

# Combinatorial and Neural Graph Vector Representations

*Evgeny Burnaev*

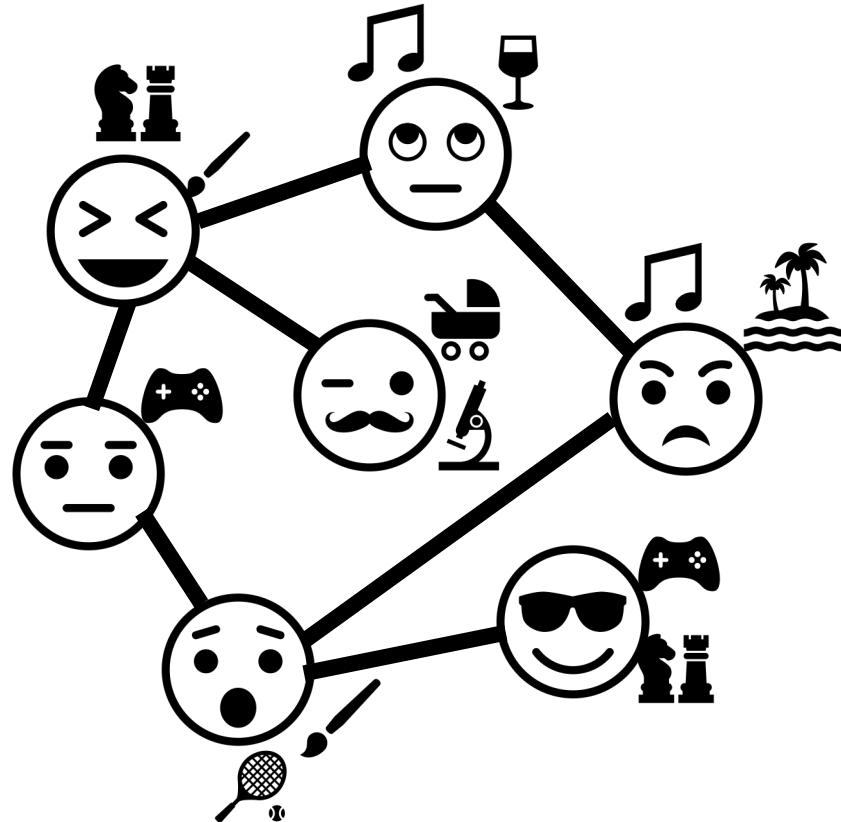
Head of ADASE group, Skoltech

joint with *Sergey Ivanov, Sergey Sviridov*

# Product recommendation

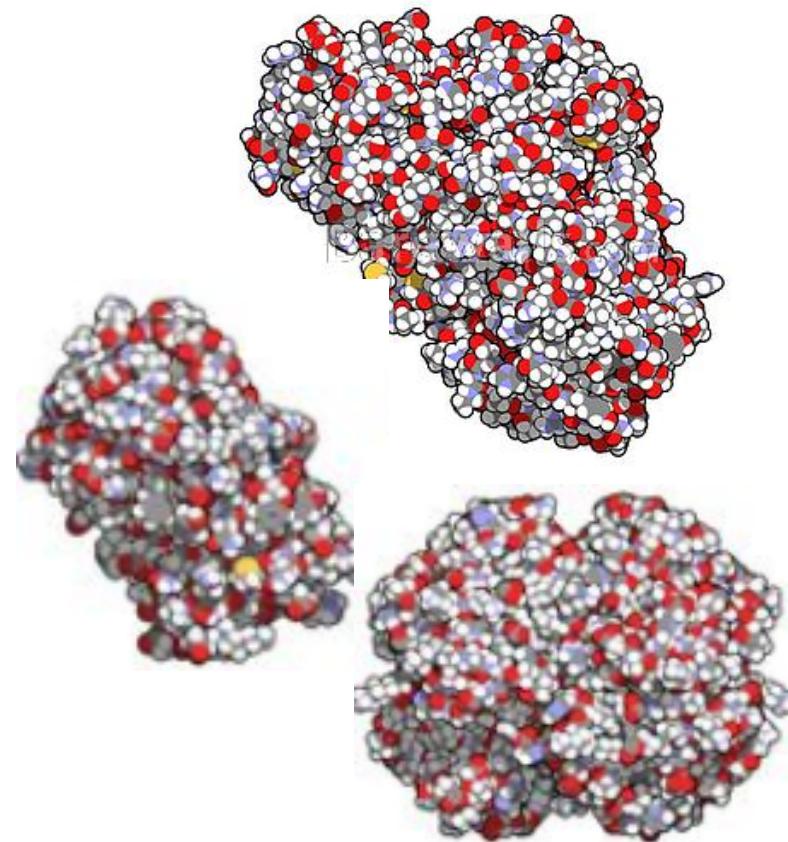
*How to find people  
that maximize product  
adoption?*

*How to scale solutions  
to billions users and  
consider user  
preferences?*



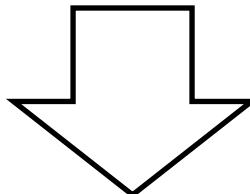
# Protein function prediction

*How to find similar proteins  
based on structural,  
physical, and chemical  
information?*



# What is common?

- We can model these *problems with graphs*;
- We can find solution for some classes of problems.



