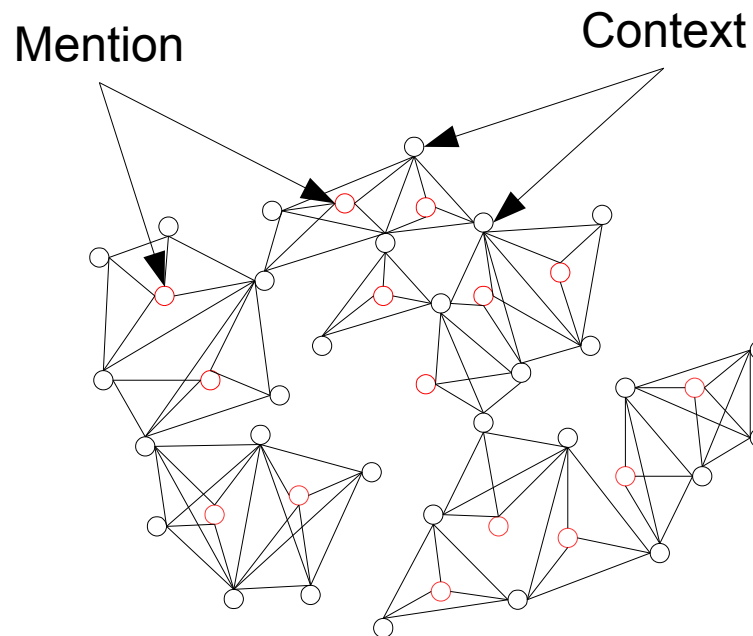
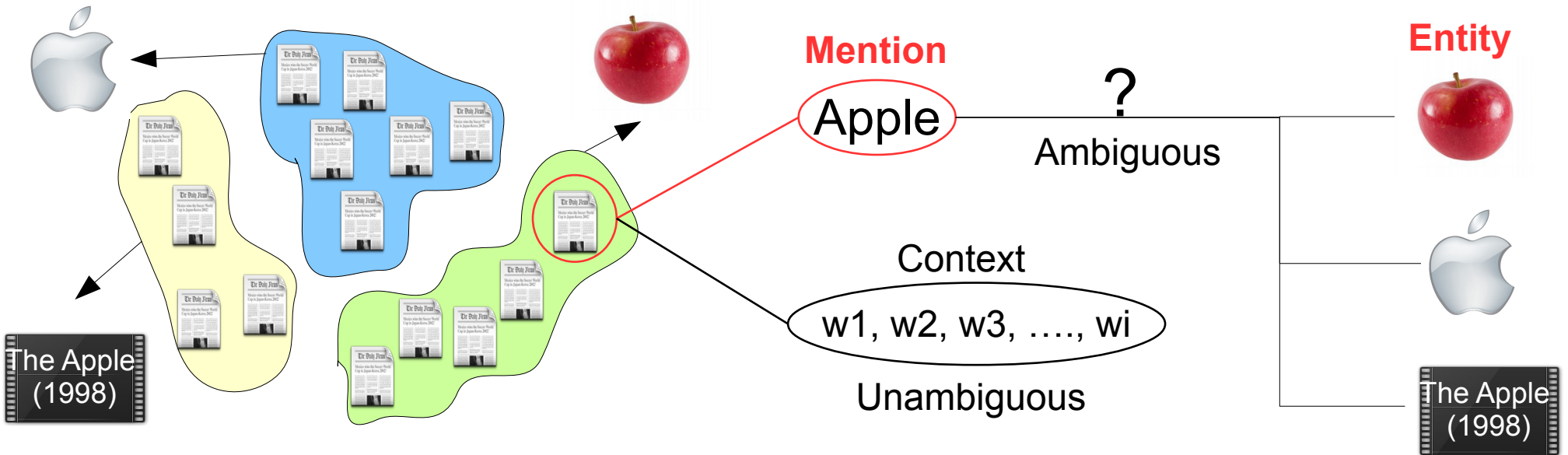
Mention Context



## Graph based modeling

- Assign a unique vertex to each **Mention**
- Assign a unique vertex to each **context word**
- Connect vertices from **same** document

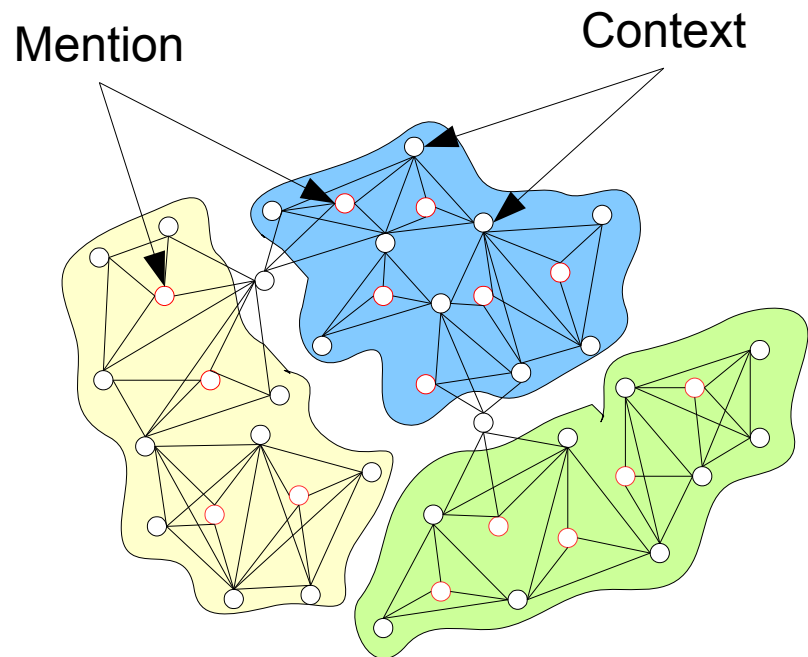
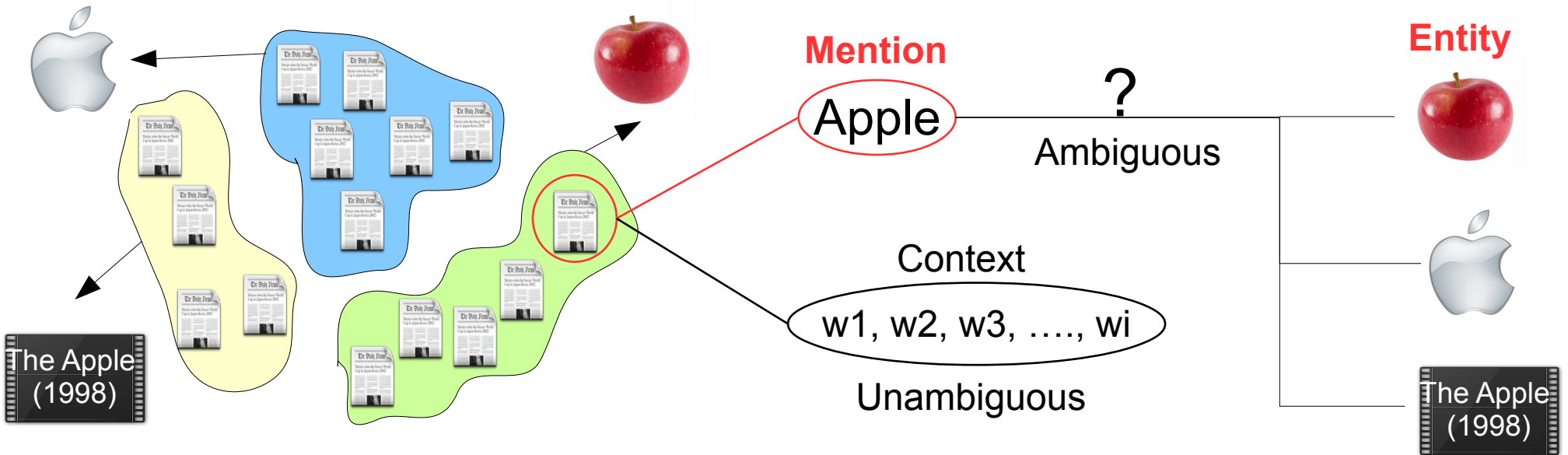
# Graph Based Modeling



## Graph based modeling

- Assign a unique vertex to each **Mention**
- Assign a unique vertex to each **context word**
- Connect vertices from **same** document
- Graph Representation

