

# Sparse Representation Method for Whole Graph Embedding

Oral Comprehensive Exam  
Kaveen Liyanage

Dr. Bradley Whitaker  
Dr. Rob Maher

# Graphs are helpful

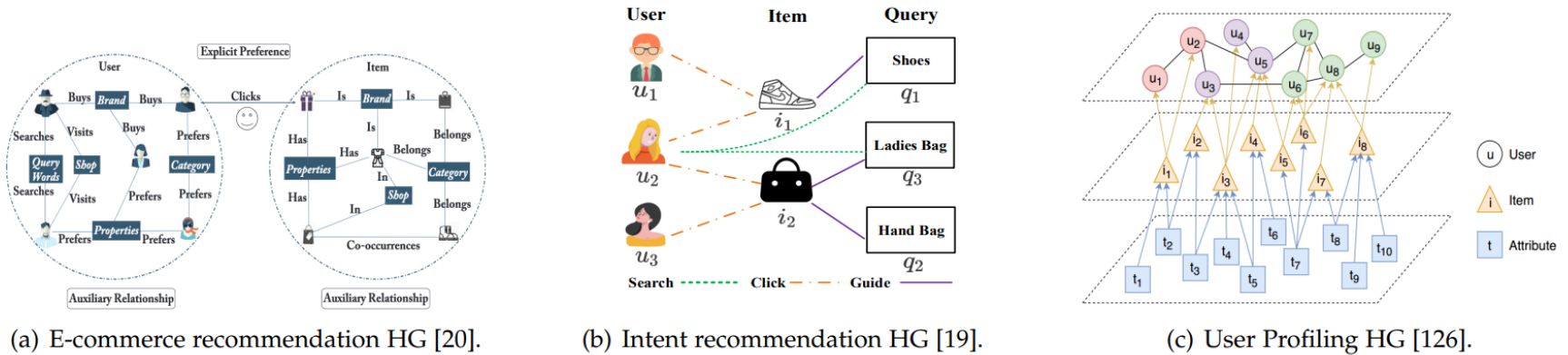


Fig. 3: The representative HGs in E-commerce.

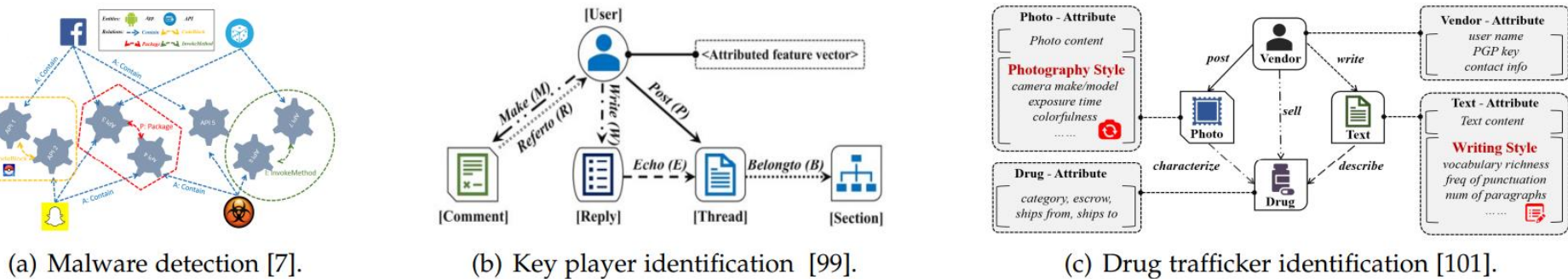


Fig. 4: The representative HGs in cybersecurity applications.

# Overview

## Introduction

- Motivation / Problem

- Goal

- Objectives

## Background

- Graph embedding

- Sparse representation

- Feature ranking

## Workflow

- Preliminary work

- Proposed timeline

- Future directions

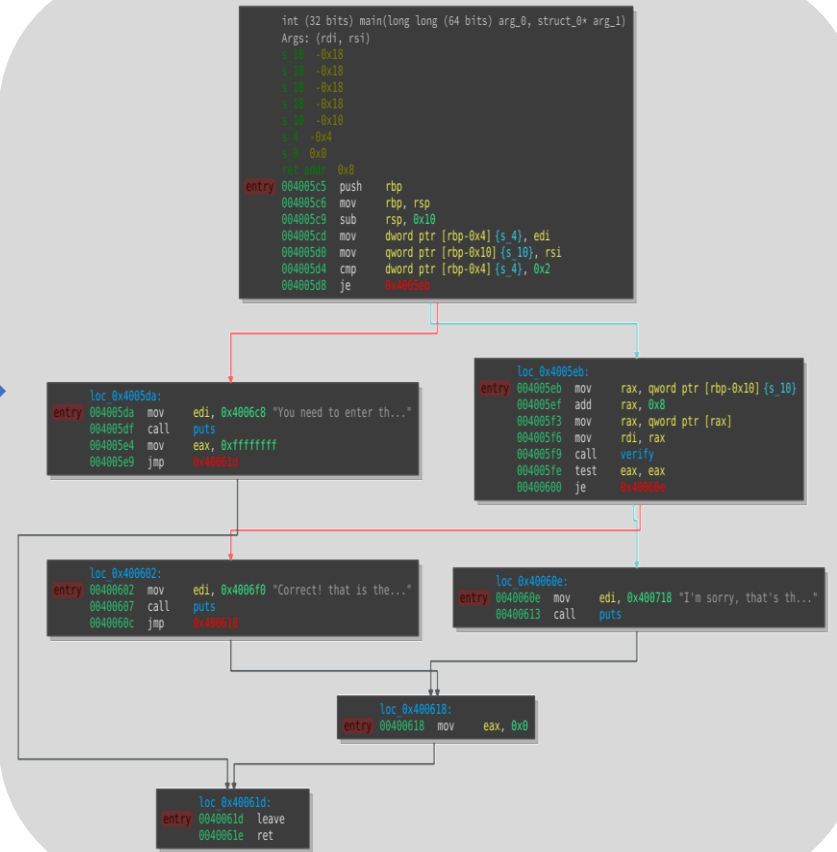
# Malware detection by Control Flow Graphs (CFG) of binary files

```

402000 endbr64
402004 sub    rsp, 0x8
402008 mov    rax, qword ptr [0x409fd0]
40200f test   rax, rax
402012 je     0x402016
402016 add    rsp, 0x8
40201a ret
402014 call   rax
402020 push    qword ptr [0x409e70]
402026 bnd jmp  qword ptr [0x409e78]
40202d nop     dword ptr [rax]
402030 endbr64
402034 push    0x0
402039 bnd jmp  sub_402020
40203f nop
402040 endbr64
402044 push    0x1
402049 bnd jmp  sub_402020
40204f nop
402050 endbr64
402054 push    0x2
402059 bnd jmp  sub_402020
40205f nop
402060 endbr64
402064 push    0x3
402069 bnd jmp  sub_402020
40206f nop
402070 endbr64
402074 push    0x4
402079 bnd jmp  sub_402020
40207f nop
402080 endbr64
402084 push    0x5

```

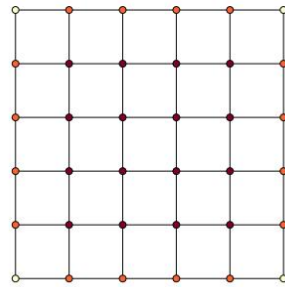
# Python *angr* library



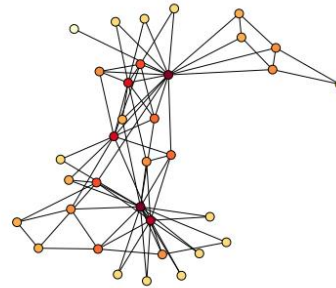
## Binary File

## Control Flow Graph (CFG)

# Graph embedding is an important process

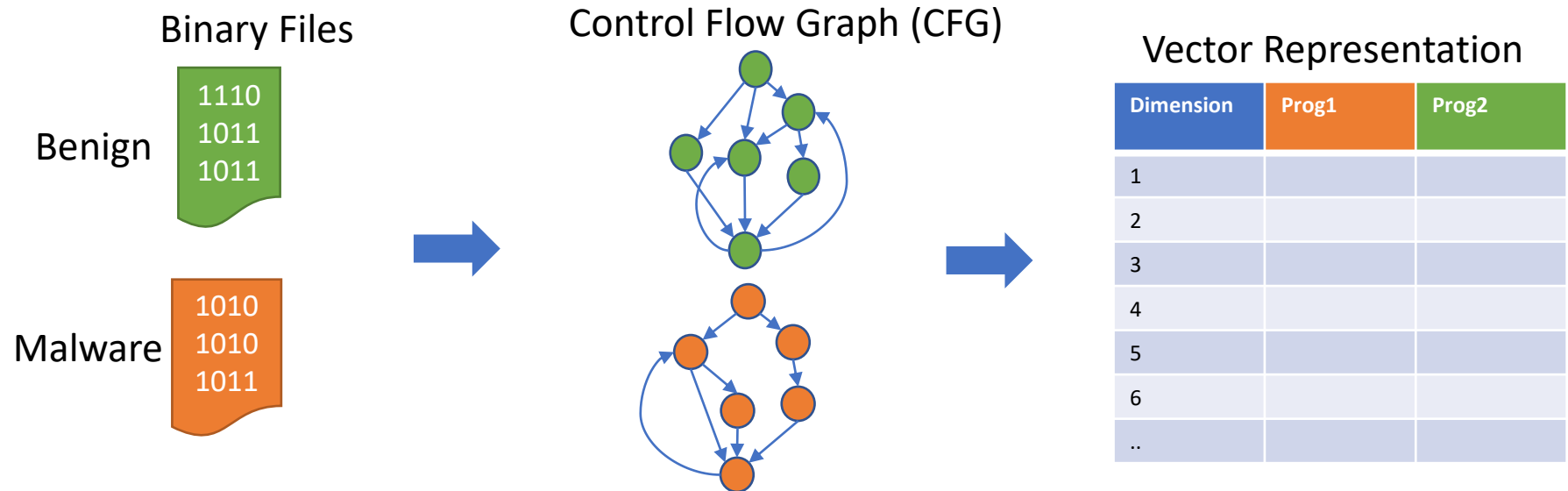


(a) Grid (Euclidean).

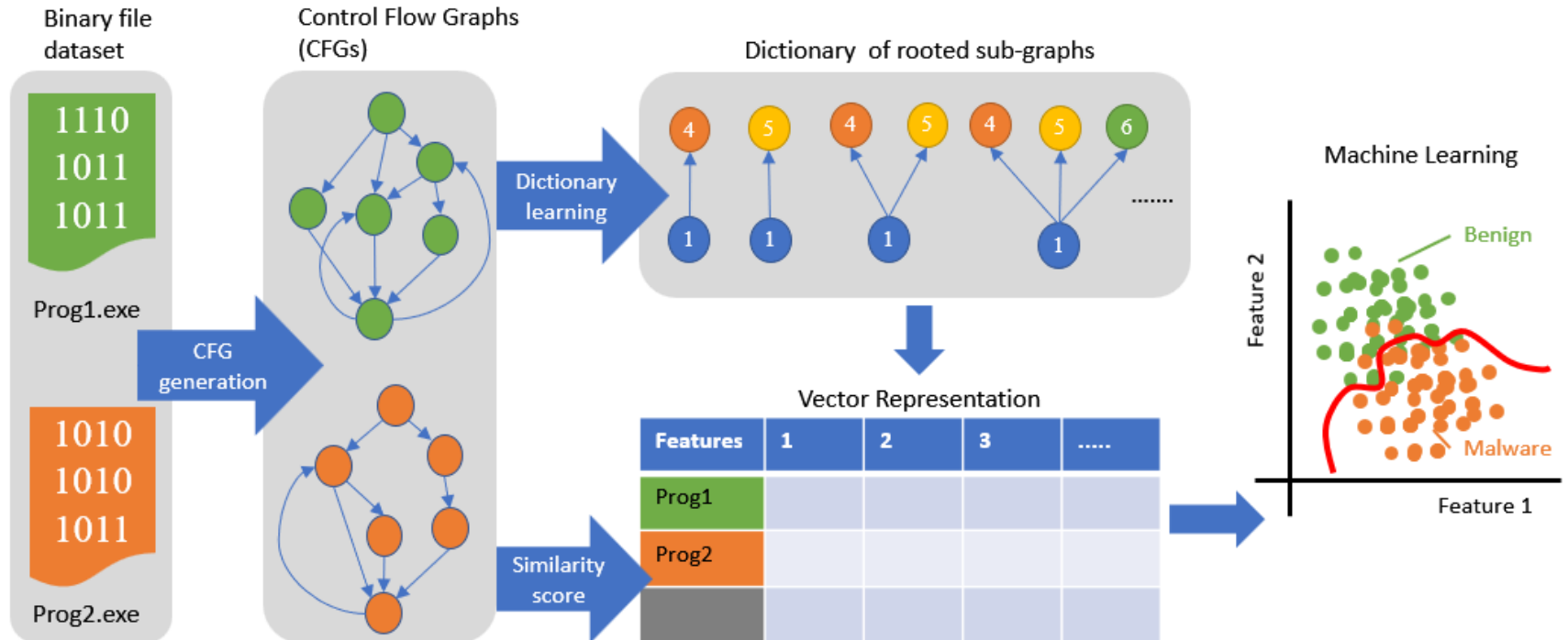


(b) Arbitrary graph (Non-Euclidean).