We can use ML methods on graphs to solve new instances.

*How to represent graphs for ML models?*

# Representation learning on graphs

Topological Descriptors  
[1970-2000]

Graph Neural Networks  
[2015-Now]

Graph Kernels  
[1999-2019]

Graph vector representation  $\nu \in R^d$  is called *embedding*.

# Topological descriptors

Simple feature vectors or a scalar number (e.g. Wiener index).

## Pros:

- Simple and inspired by chemical properties of studied networks

## Cons:

- Very limited scope
- Ad-hoc design
- Prediction is not efficient (e.g. via knn)

[1] Handbook of molecular descriptors, Wiley & Sons 2008

# Graph kernels

Symmetric, positive semidefinite function that maps two graphs to a real number:

$$K(G_1, G_2) \mapsto R$$

# Graph kernels: pros

- More expressive than topological descriptors
- Includes different substructures
- Suitable for kernel machines (e.g. SVM)
- Non-linearity achieved by kernel trick
- Allows skipping explicit computation of vector representation
- Backed by theoretical developments in graph theory

# Graph kernels: cons

- Not scalable (often,  $O(n^3)$  running time)
- Requires  $N \times N$  kernel matrix
- Do not preserve graph isomorphism in feature space



Addressed in this work

# Isomorphism property

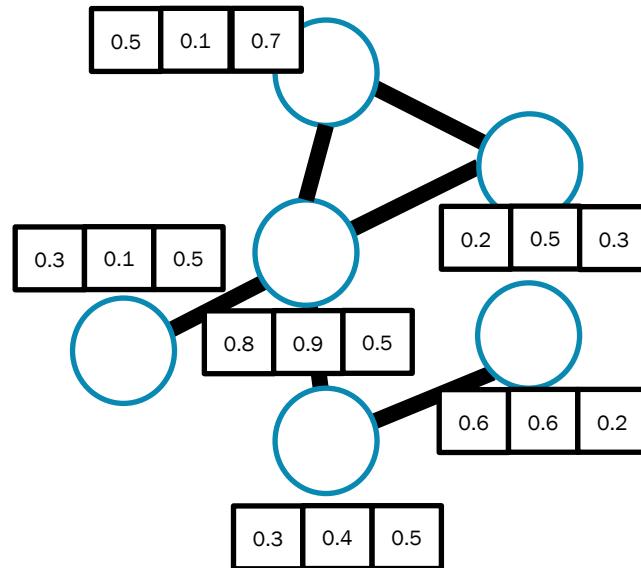
$$K(G_1, G_2) = \langle \varphi_1, \varphi_2 \rangle,$$

where  $\varphi: G \mapsto R^d$

If  $\varphi$  is bijective, then we say graph kernel has *isomorphism property*. Such graph kernel minimizes loss of information.

# Graph neural networks

Graph embedding is initialized with random vector, which is updated to fit the given data (e.g. CNN-like models).



# Graph neural networks: pros & cons

## Pros:

- SOTA performance in classification task
- Similar theoretical background

## Cons:

- Complex models
- Hardly interpretable

Addressed in this work

[1] A Comprehensive Survey on Graph Neural Networks, Wu et al. 2019

# Main goal

To develop efficient graph representation that:

- Has isomorphism property
- Inherits strong graph kernel and neural network sides
- Test efficiency on real-world problems

# Main results

- Proposed and justified a new graph representation that **strictly provides isomorphism property**;
- Designed two approaches for **approximate efficient computation of embeddings**;
- Demonstrated **superior quality of embeddings** in graph classification problem.