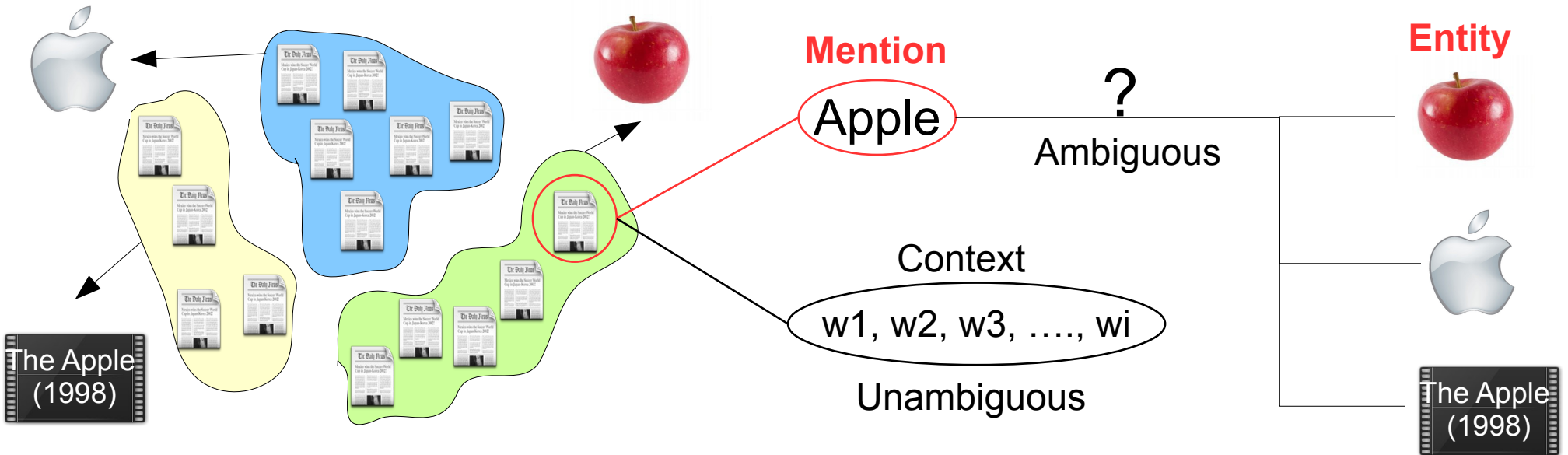
# Graph Based Modeling



## Graph based modeling

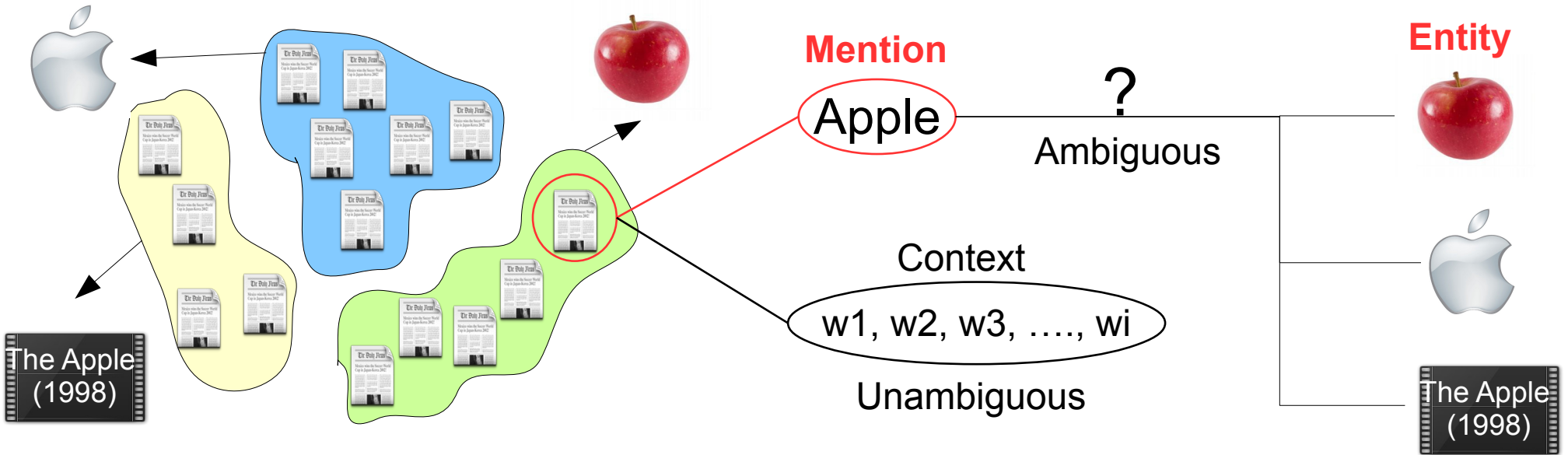
- Assign a unique vertex to each **Mention**
- Assign a unique vertex to each **context word**
- Connect vertices from **same** document
- Graph Representation