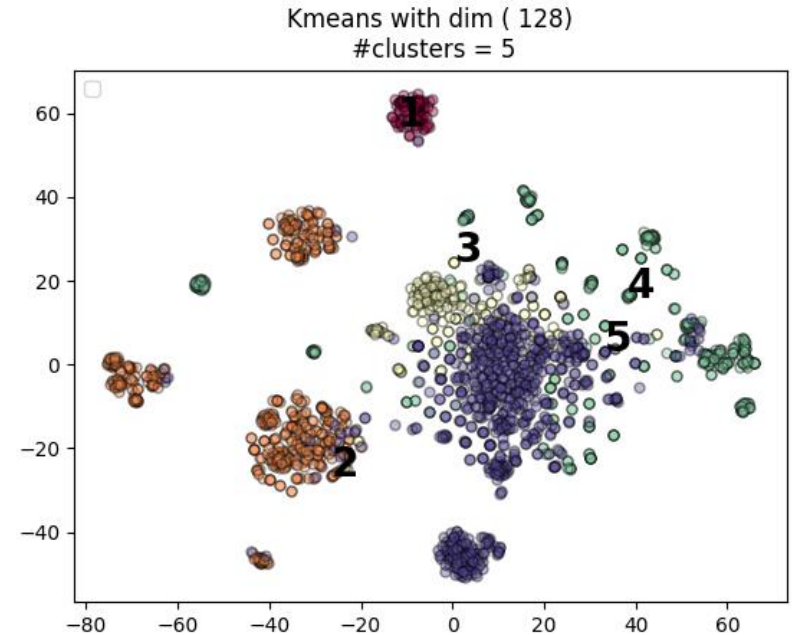
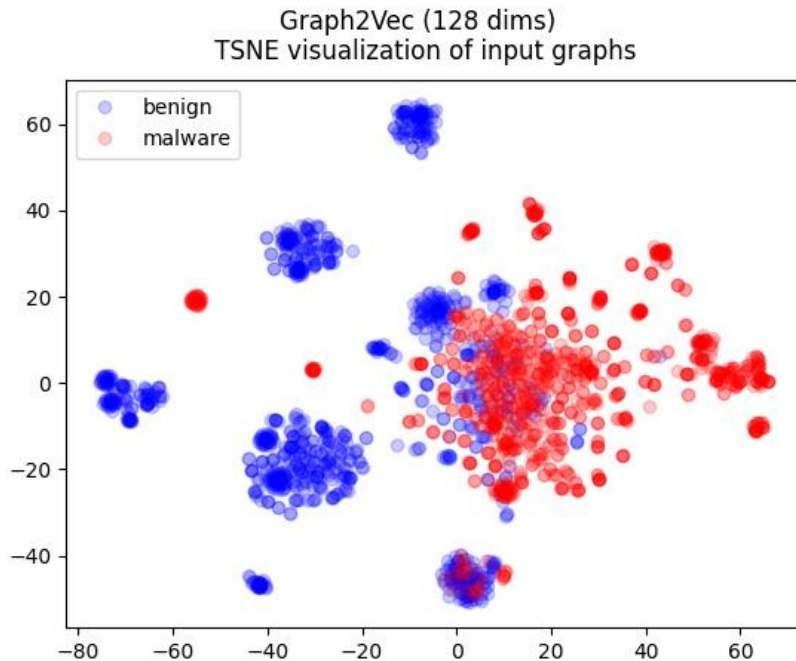
Figure 1: An illustration of Euclidean vs. non-Euclidean graphs.



# Overview of binary file analysis using CFGs



# Unsupervised clustering of binary files using CFGs



# Goal

Incorporate Sparse representation and its features for whole graph embedding.

A Graph is represented as one sparse vector and two graphs with similar sub- structure are embedded to be closer.

The Graph can be

- Directed

- Cyclic

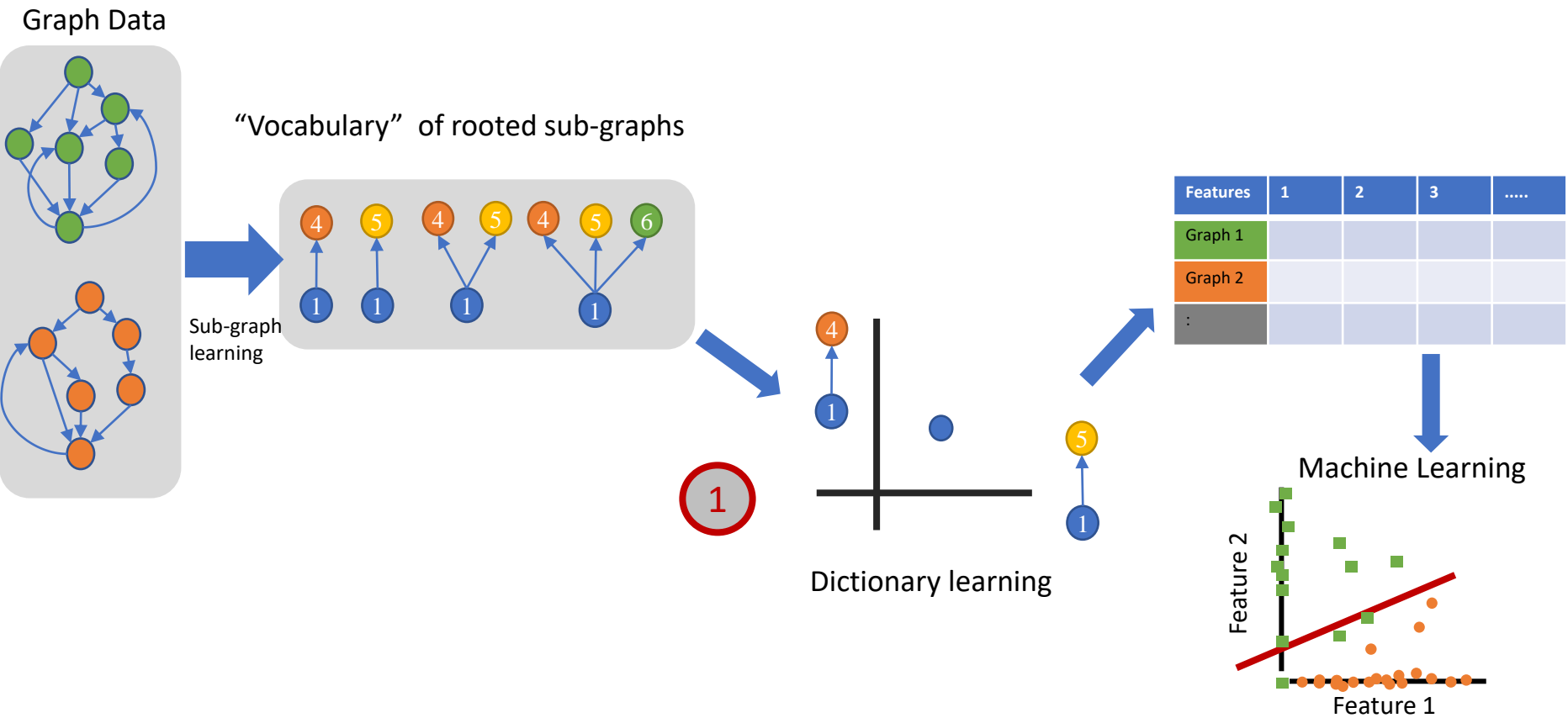
- Have node feature



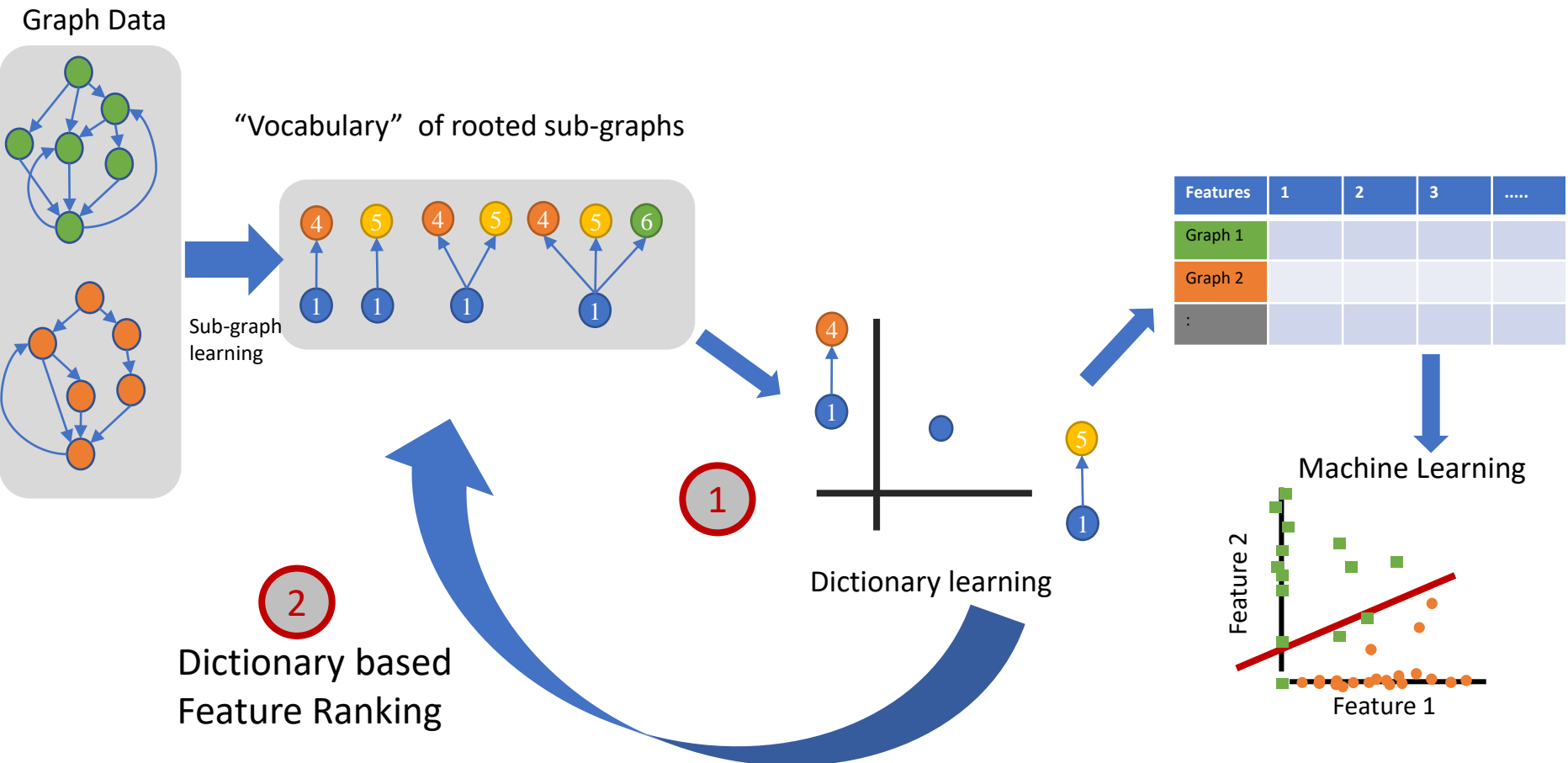
# Objectives

1. To develop a sub-tree pattern based sparse graph embedding method  
WL+KSVD
2. Framework for Identifying important sub-tree patterns
3. Develop a sparse graph representation in hyperbolic space