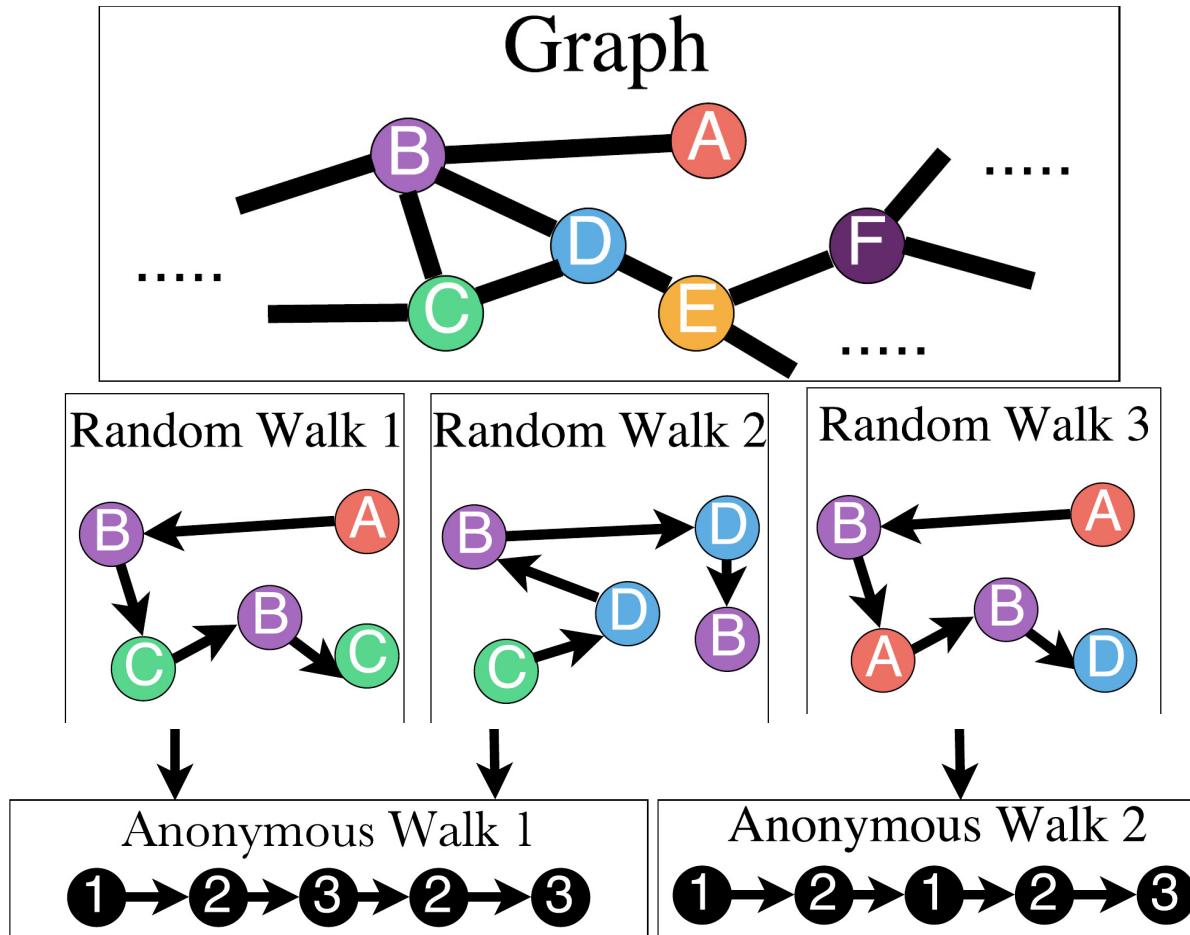
# 1. Anonymous Walks

# Anonymous Walks

Definition:

If  $w = (v_1, v_2, \dots, v_l)$  is a random walk then anonymous walk is the sequence  $a = (a_1, a_2, \dots, a_l)$  where  $a_i$  is the first position of  $v_i$  in  $w$ .

# Anonymous Walks



# Reconstruction property

Theorem [Zhu & Micali, 2015]:

Let  $B(\nu, r)$  be the induced graph at node  $\nu$  of radius  $r$  containing  $m$  edges, and  $D_l$  is a set of all possible anonymous walks of radius  $l = 2(m + 1)$ , that start at node  $\nu$ .

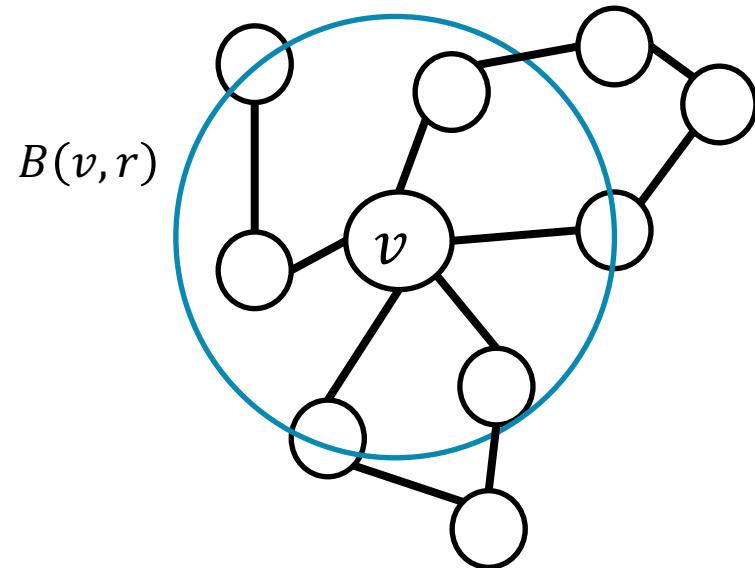
There is an algorithm to reconstruct a graph  $G$  that is isomorphic to  $B(\nu, r)$ .