# Diffusion based Clustering

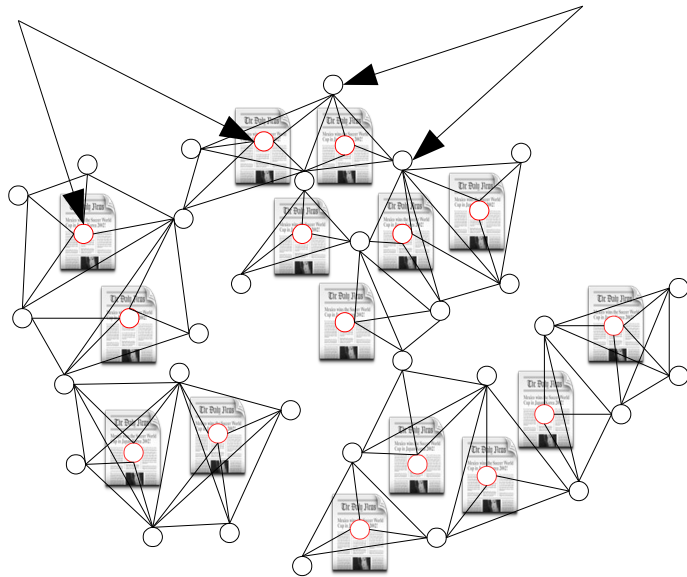


## Diffusion based Clustering

# Diffusion based Clustering



Mention                      Context

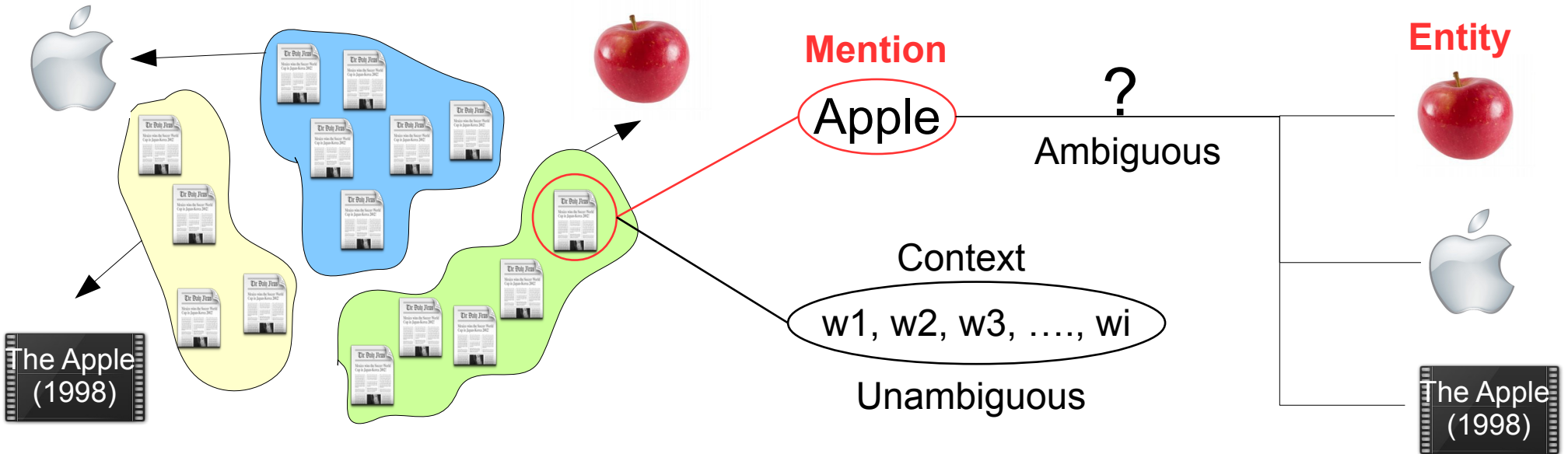


Diffusion Clustering

- Initialization
- Diffusion



# Diffusion based Clustering

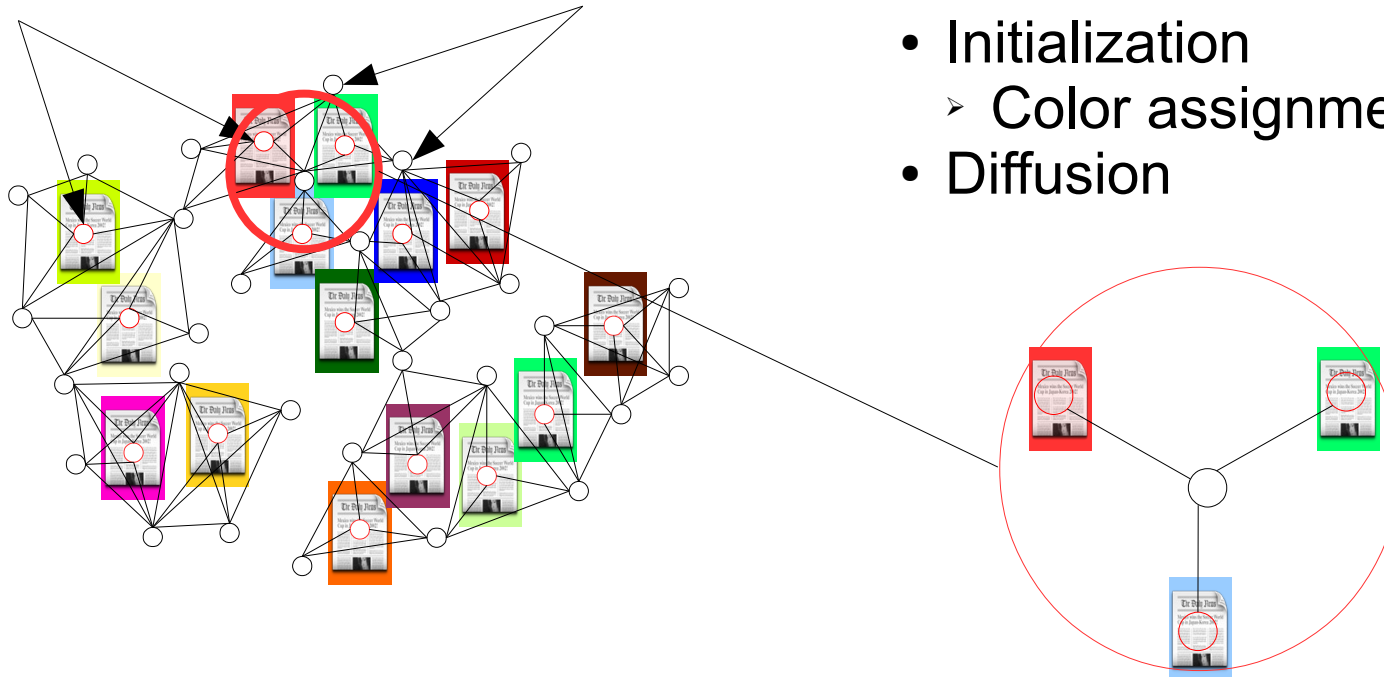


Mention

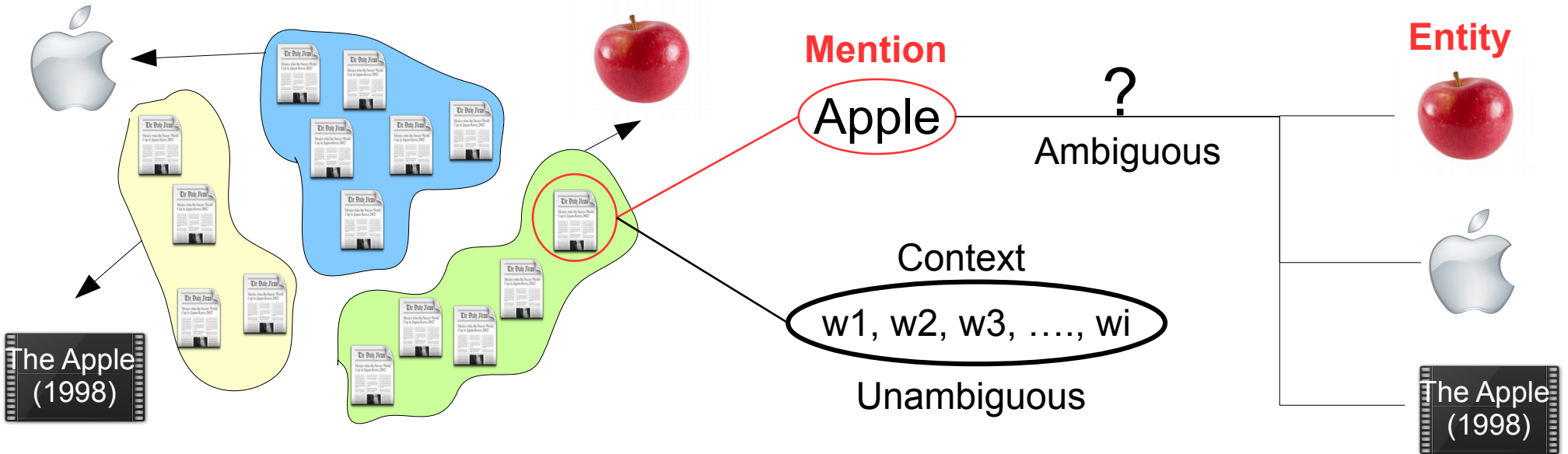
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion



# Diffusion based Clustering

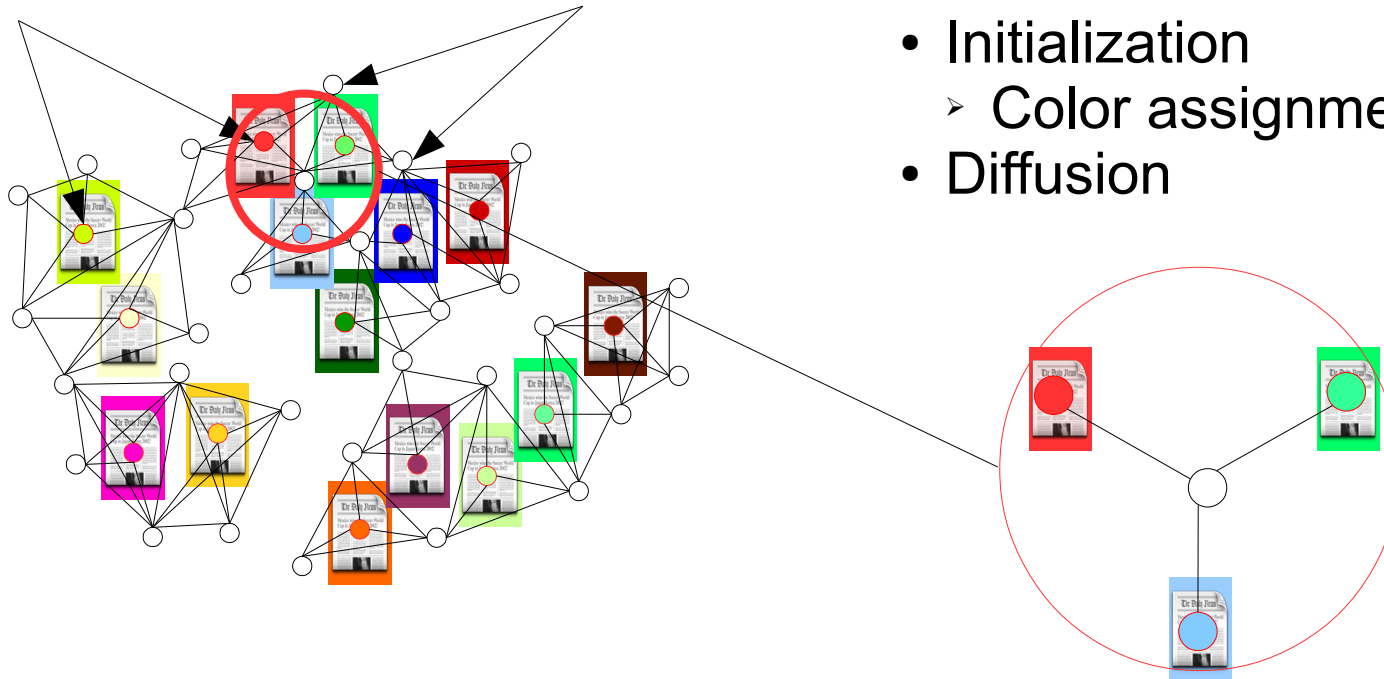


Mention

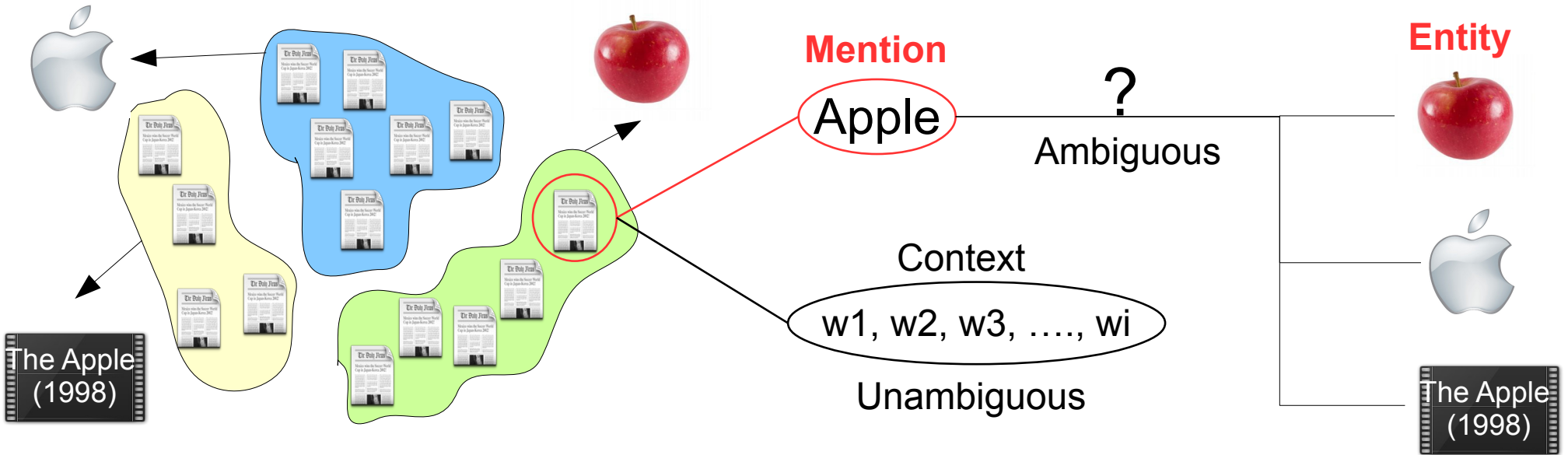
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion



# Diffusion based Clustering

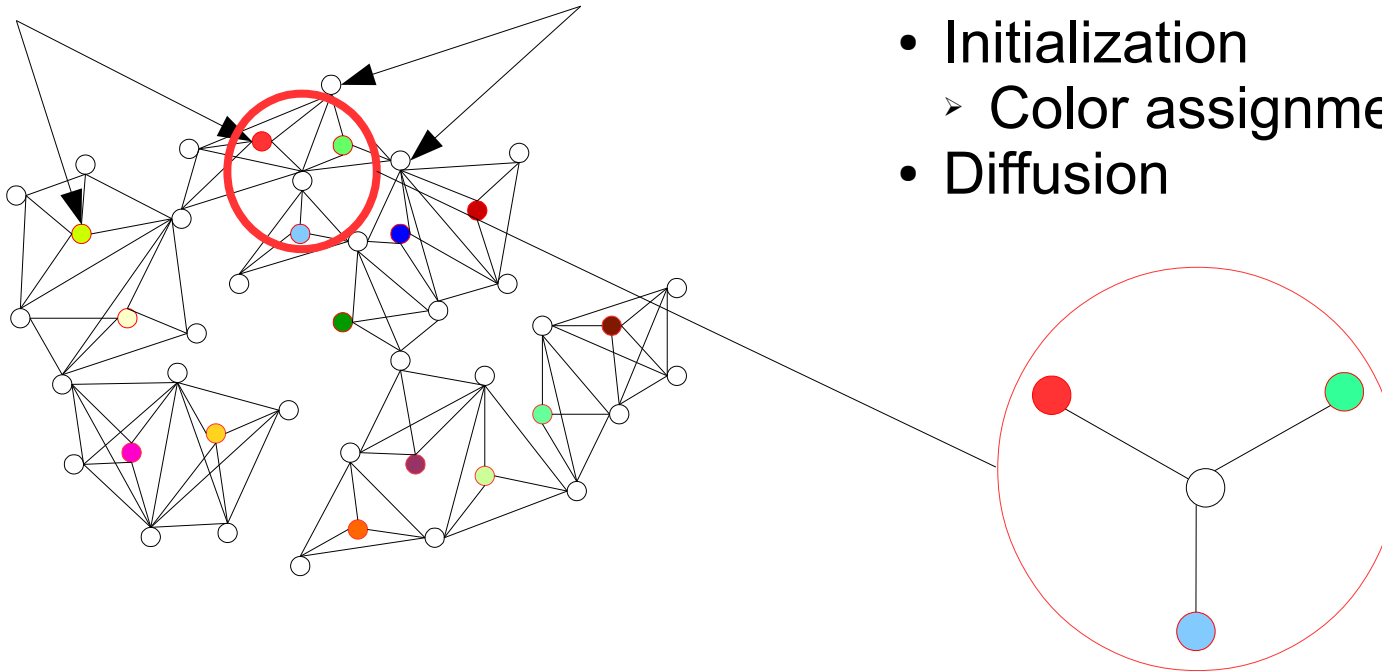


Mention

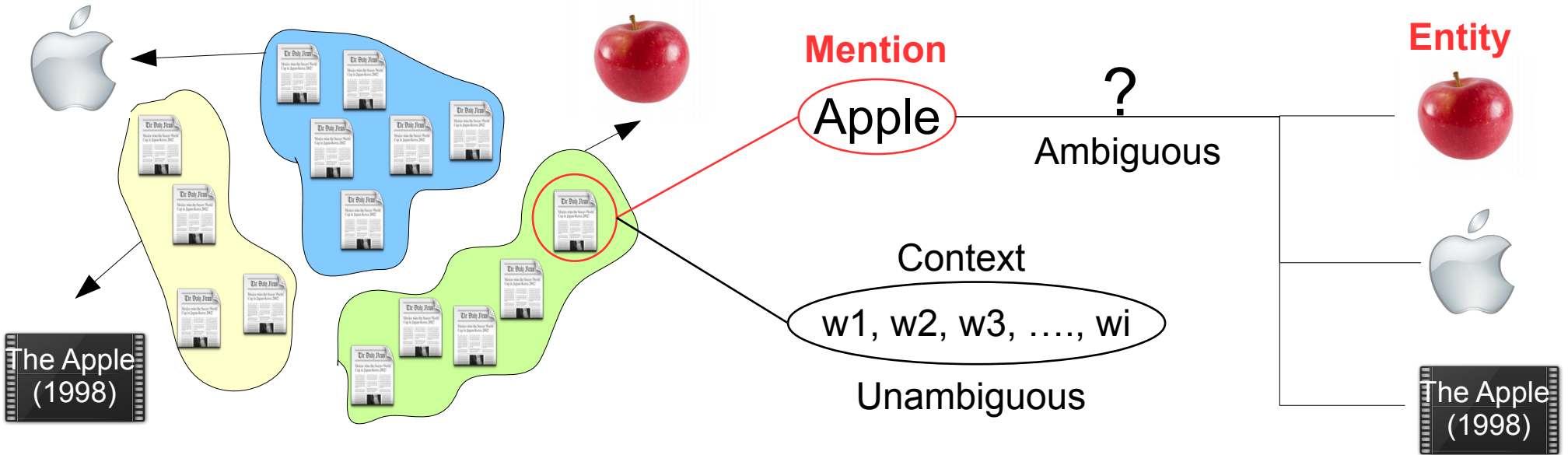
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion



# Diffusion based Clustering

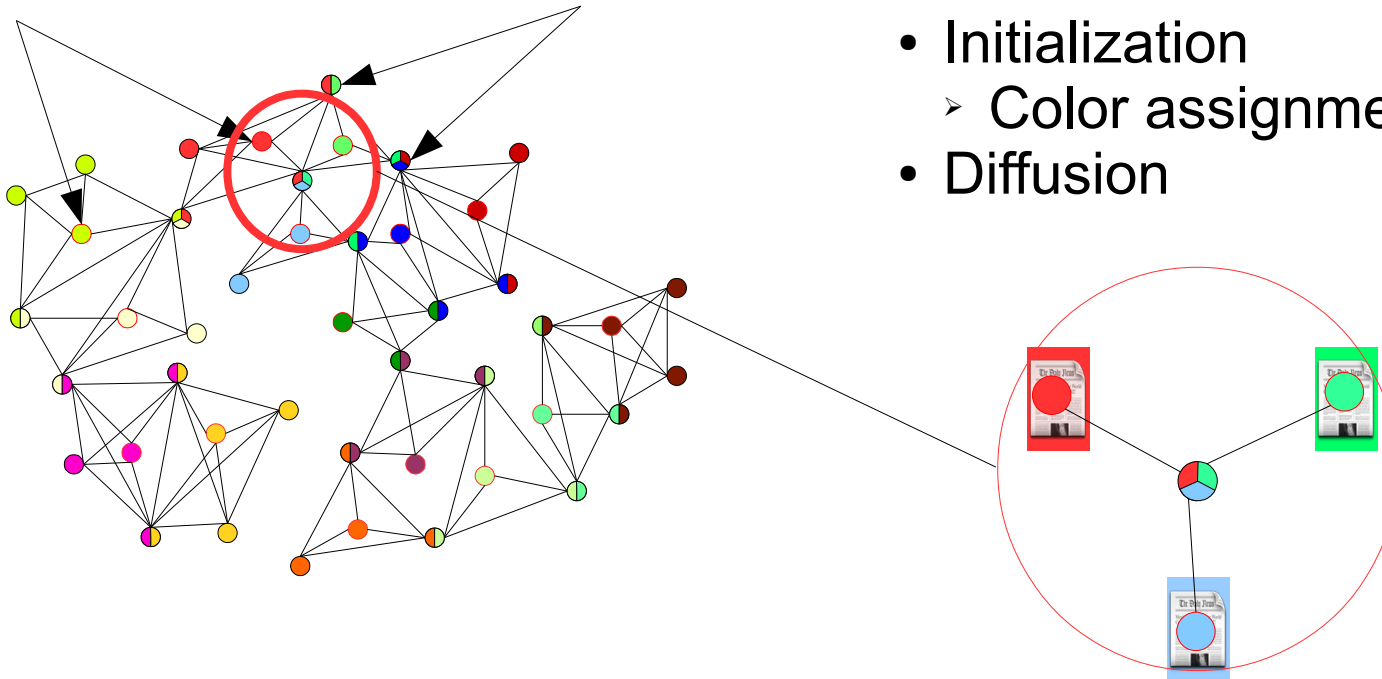


Mention

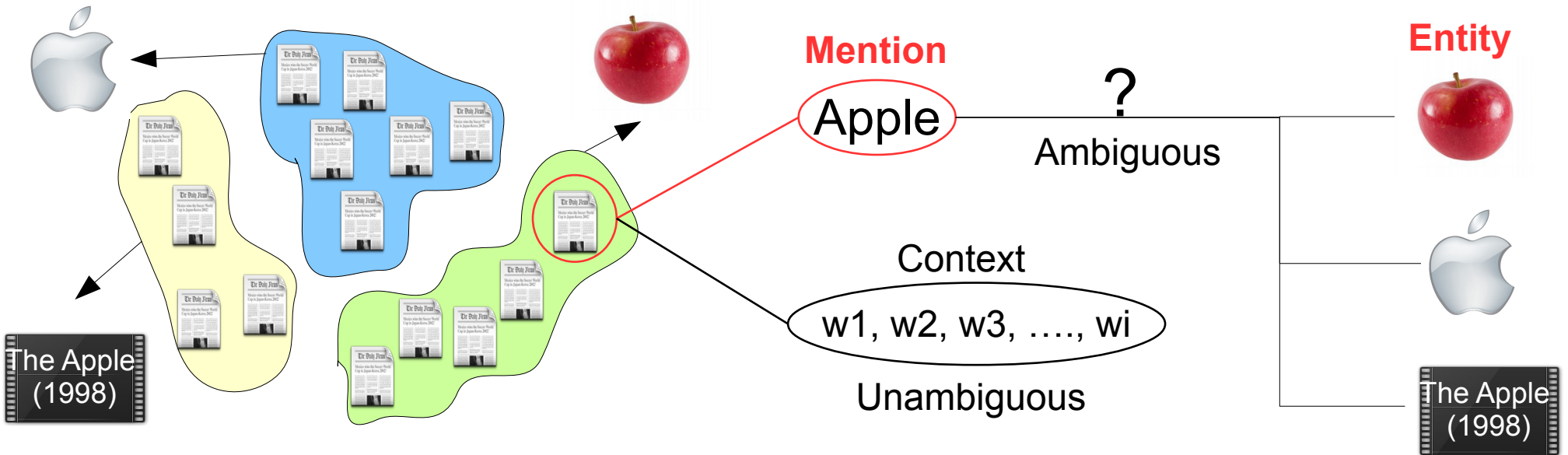
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion



# Diffusion based Clustering

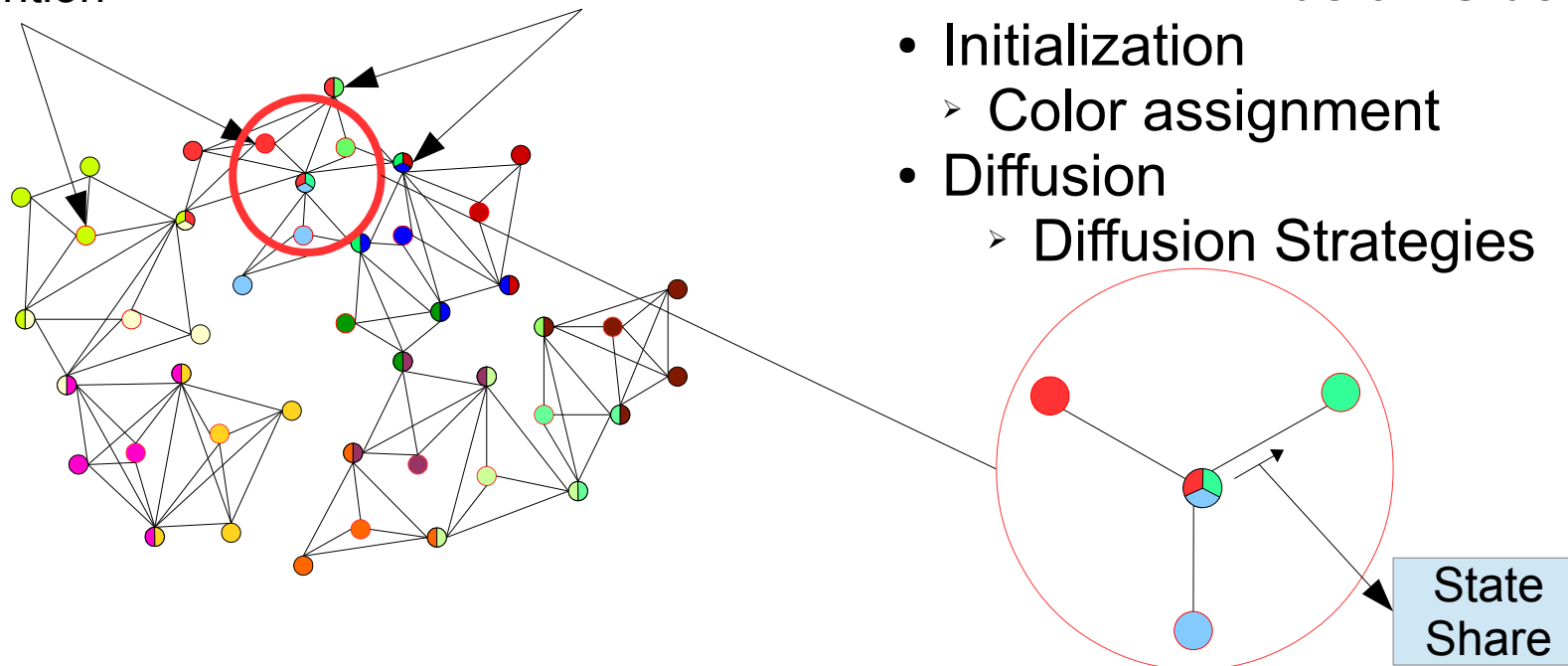


Mention

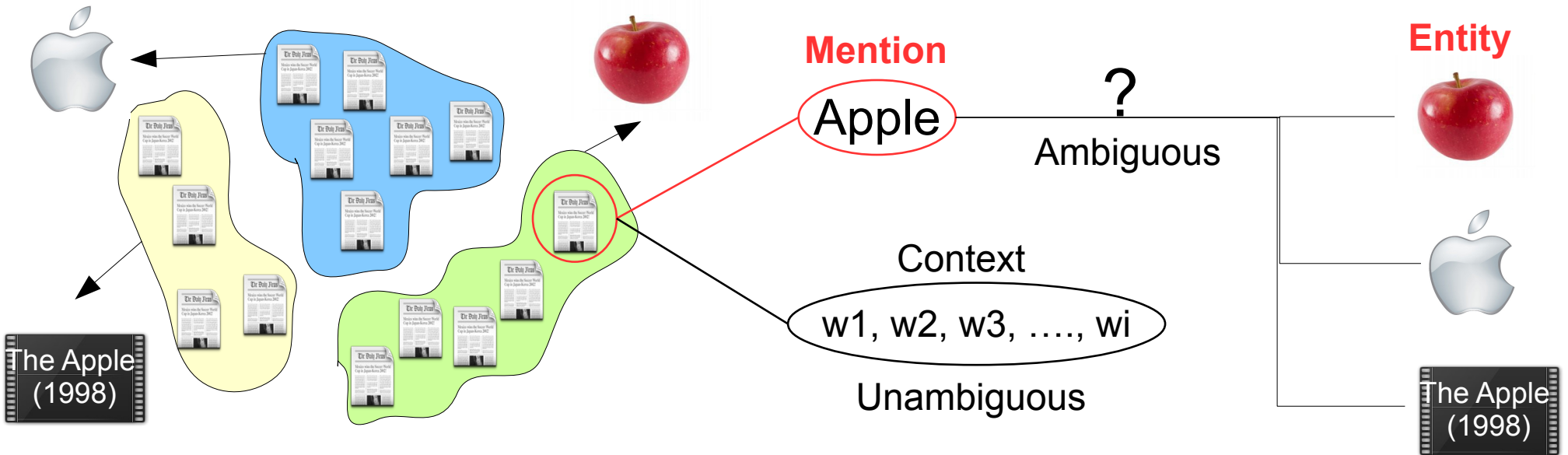
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion
  - Diffusion Strategies