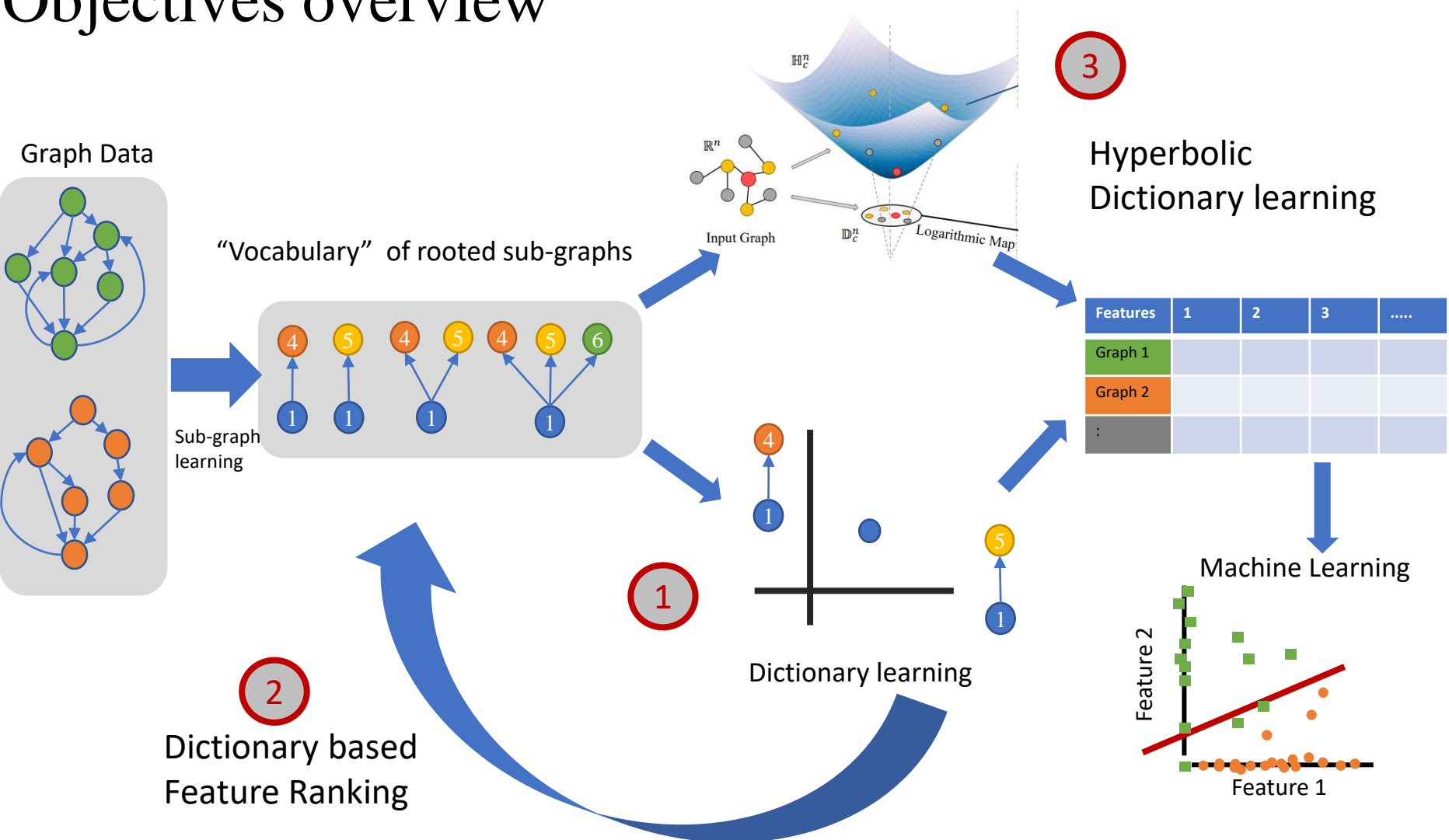
# Objectives overview



# Objectives overview



# Objectives overview



# Overview

## Introduction

Motivation / Problem

Goal

Objectives

## Background

Graph embedding

Sparse representation

Feature ranking

## Workflow

Preliminary work

Proposed timeline

Future directions

# Graph embedding

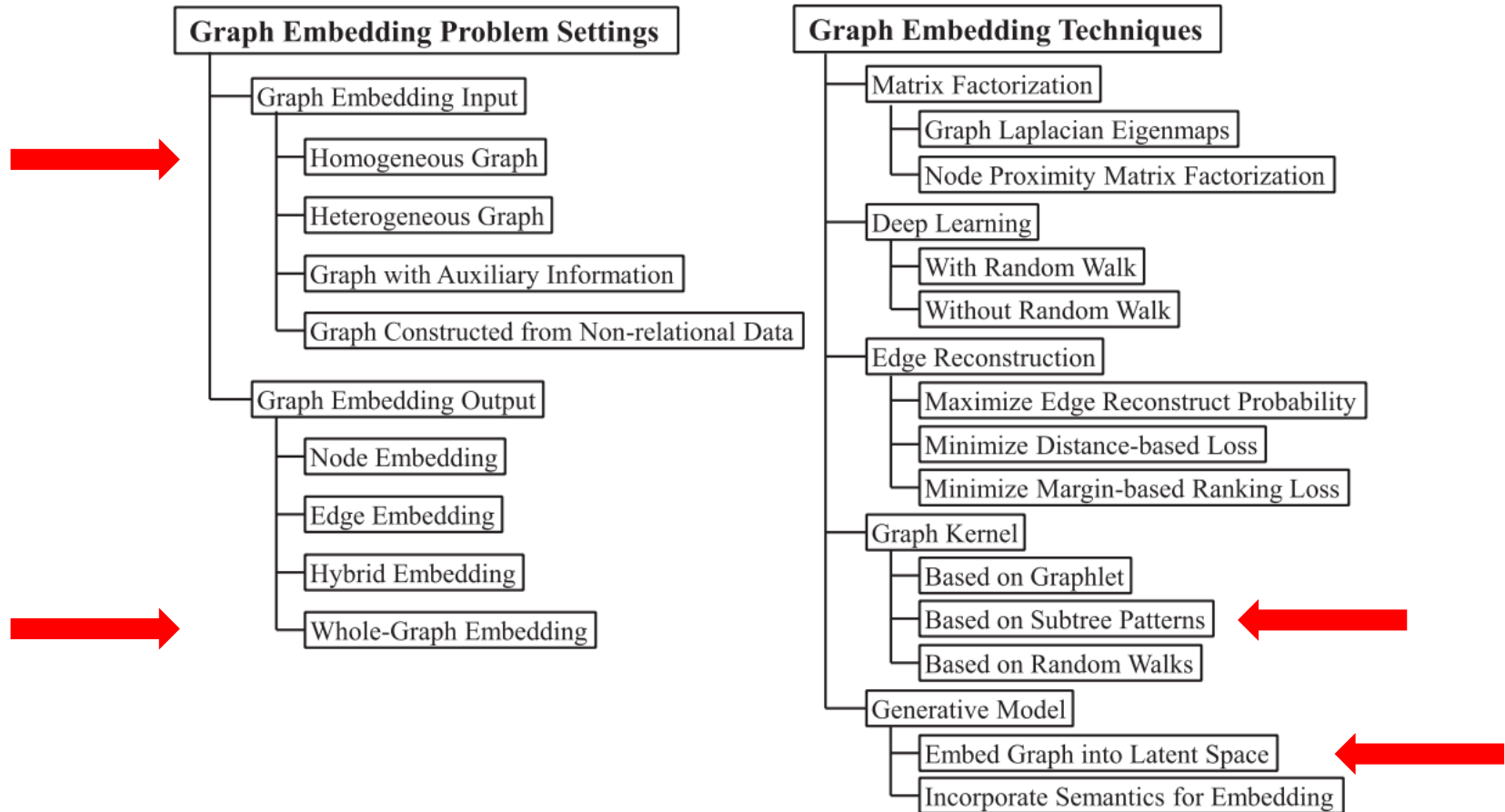
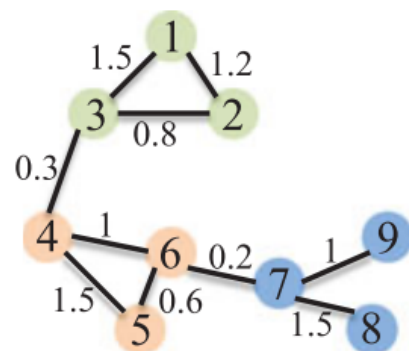
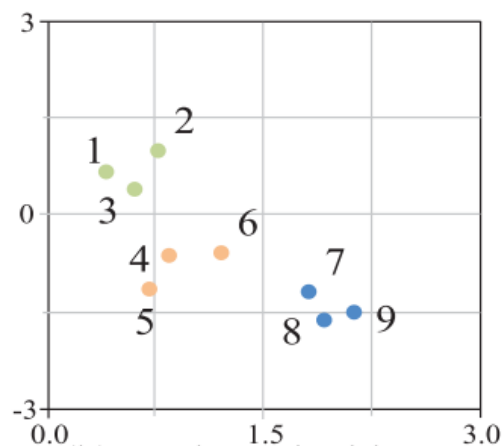


Fig. 2. Graph embedding taxonomies by problems and techniques.