## Theorem illustration

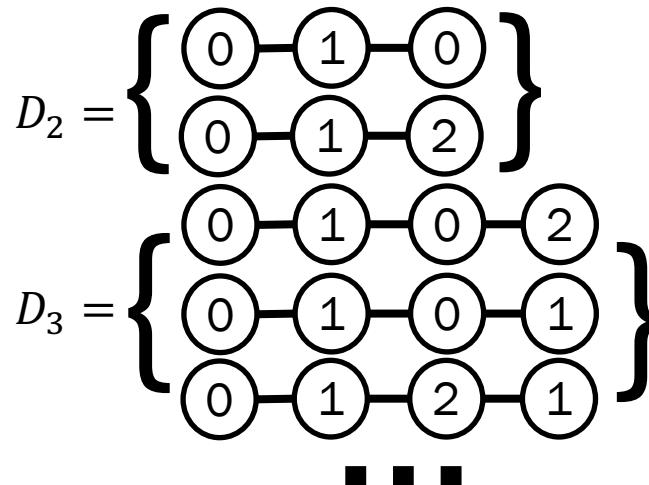
Radius  $r = 2$

Edges  $m = 6$

Length  $l = 2(m + 1) = 17$



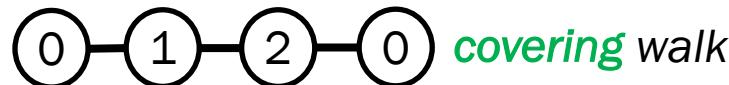
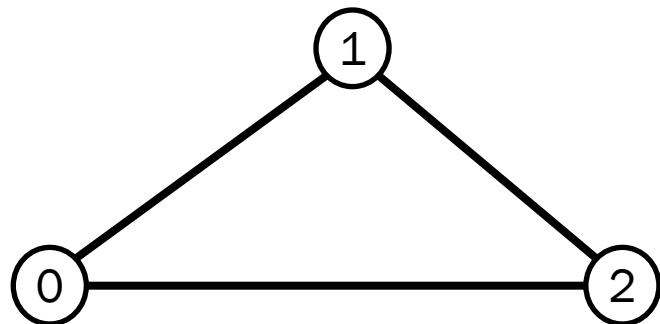
*Knowing distributions  $D_2, D_3, \dots, D_{17}$  we can obtain graph  $G$  that is isomorphic to  $B(v, r)$ .*



# Covering walks

Definition:

If anonymous walk traverses each edge of the graph at least once, we call it *covering walk*.



*covering walk*



*not covering walk*

# Isomorphism test

Theorem [This work]:

Let  $S_1$  and  $S_2$  be the sets of all covering walks of length  $l = 2(m + 1)$  for graphs  $G_1$  and  $G_2$  with  $m$  edges. Two graphs are isomorphic if and only if  $S_1 \cap S_2 \neq \emptyset$ .

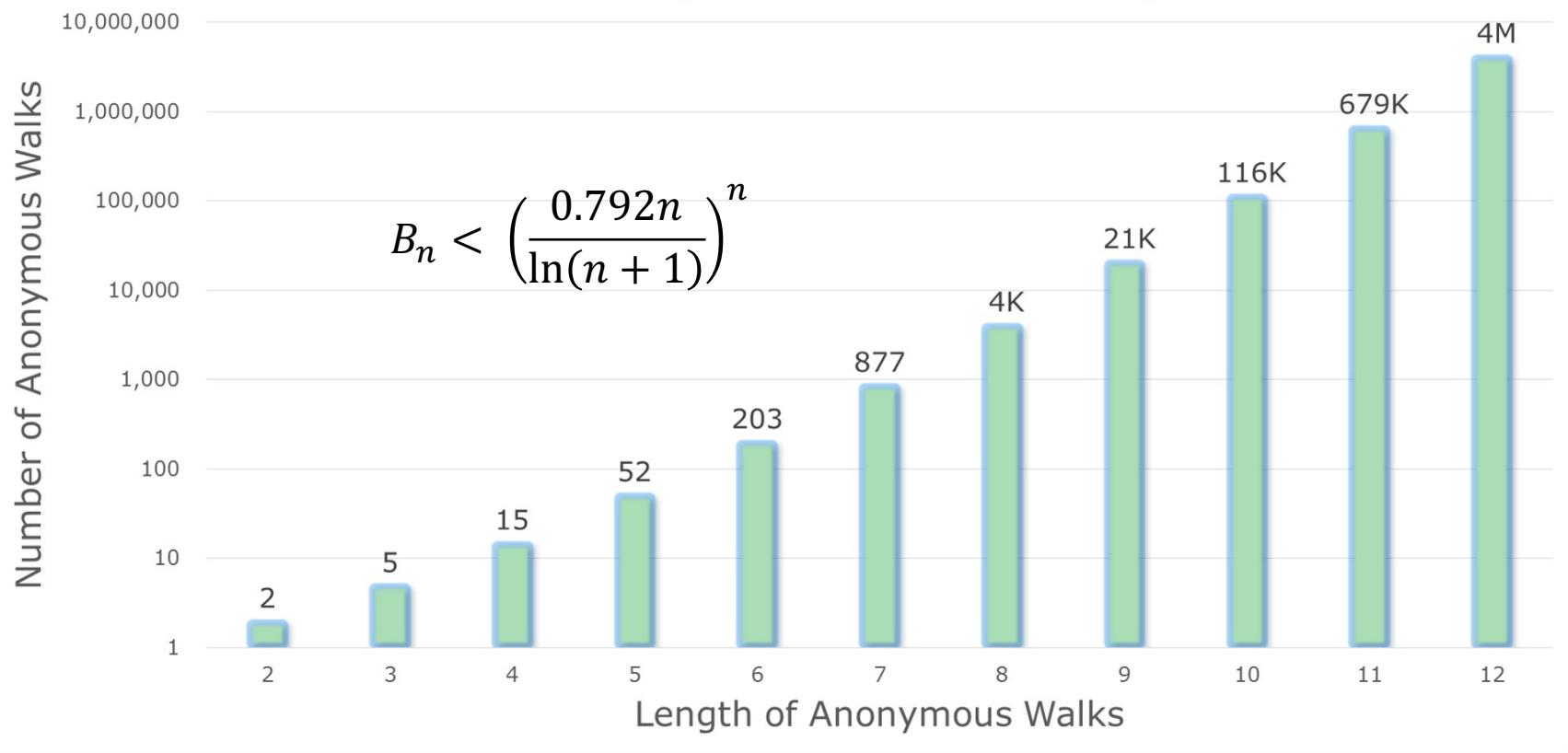
# Running time complexity

Theorem [This work]:

The number of possible anonymous walks  $|D_l|$  of length  $l$  in a graph is upper bounded by the Bell number  $B_{l-1}$ . The maximum for  $|D_l|$  is attained for complete graph with  $l$  nodes.

$$|D_l| \leq B_{l-1}$$

### Growth of Anonymous Walks with Length



# Combinatorial graph embeddings

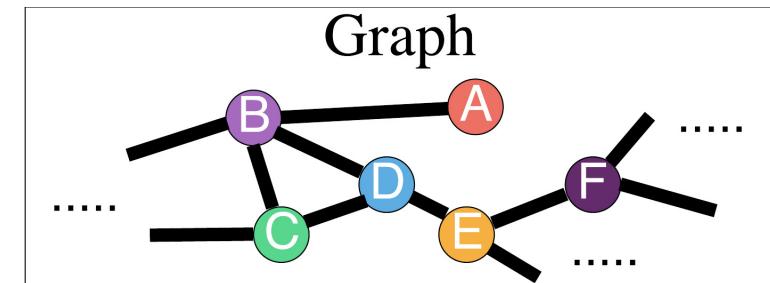
- $(a_1, a_2, \dots, a_\eta)$  all possible anonymous walks of length  $l$ .

$l$  is a parameter for embeddings.

## Combinatorial Graph Embedding

$$\text{AWE}(G) = (p(a_1), p(a_2), \dots, p(a_\eta))$$

Where  $p(a_i)$  is frequency of AW  $a_i$