# Diffusion based Clustering

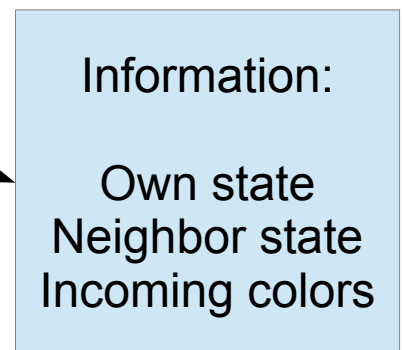
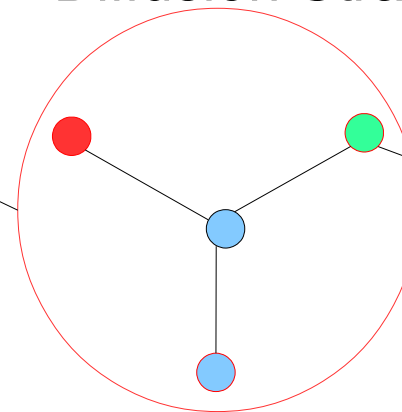
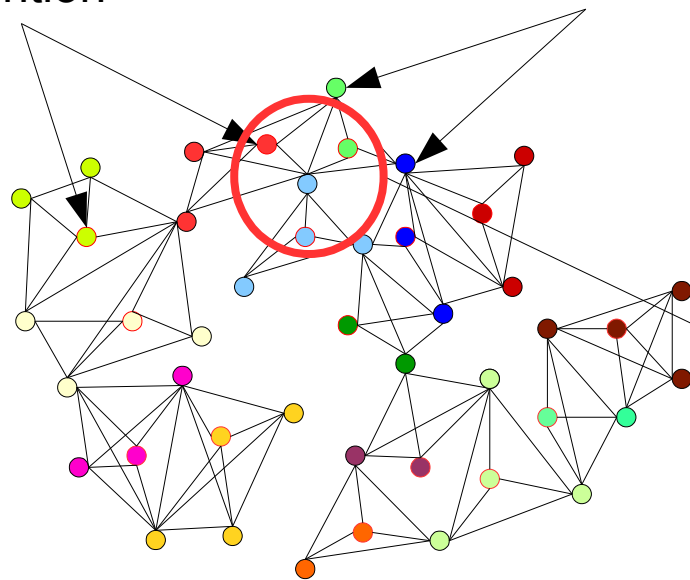


Mention

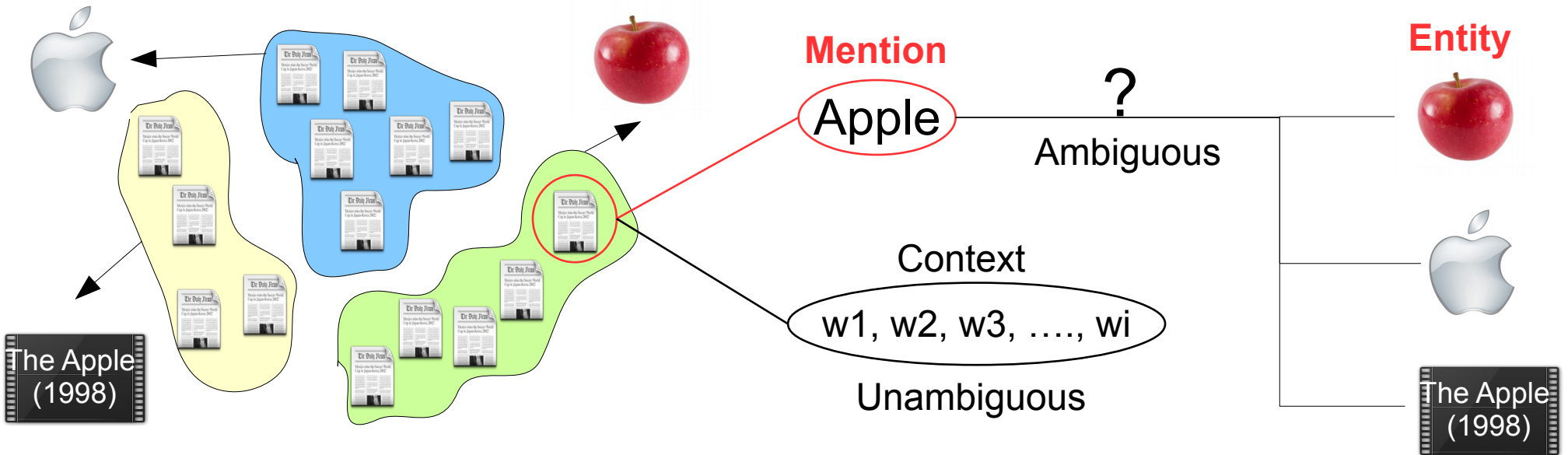
Context

Diffusion Clustering

- Initialization
  - Color assignment
- Diffusion
  - Diffusion Strategies



# Diffusion based Clustering

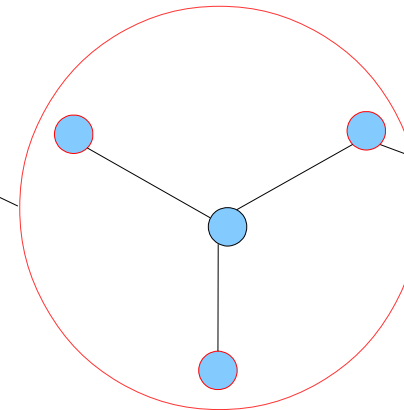
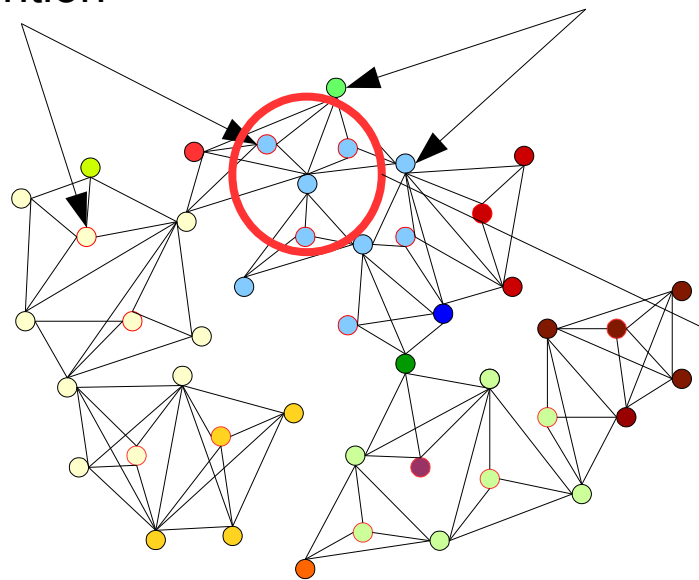


Mention

Context

Diffusion Clustering

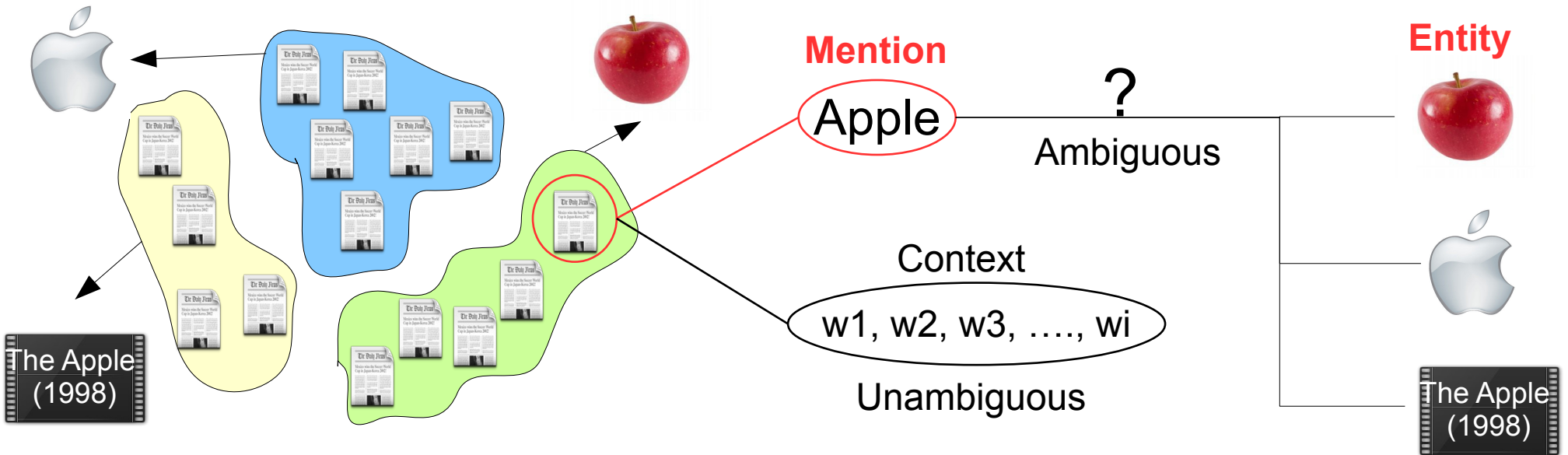
- Initialization
  - Color assignment
  - Diffusion
    - Diffusion Strategies