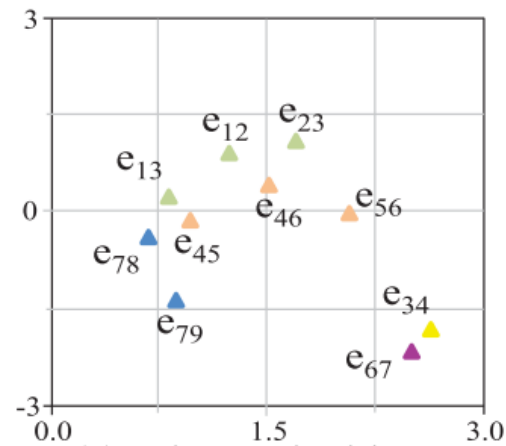
# Graph embedding types



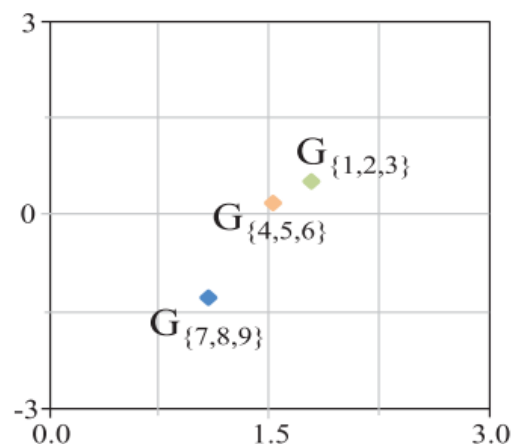
(a) Input Graph  $G_1$



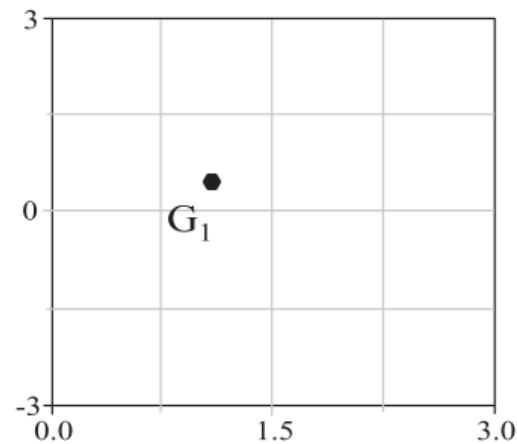
(b) Node Embedding



(c) Edge Embedding

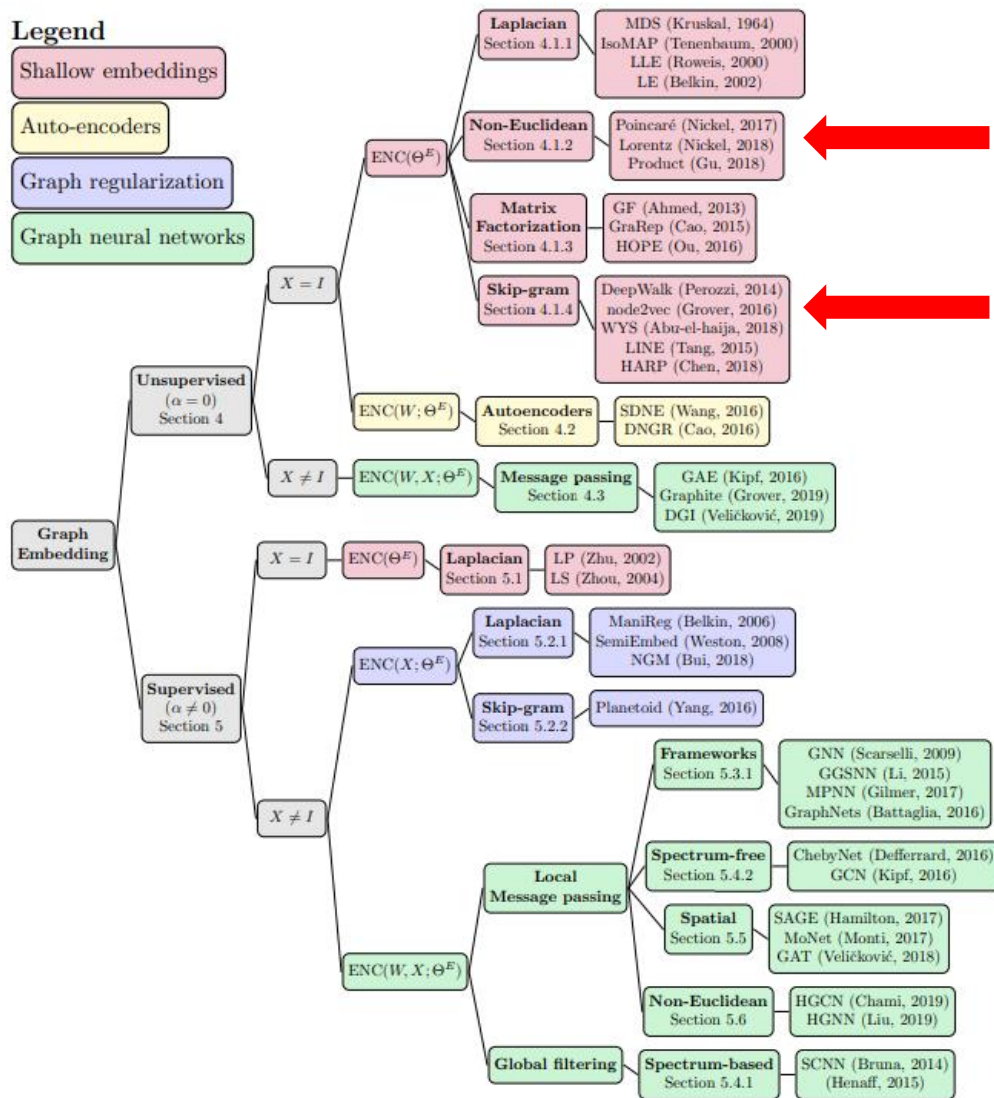


(d) Substructure Embedding



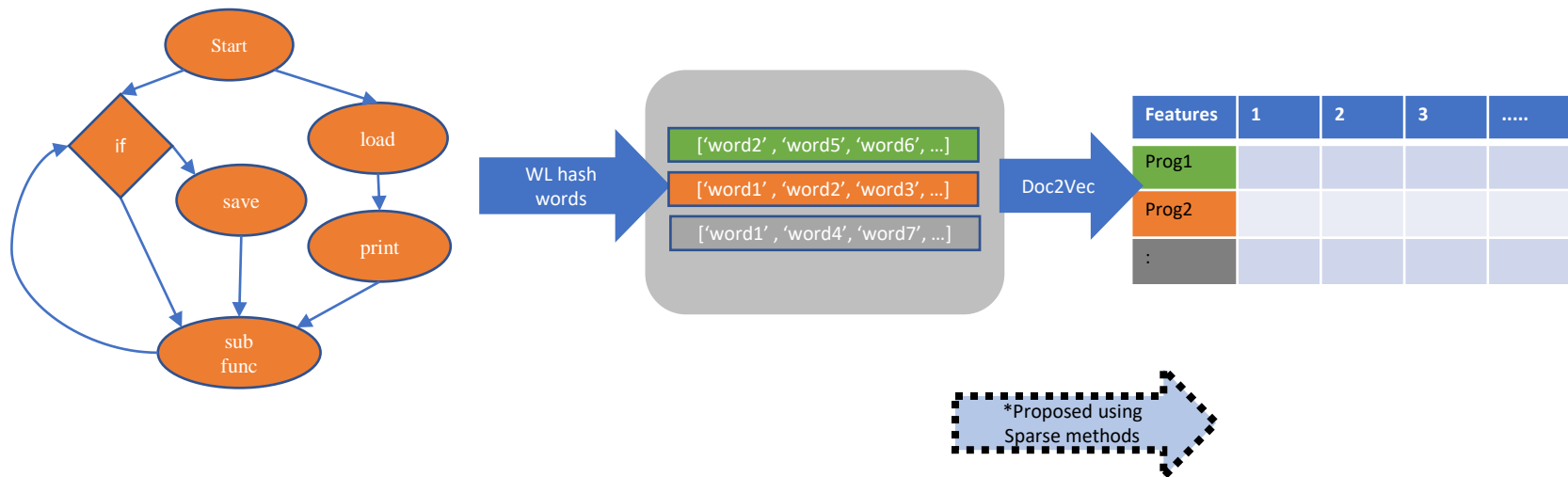
(e) Whole-Graph Embedding

# Graph embedding techniques

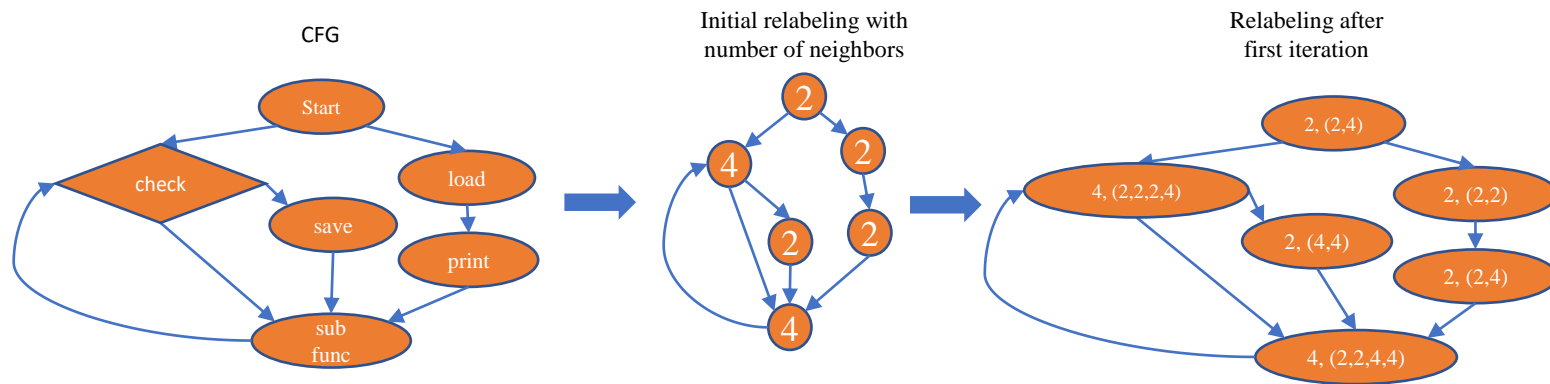




# Graph2Vec overview



# Weisfeiler-Lehman subtree hash



Node	Word
Start	2, hash(2,4)
check	4, hash(2,2,2,4)
Load	2, hash (2,2)
Save	2, hash(4,4)
Print	2, hash(2,4)
Sub_func	4, hash(2,2,4,4)
...	....

# Word2Vec

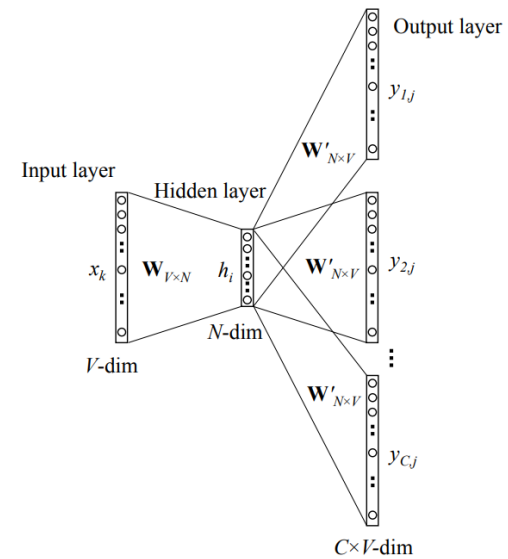
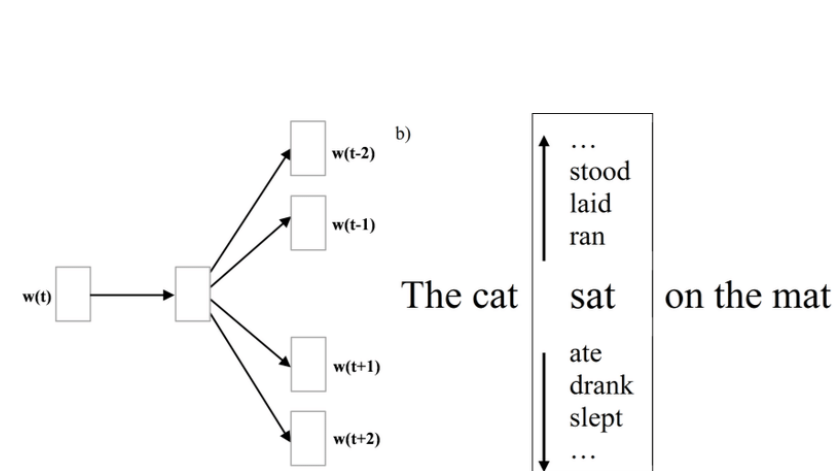
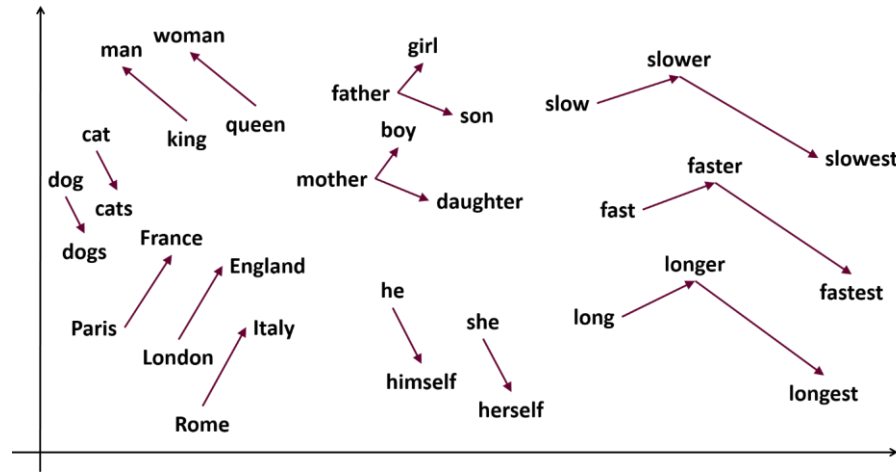
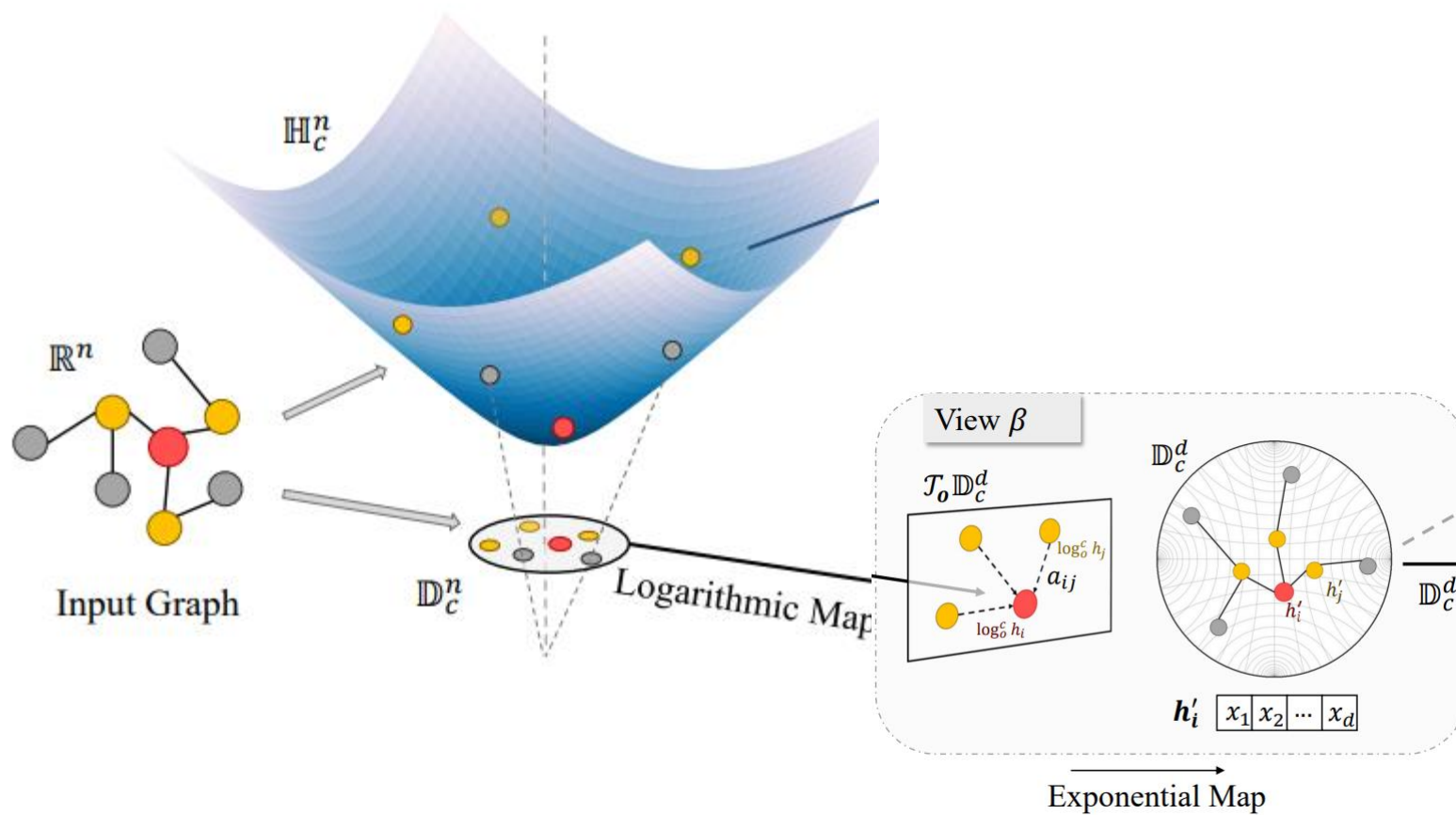


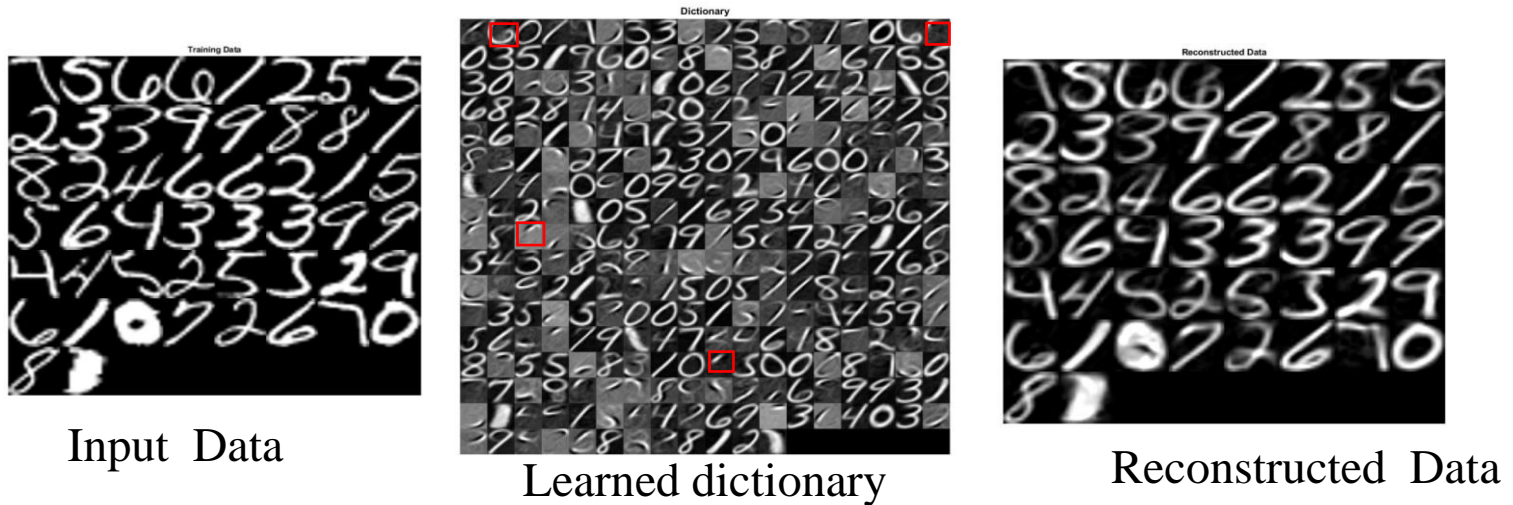
Figure 3: The skip-gram model.