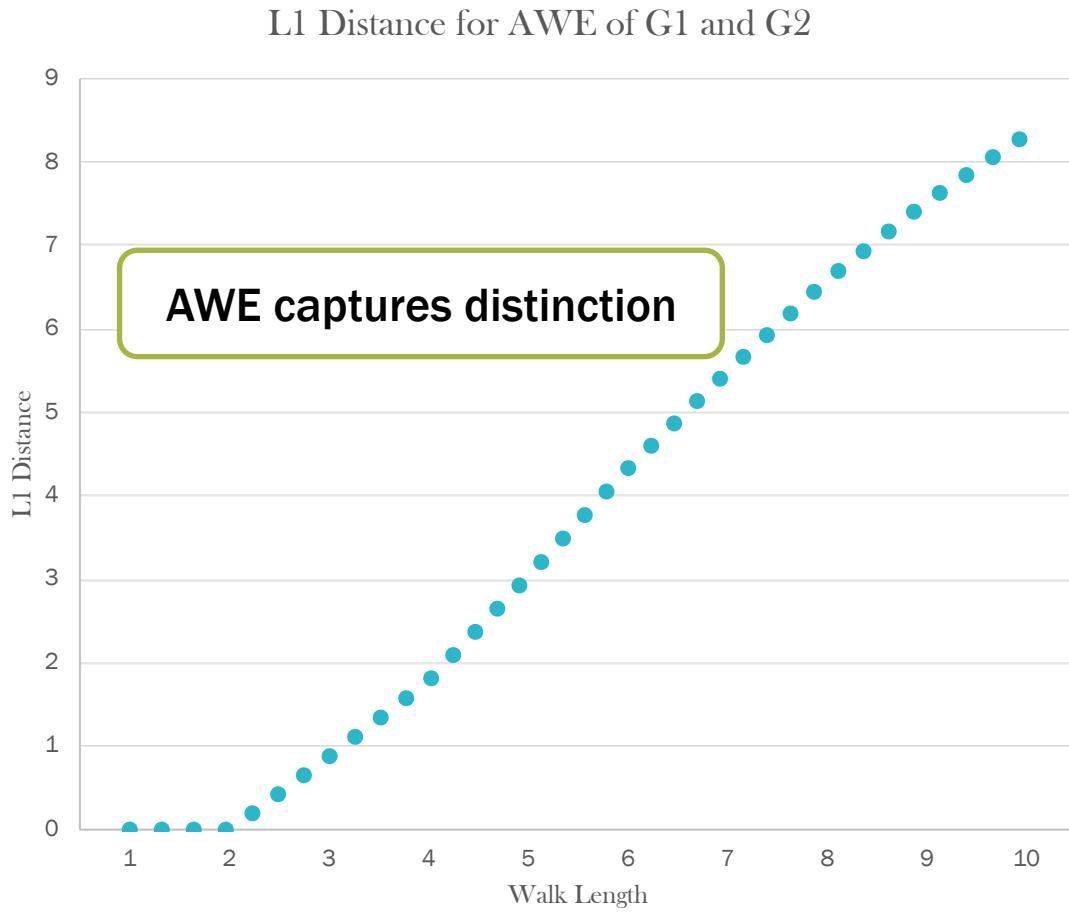


Graph

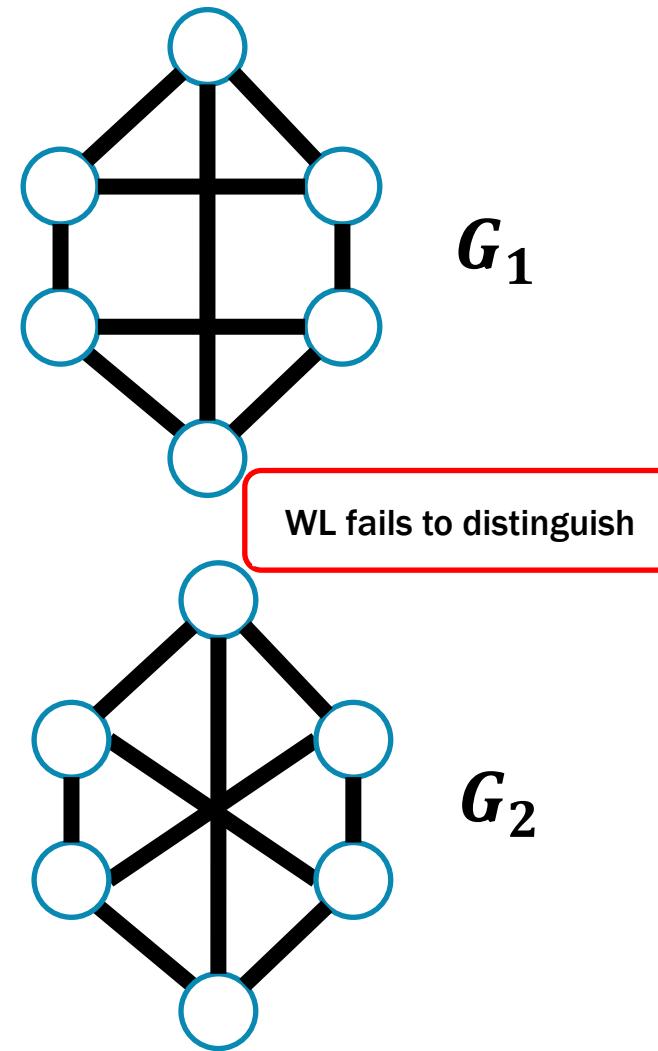
Embedding

$$l = 2(m + 1)$$

# Example of Isomorphism Property



$$\|AWE(G)\|_1 = 1$$



# Resolving computation complexity

Finding all possible anonymous walks  $(a_1, a_2, \dots, a_\eta)$  of length  $l$  can be expensive.

Instead, we can sample  $\mu$  anonymous walks and compute embeddings from them.

Can we guarantee the quality of sampled embeddings?

# Sampling inequality

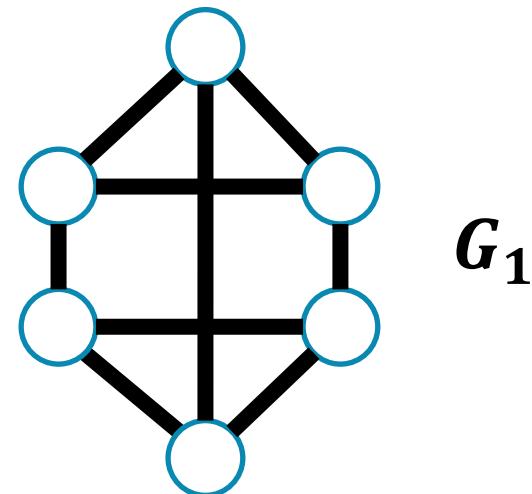
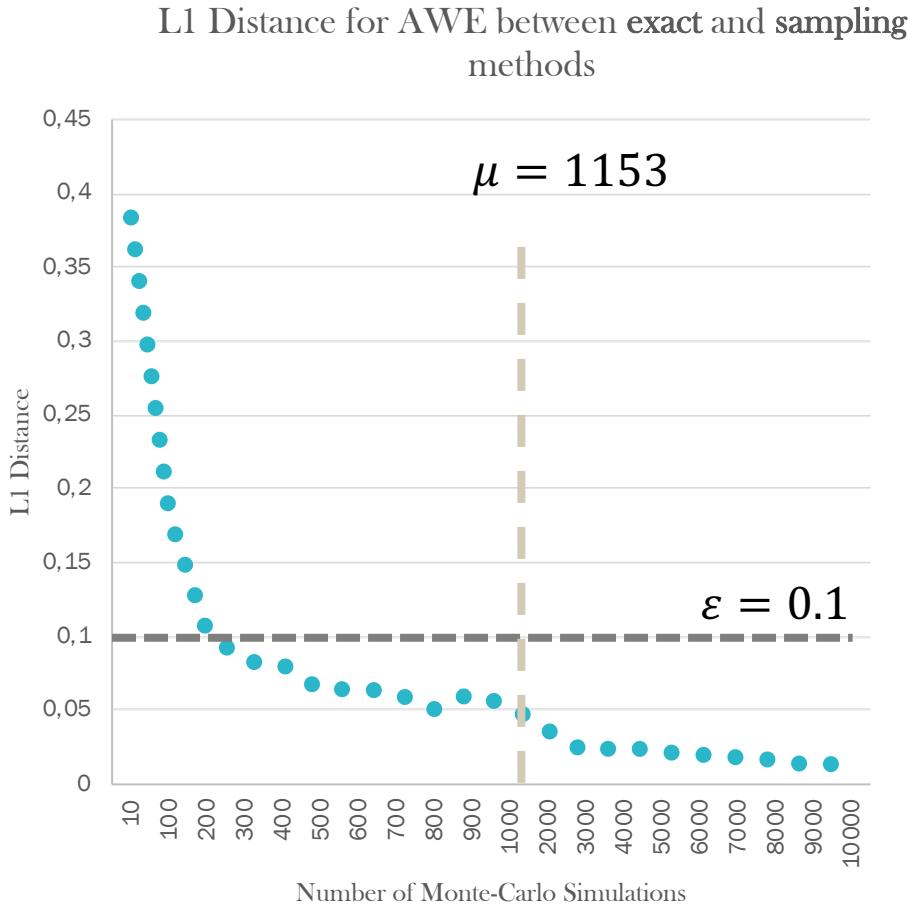
Theorem [This work]:

- $D_l$  is *true* distribution of anonymous walks of length  $l$  in a graph  $G$
- $\widehat{D}_l$  is *sampling* distribution of  $\mu$  anonymous walks of length  $l$  from graph  $G$ .

For all  $\varepsilon > 0$  and  $\delta \in [0,1]$ , the number of samples  $\mu$  to satisfy  $P\left(\|D_l - \widehat{D}_l\|_1 \geq \varepsilon\right) \leq \delta$  equals to:

$$\mu = \lceil \frac{2}{\varepsilon^2} (\log(2^\eta - 2) - \log(\delta)) \rceil$$

## Example of Sampling Bound Property



$\varepsilon = 0.1$

$\delta = 0.1$

$$P \left( \|D_l - \widehat{D}_l\|_1 \geq 0.1 \right) \leq 0.1$$

$$\mu = \left\lceil \frac{2}{\varepsilon^2} (\log(2^\eta - 2) - \log(\delta)) \right\rceil = 1153$$

# Neural network embeddings

Sample a corpus of anonymous walks that start from the same node.

Initialize randomly graph embedding  $d$  and a matrix of embeddings  $W$  for each anonymous walk.

Maximize the average log probability of observing the corpus

# Neural network embeddings

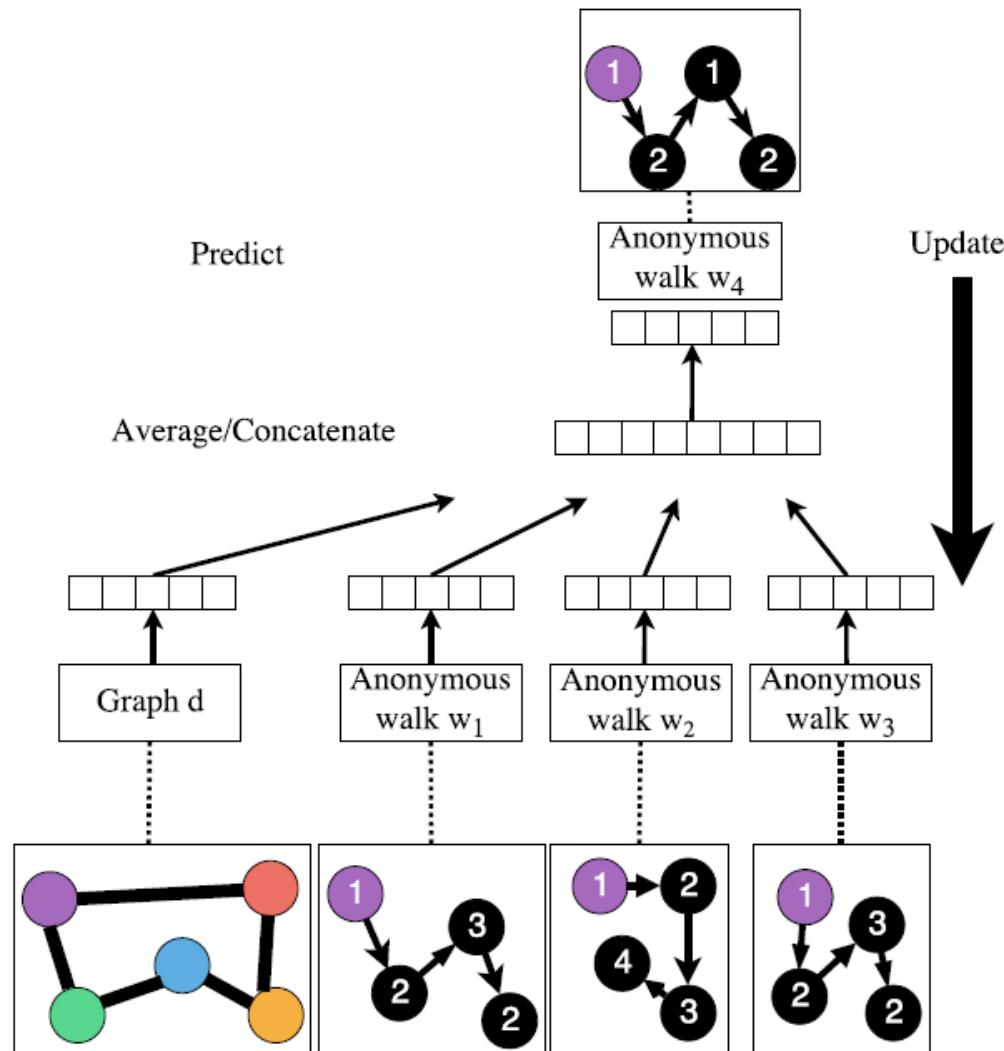
We optimize the objective

$$\frac{1}{T} \sum_{t=\Delta}^{T-\Delta} \log p(w_t | w_{t-\Delta}, \dots, w_{t+\Delta}, d) \mapsto_{W,d} \max$$