Information:

Own state  
Neighbor state  
Incoming colors

# Diffusion based Clustering

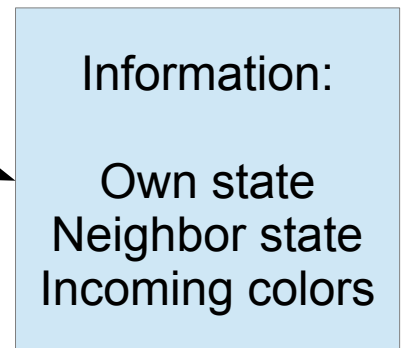
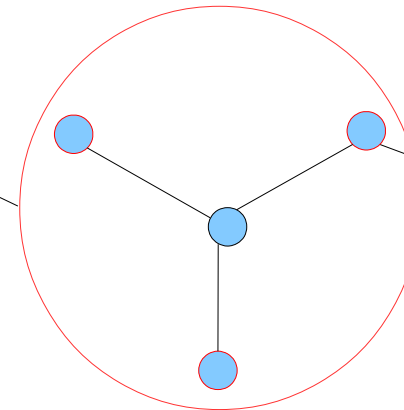
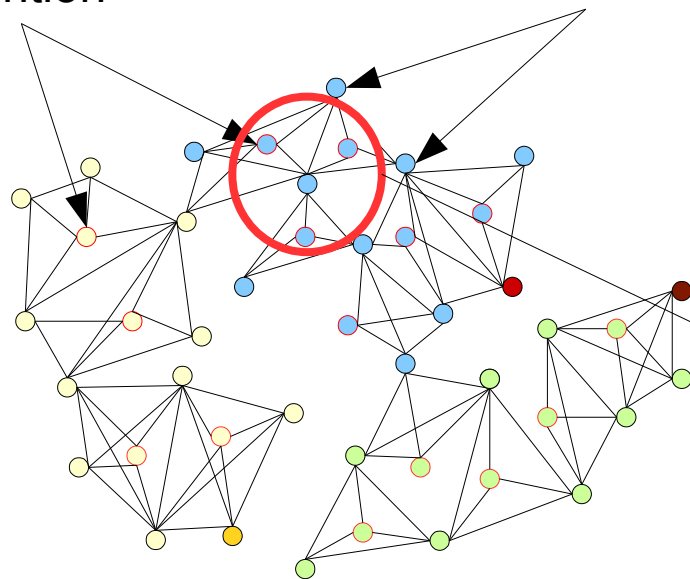


Mention

Context

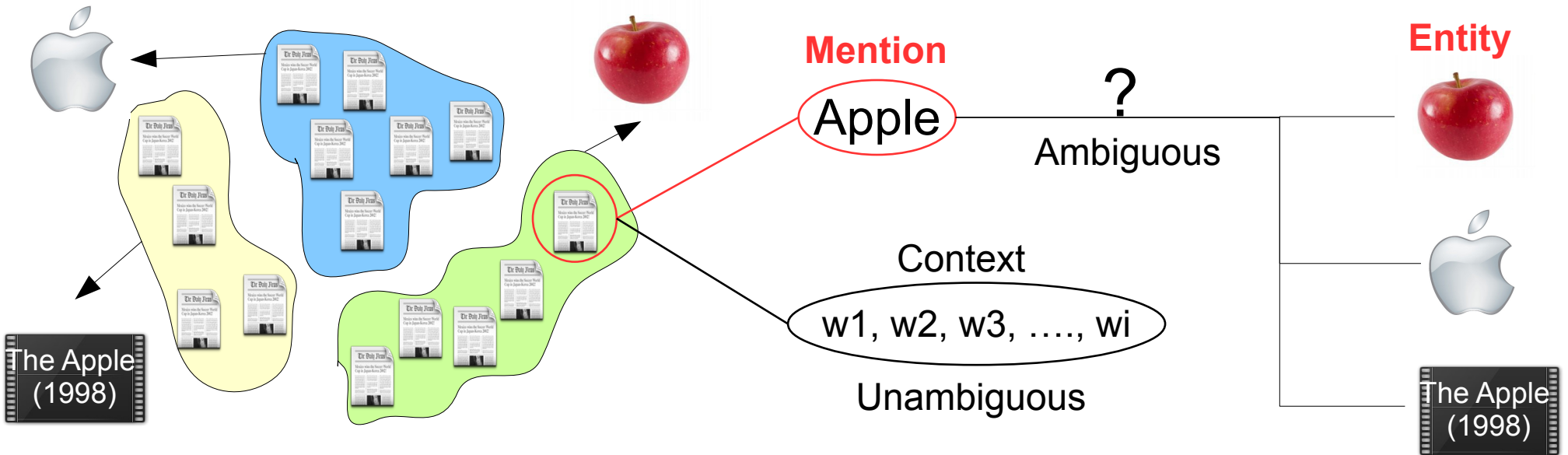
Diffusion Clustering

- Initialization
  - Color assignment
  - Diffusion
    - Diffusion Strategies





# Diffusion based Clustering

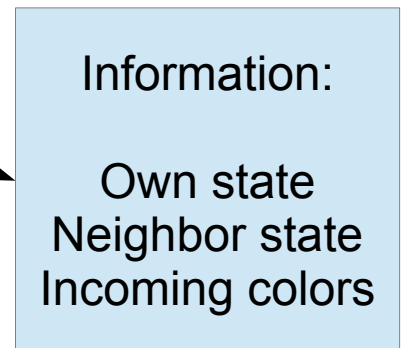
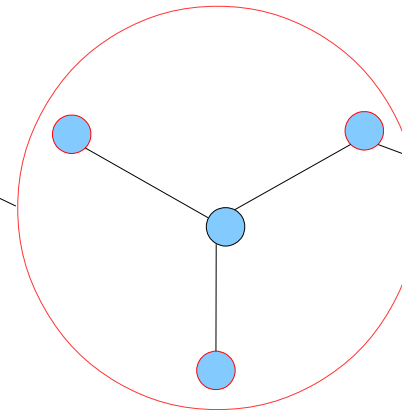
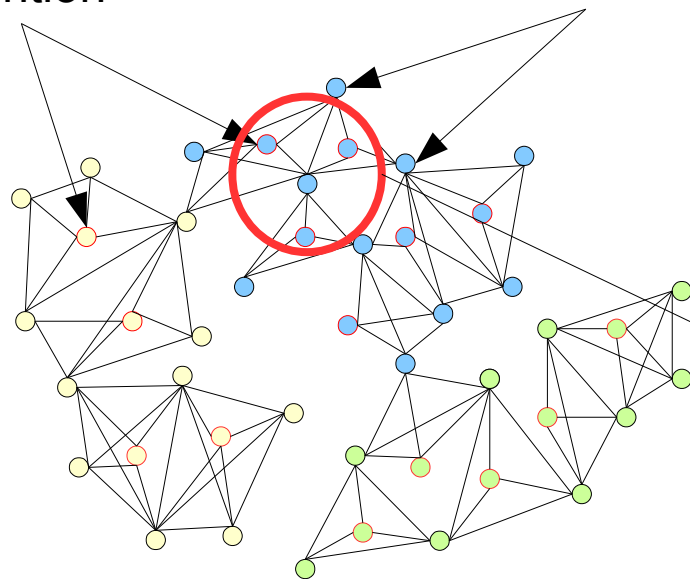