<https://ronxin.github.io/wevi/>

# Poincare map



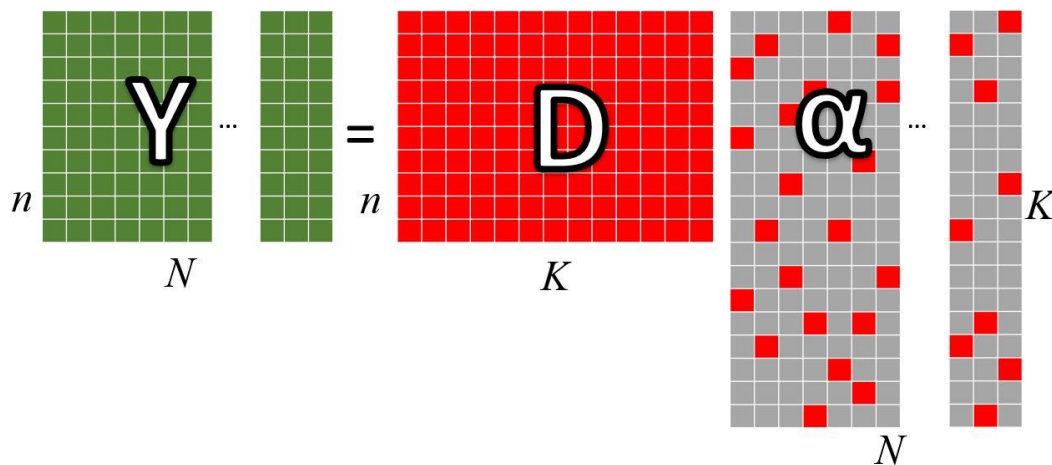
# Sparse Representation



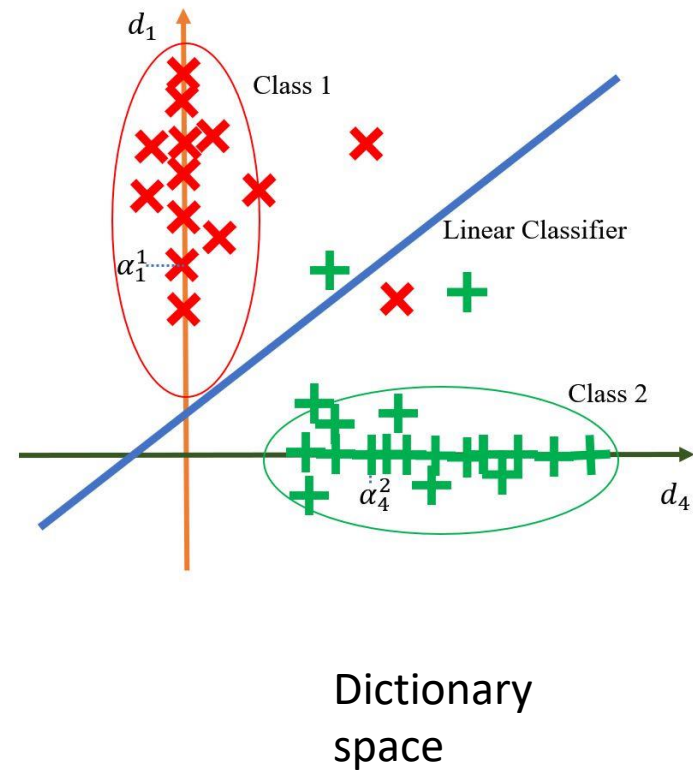
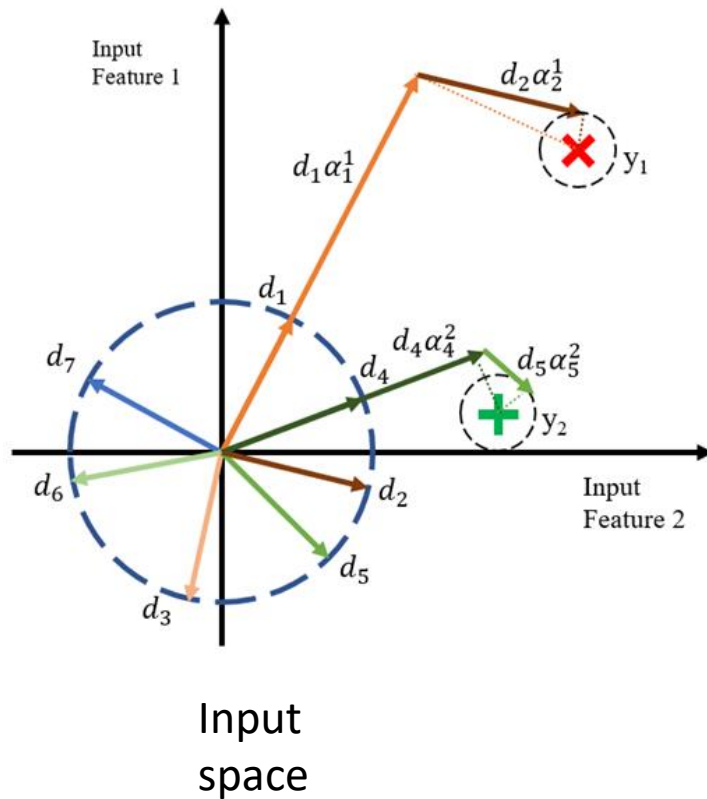
# Sparse Representation

- Input signals :  $\mathbf{Y} = [y_1, y_2, \dots, y_N] \in R^{n \times N}$
- Dictionary elements :  $d_i \in R^n$
- Dictionary :  $\mathbf{D} = [d_1, d_2, \dots, d_K] \in R^{n \times K}$
- Resulting signal :  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N] \in R^{K \times N}$
- Sparsity :  $S$

$$\underset{\mathbf{D}, \boldsymbol{\alpha}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \text{ s.t. } \forall i, \|\alpha_i\|_0 \leq S,$$



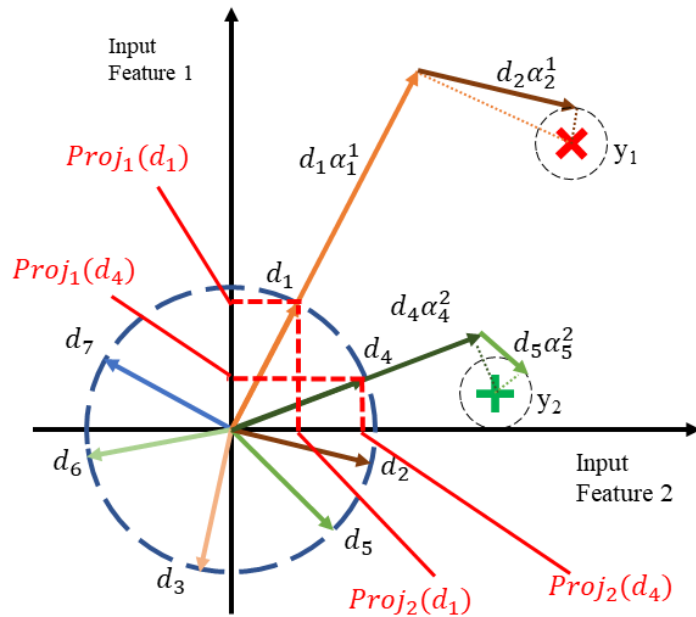
# Learned dictionary is in the same input space



# Dictionary based Feature Ranking metrics

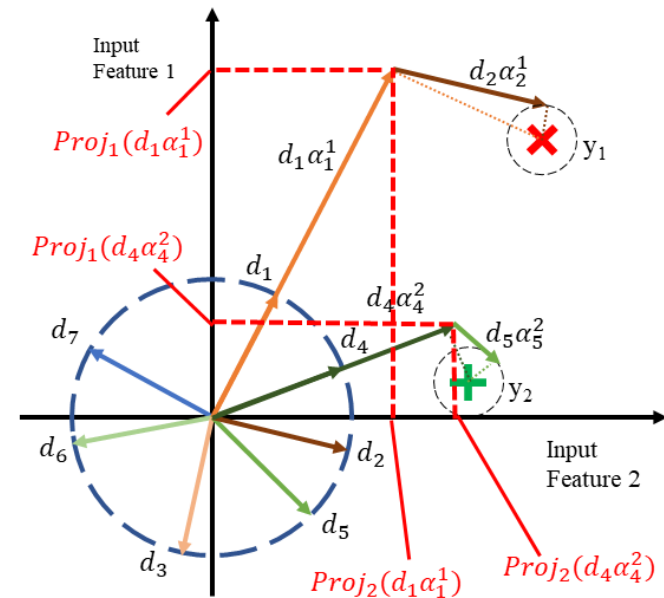
## Dictionary mapping

$$\mathbf{D}_{\text{map}}(j) = \sum_{i=1}^m \text{Proj}_j(D_i)^2 = \sum_{i=1}^m \mathbf{D}_{(j,i)}^2$$



## Dictionary utilization

$$\mathbf{D}_{\text{util}}(j) = \sum_{i=1}^m \text{Proj}_j(D_i \cdot \sum_{k=1}^n |\alpha_{j,k}|)$$





# Sparse coding-based FR

## Pros

### Model agnostic

Calculated Sparse coefficients can be used with different models

### Simple relationship

More intuitive mapping

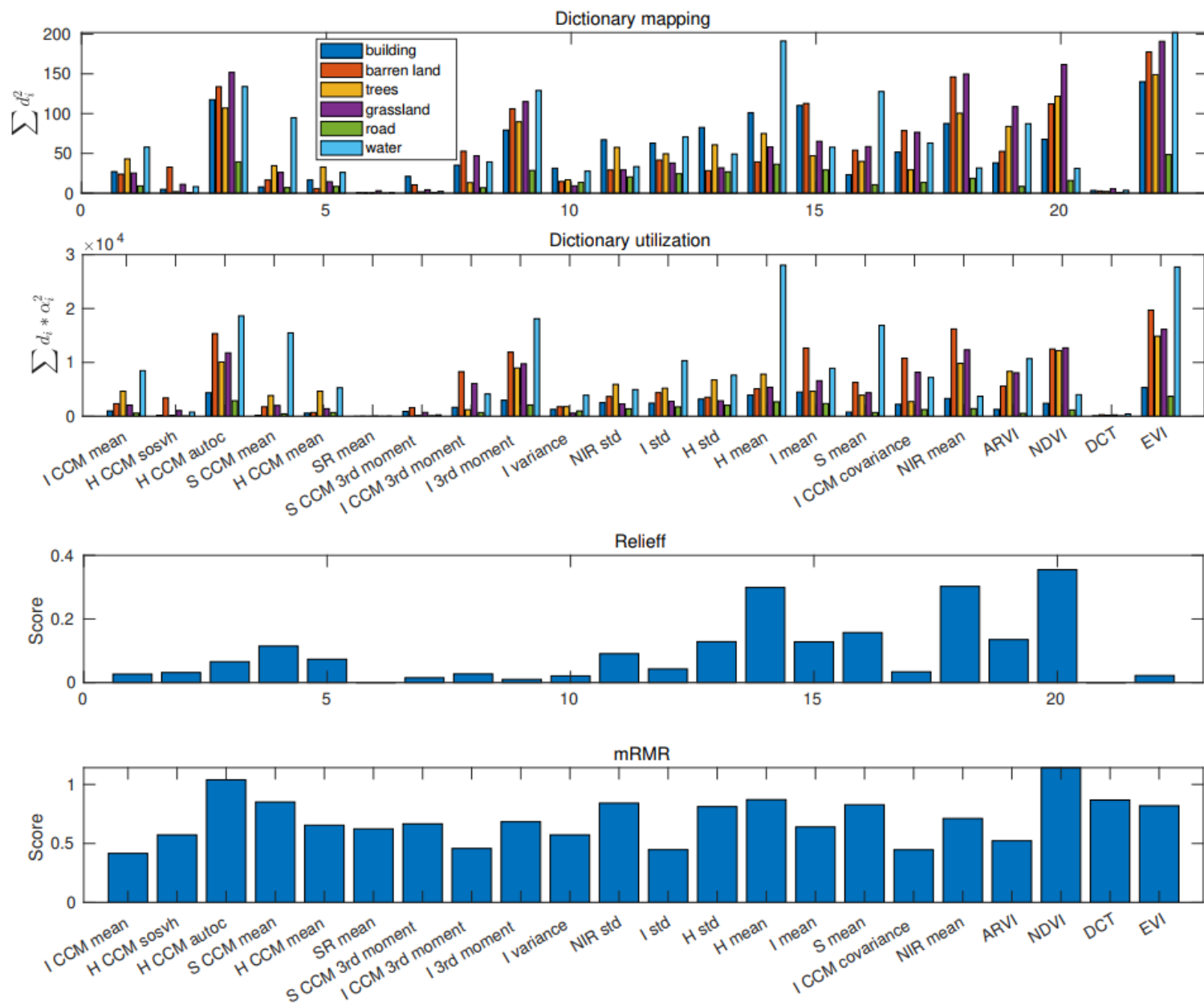
### Non-myopic manner

Take into consideration correlations and redundancies

## Cons

### Time must be spent to compute the sparse coefficients

Time spent learning sparse coefficients is not wasted



# Overview

## Introduction

Motivation / Problem

Goal

Objectives

## Background

Graph embedding

Sparse representation

Feature ranking

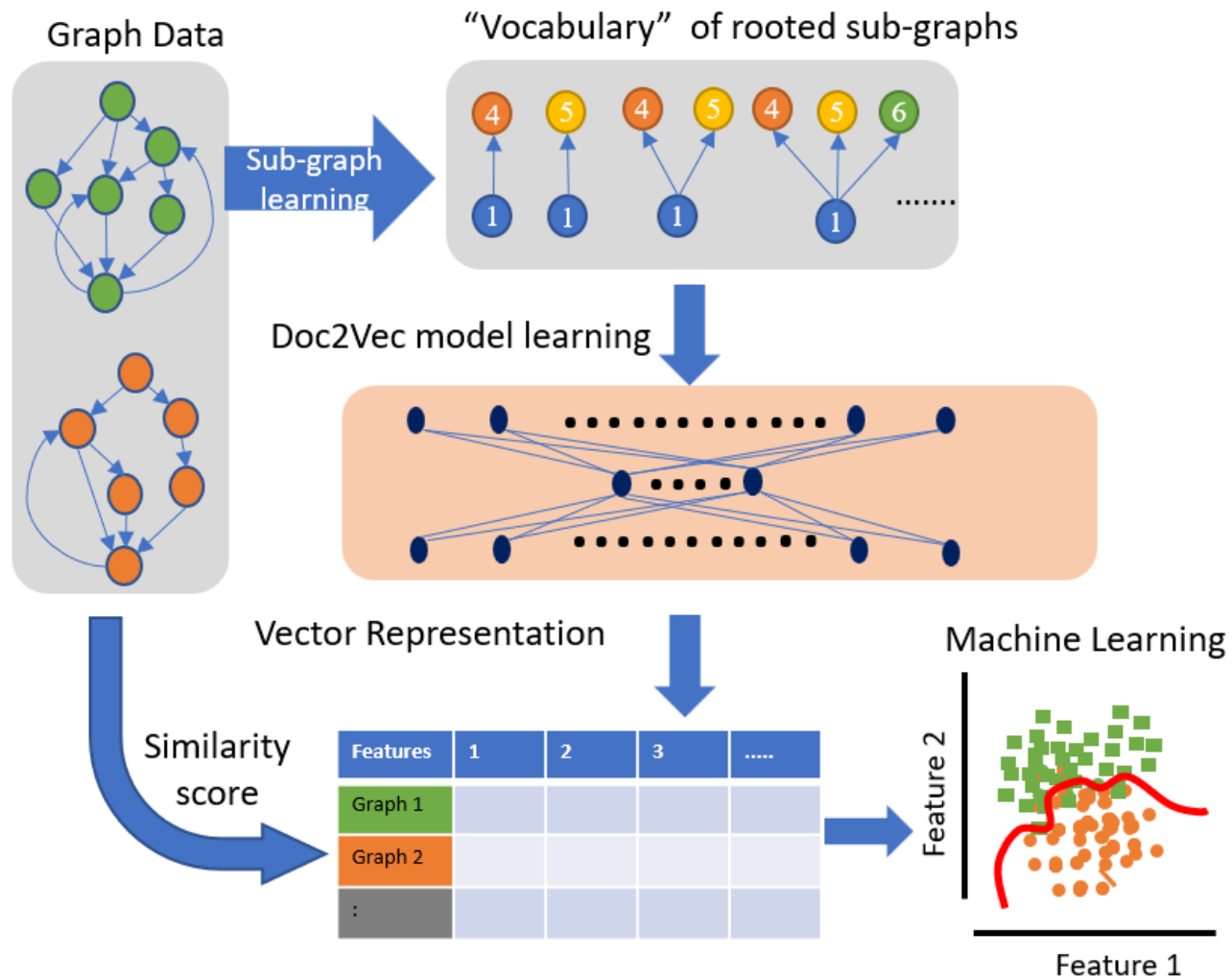
## Workflow

Preliminary work

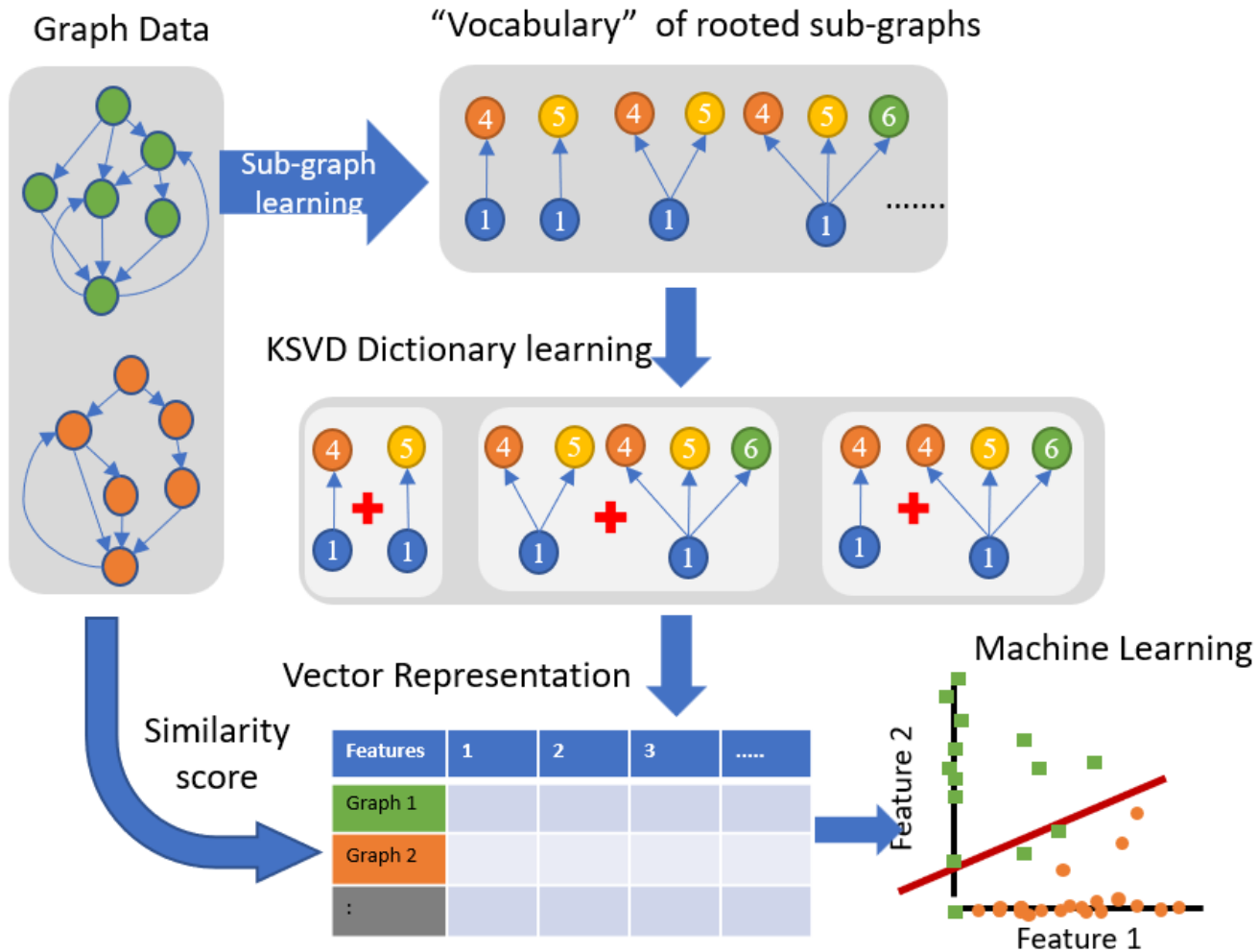
Proposed timeline

Future directions

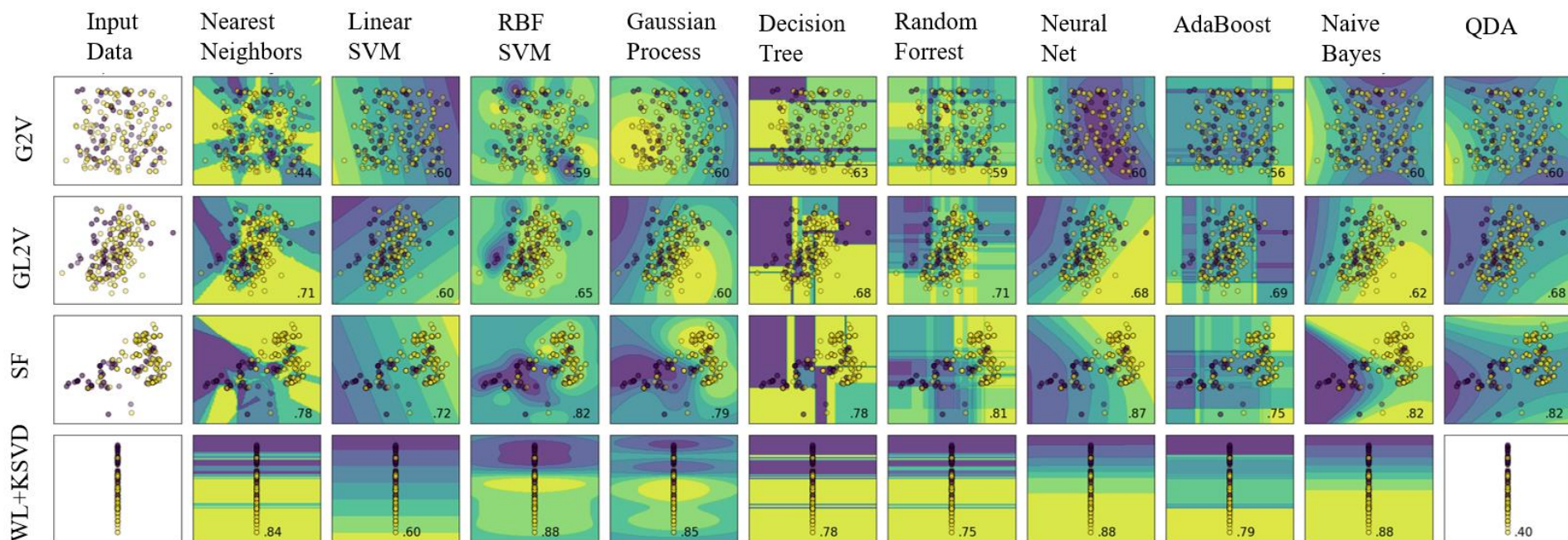
# Graph2Vec



# WL+KSVD



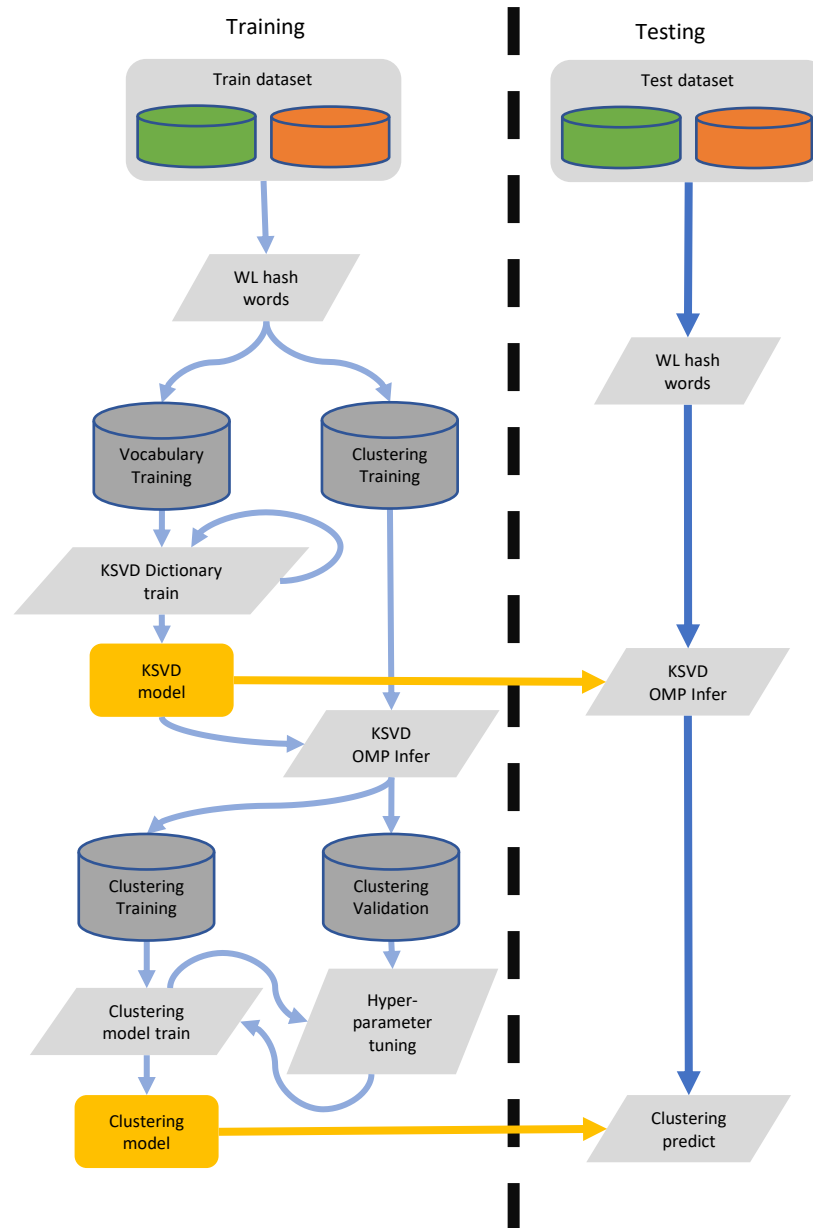
# Preliminary results



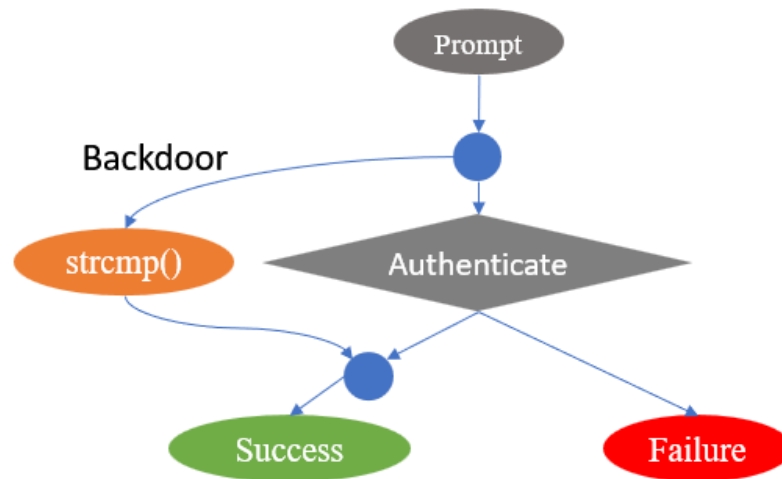
**Table 2.** Linear SVM accuracy with  $N = 1024$  embedding

	MU	PTC	PROT	NCI1	NCI109
G2V	68.55 ±10.03	55.23 ±5.9	67.30 ±0.87	59.30 ±4.46	56.46 ±2.82
GL2V	74.92 ±7.8	52.04 ±6.5	69.09 ±1.38	64.52 ±1.99	62.98 ±2.97
SF	83.47 ±4.15	57.59 ±9.34	70.98 ±1.00	61.90 ±3.24	61.96 ±2.40
WL+KSVD	72.38 ±3.20	54.36 ±2.31	64.60 ±2.00	64.16 ±2.19	62.93 ±0.24

# Pipeline



# Identifying vulnerable subtree structures



Conceptual authentication bypass vulnerability



# Hyperbolic space

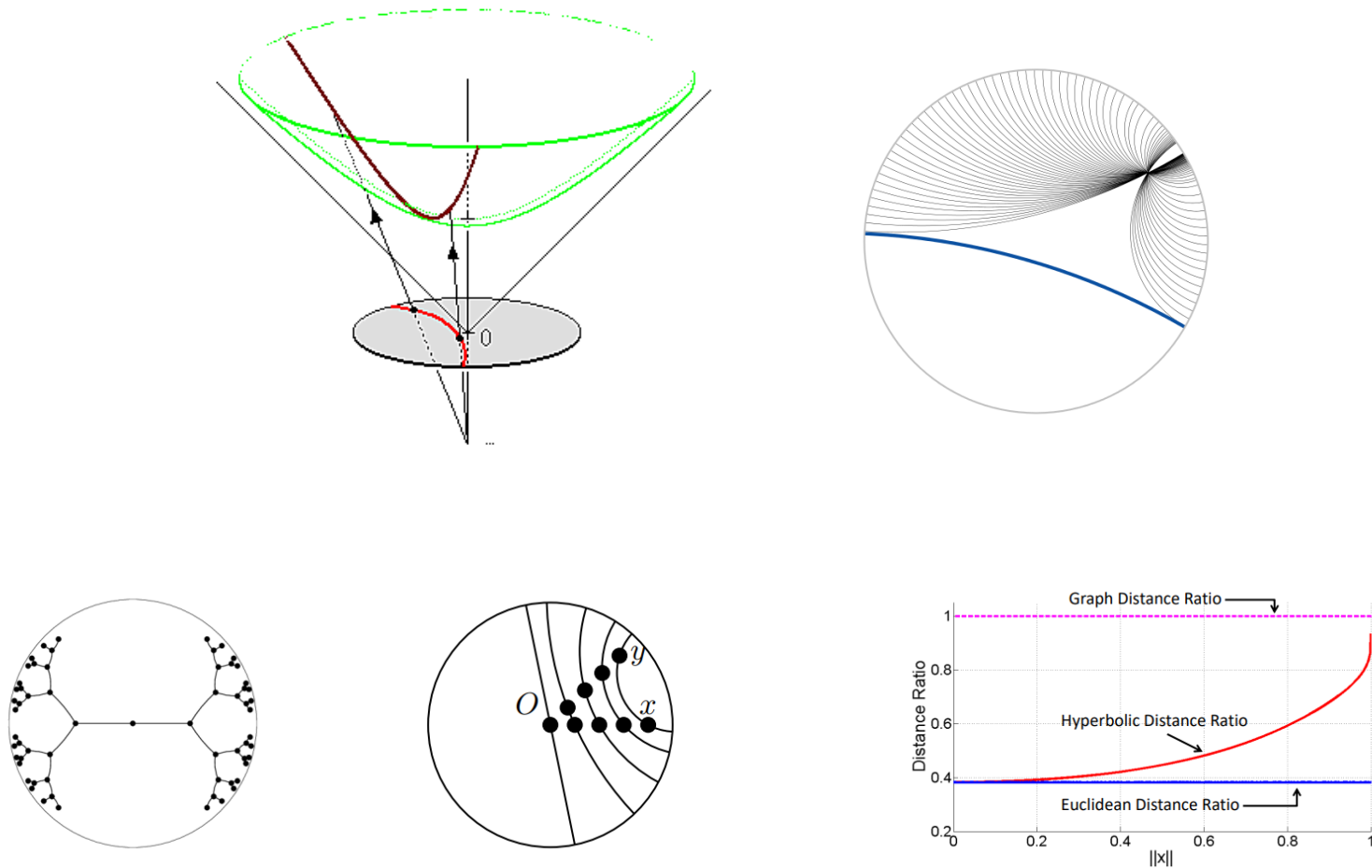