Mention

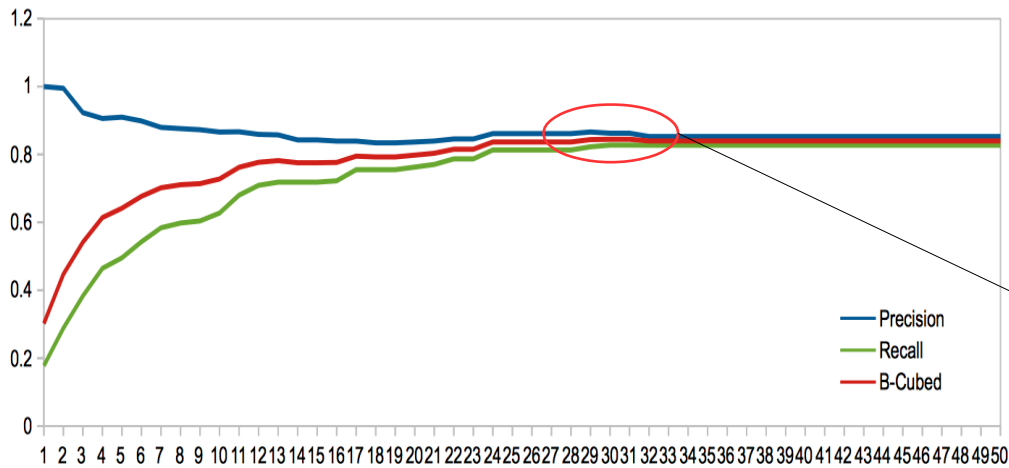
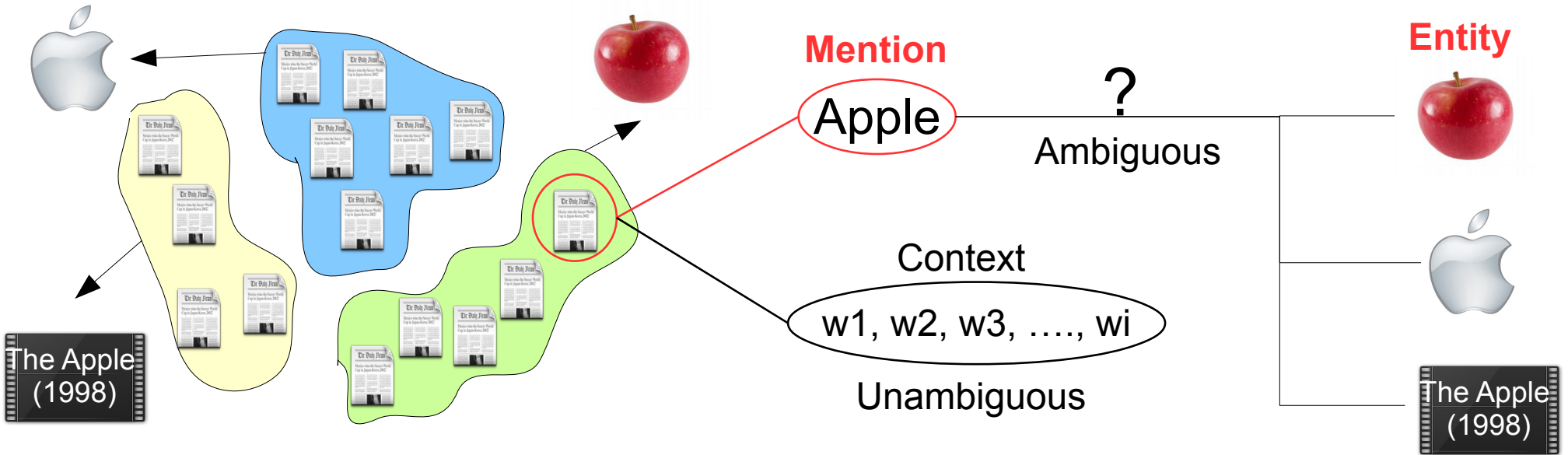
Context

Diffusion Clustering

- Initialization
  - Color assignment
  - Diffusion
    - Diffusion Strategies



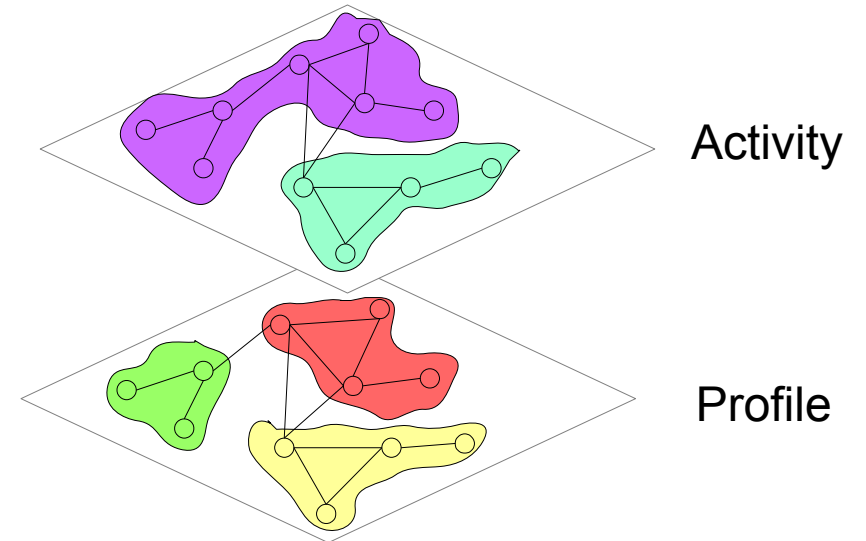
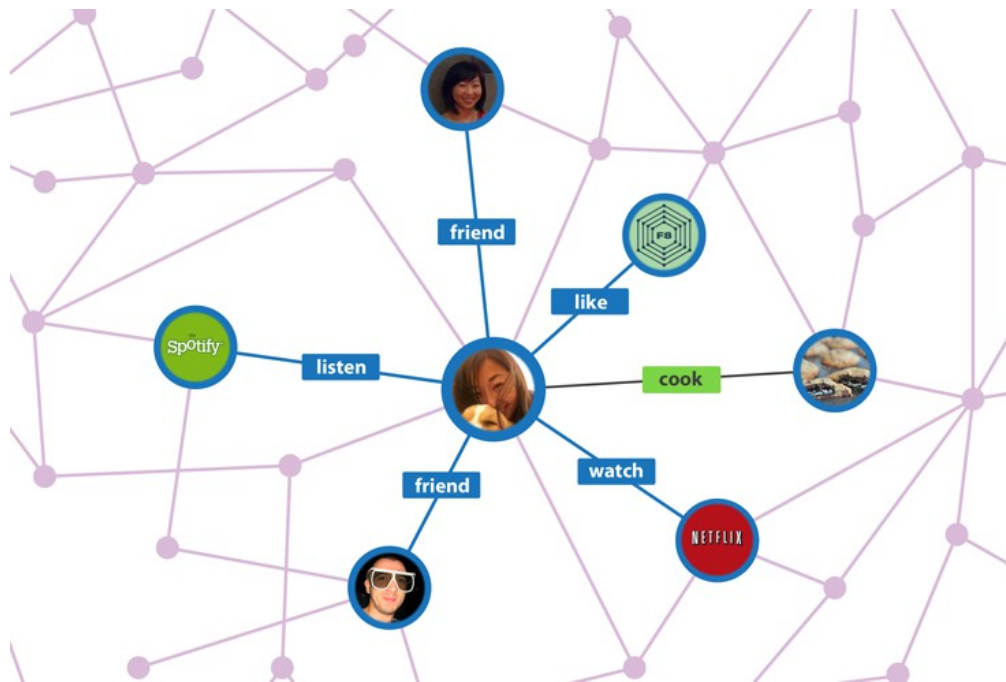
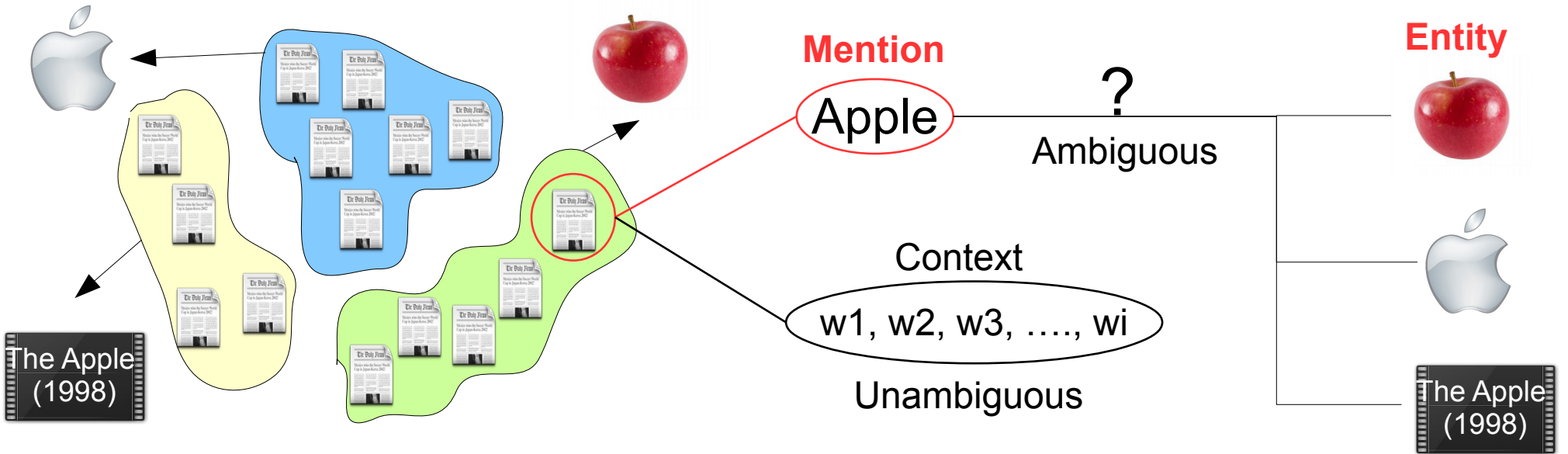
# Results



$\alpha = 0, \beta = 0.7, \gamma = 0.96, \delta = 0.05$

Model	Accuracy(F1%)
Bagga & Baldwin	84.6
Google	66.4
Basic	84.5
Extended	<b>88.9</b>

# What Next?



Thank You!