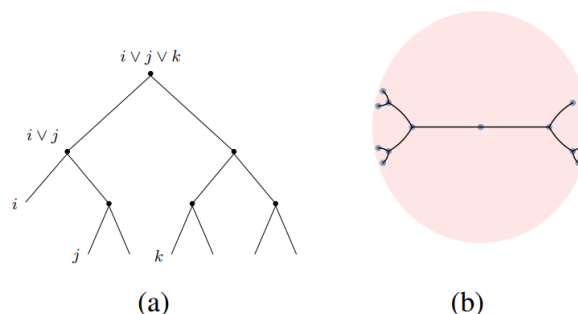


Figure 1. Left: Embedding of a binary tree in the Poincaré disk. Right: Geodesics and distances. As  $x$  and  $y$  move towards the outside of the disk (i.e., letting  $||x||, ||y|| \rightarrow 1$ ), the distance  $d_H(x, y)$  approaches  $d_H(x, O) + d_H(O, y)$ .

# Hyperbolic space dictionary learning

## Subtree embedding in Hyperbolic

Chami, I., Gu, A., Chatziafratis, V. and Ré, C., 2020. From trees to continuous embeddings and back: Hyperbolic hierarchical clustering. *Advances in Neural Information Processing Systems*, 33, pp.15065-15076.



## SVD in hyperbolic

Onn, R., Steinhardt, A.O. and Bojanczyk, A., 1989, August. The hyperbolic singular value decomposition and applications. In *Proceedings of the 32nd Midwest Symposium on Circuits and Systems*, (pp. 575-577). IEEE.

## K-means in hyperbolic

Djeddal, H., Touzari, L., Giovanidis, A., Phung, C.D. and Secci, S., 2021. Hyperbolic K-means for traffic-aware clustering in cloud and virtualized RANs. *Computer Communications*, 176, pp.258-271.

## Pursuit algorithms in hyperbolic

Tabaghi, P. and Dokmanić, I., 2020, August. Hyperbolic distance matrices. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (pp. 1728-1738).

# Q&A

## Summary

### Objectives

- Sparse dictionary learning on subtree patterns
- Feature ranking of subtree patterns
- Sparse representation on Hyperbolic space

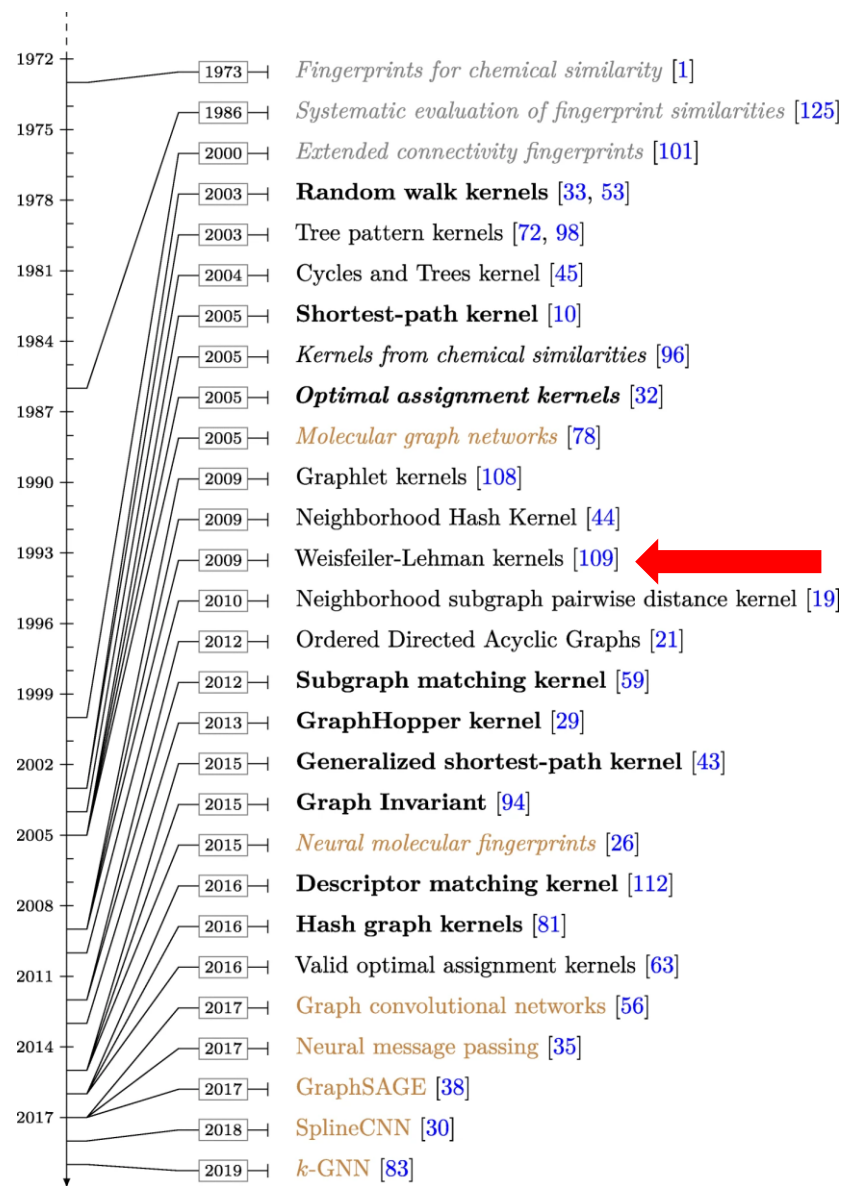
### Advantage

- Linear relationship
- Simpler (Low order) ML model
- Intuitive

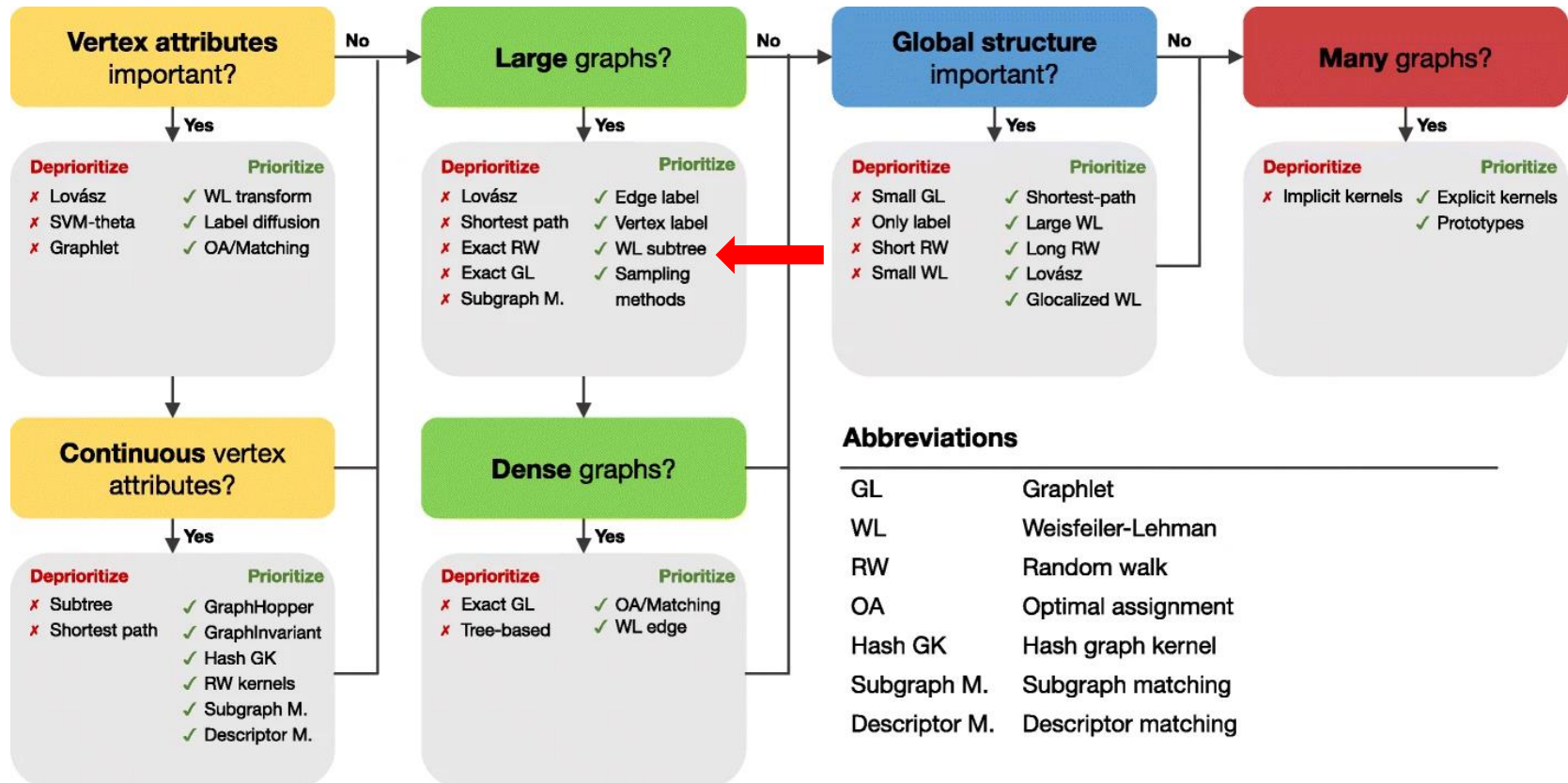
### Outcomes

- Publicly available python package

# Graph Kernels



# Graph Kernel selection guideline



Guidelines for prioritizing kernels for consideration based on known properties of the graph learning problem.

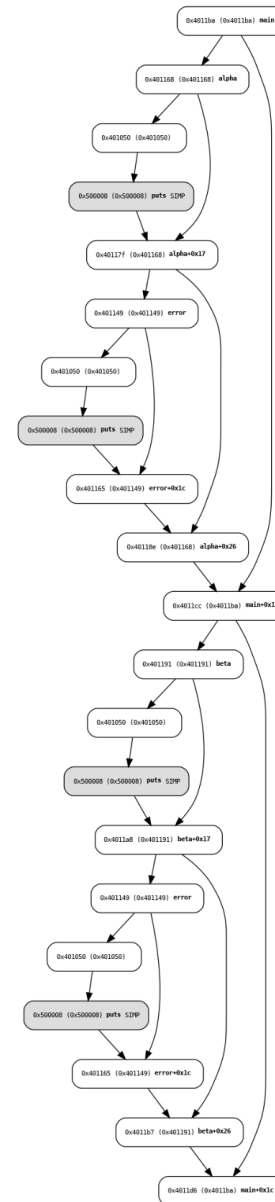
# Typical CFG for a given program

```
void error(char *error)
{
    puts(error);
}

void alpha()
{
    puts("alpha");
    error("alpha!");
}

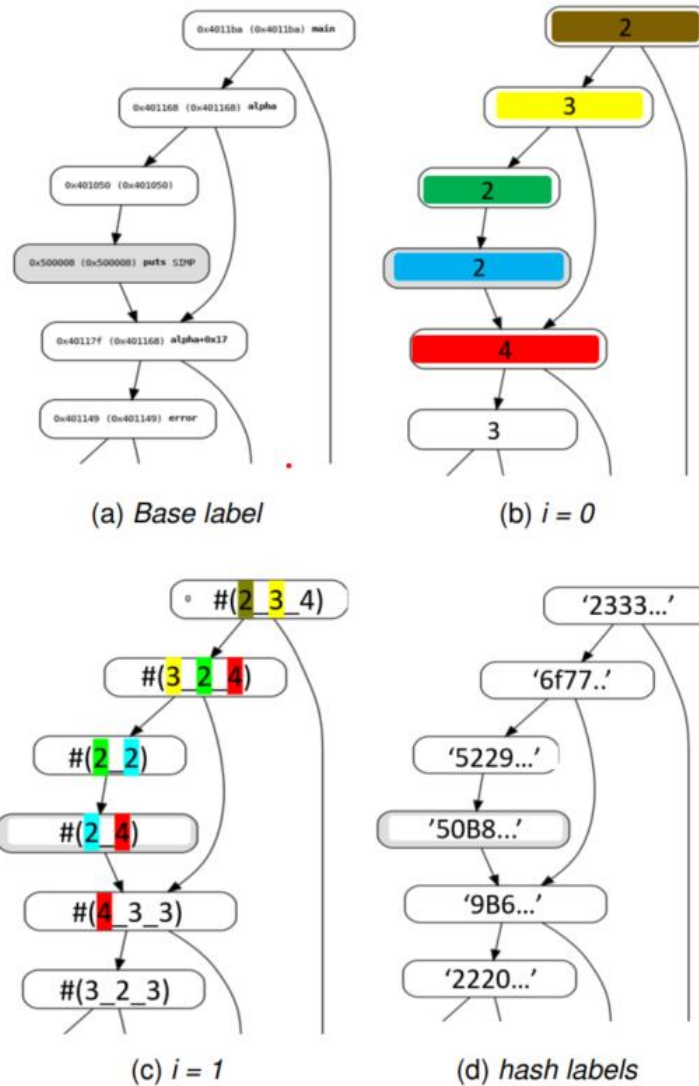
void beta()
{
    puts("beta");
    error("beta!");
}

void main()
{
    alpha();
    beta();
}
```



Conceptual authentication bypass vulnerability

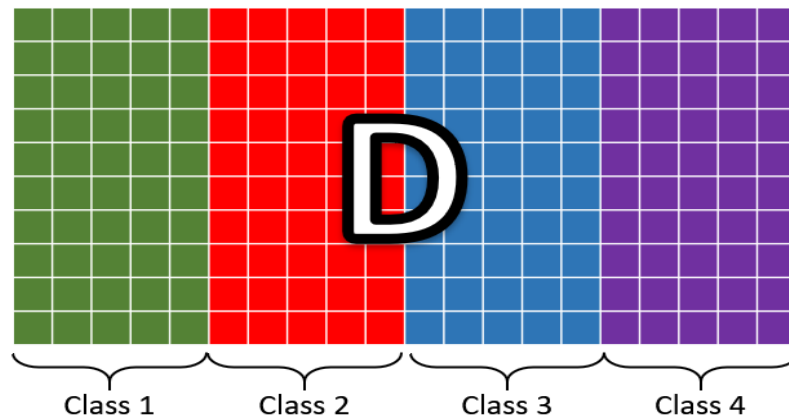
# WL subtree hash relabeling



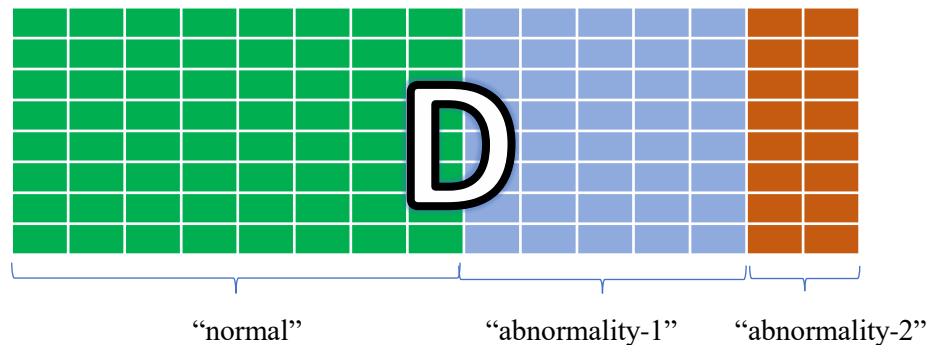
WL graph hash algorithms' first iteration in the context of CFG

# Supervised dictionary learning

Label Consistent-KSVD

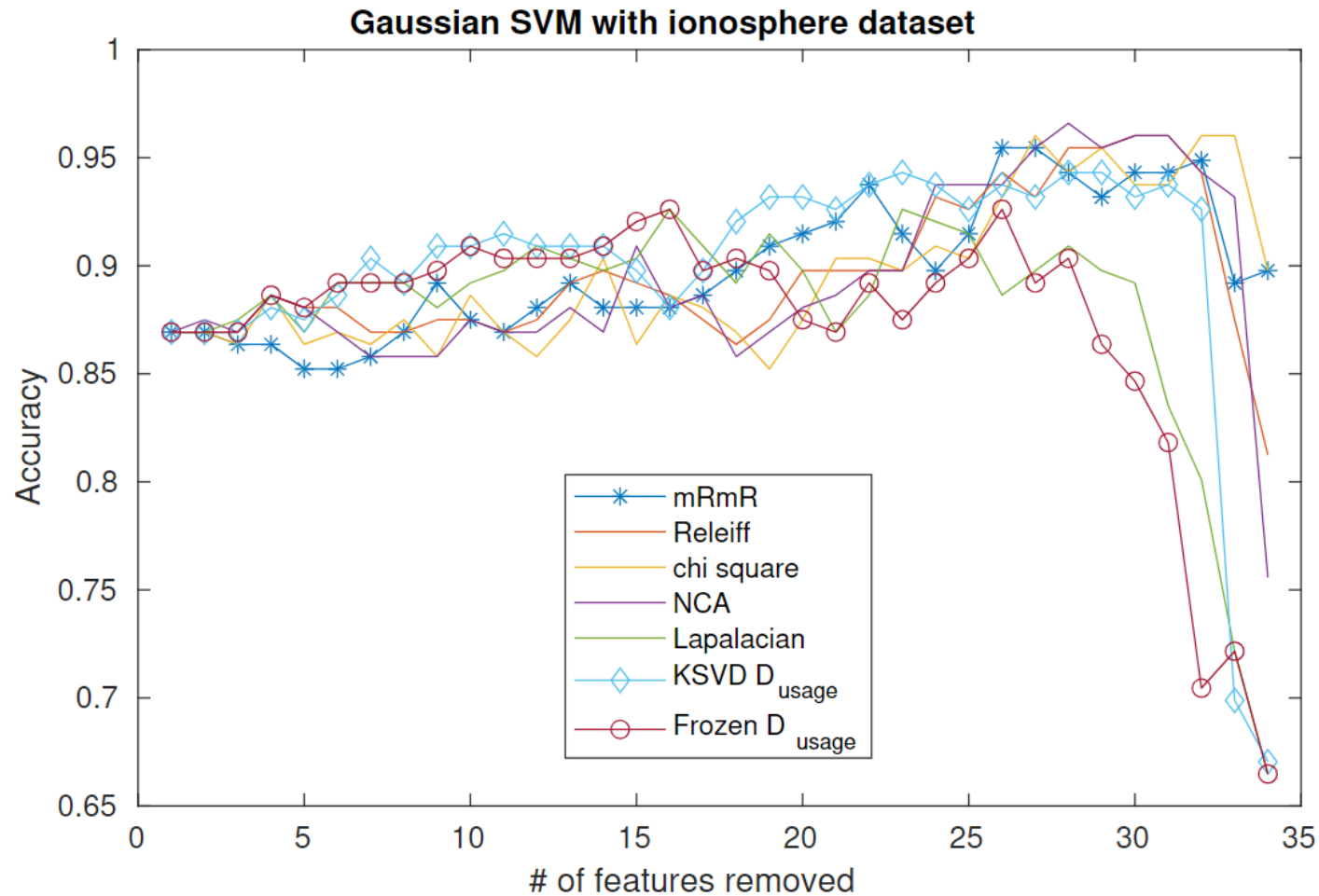


Frozen-KSVD





# Feature selection



# Overview

## Introduction

Motivation

Problem

Goal

## Background

Graph embedding

Sparse representation

Feature ranking

## Objectives

Sparse dictionary learning on subtree patterns

Feature ranking of subtree patterns

Sparse representation on Hyperbolic space