where

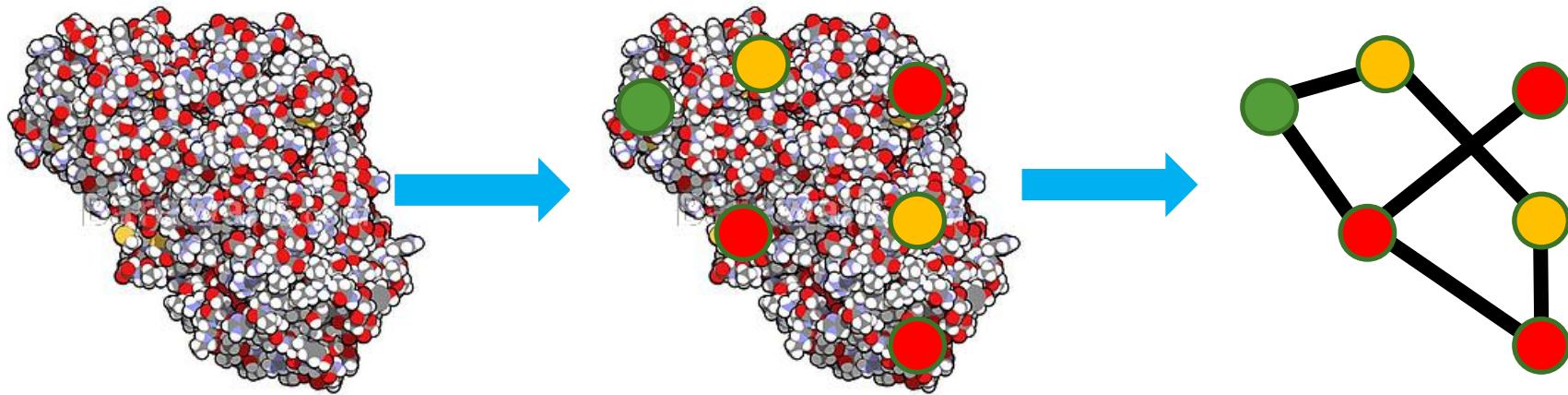
$$p(w_t | w_{t-\Delta}, \dots, w_{t+\Delta}, d) = \frac{e^{y(w_t)}}{\sum_{i=1}^{\eta} e^{y(w_i)}}$$

is the softmax probability of seeing anonymous walk in a graph and  $y(w_i) = \langle w_i, [\frac{1}{2\Delta} \sum_{j=t+\Delta}^{t+\Delta} w_j ; d] \rangle$  is similarity between walk  $w_i$  and its neighborhood.



## 2. Graph Classification

# Protein function prediction



Nodes: SSEs (helices, sheets, turns)

Edges: sequential or structural neighbors

Labels: length between atoms, polarity of SSE, etc. [1]

[1] Protein function prediction via graph kernels, Borgwardt et al. 2005

# Graph classification problem

Given

- $T = \{(G_i, y_i)\}_1^N$  - train graph data set
- $Q = \{(G_i, y_i)\}_1^M$  - test graph data set

Using the train set  $T$ , find a function  $f \in F = \{\phi: G \mapsto Y\}$  such that

$$acc = \frac{1}{M} \sum_1^M [f(G_i) = y_i] \mapsto max$$

# Classification pipeline

Prepare k-fold Cross-Validation splits.



Train model and choose the best hyperparameters for embeddings and classification models.



Evaluate the best model on test instances.

## Embedding parameters

- Length of a walk
- Window size
- Embedding size

## Model parameters

- Penalty term C
- Kernel type (e.g. Gaussian, Polynomial)
- Batch size

# Datasets

	Dataset	Source	Graphs	Classes (Max)	Nodes Avg.	Edges Avg.
[1]	COLLAB	Social	5000	3 (2600)	74.49	4914.99
	IMDB-B	Social	1000	2 (500)	19.77	193.06
	IMDB-M	Social	1500	3 (500)	13	131.87
	RE-B	Social	2000	2 (1000)	429.61	995.50
[2]	RE-M5K	Social	4999	5 (1000)	508.5	1189.74
	RE-M12K	Social	12000	11 (2592)	391.4	913.78
	Enzymes	Bio	600	6 (100)	32.6	124.3
[3]	DD	Bio	1178	2 (691)	284.31	715.65
	Mutag	Bio	188	2 (125)	17.93	19.79

[1] Deep Graph Kernels, Yanardag et al. 2012.

[2] Protein function prediction via graph kernels, Borgwardt, 2005.

[3] Distinguishing enzyme structures from non-enzymes without alignments, Dobson & Doig, 2003

# Evaluation accuracy

	Algorithm	IMDB-M	IMDB-B	COLLAB
Neural	DGK	$44.55 \pm 0.52$	$66.96 \pm 0.56$	$73.09 \pm 0.25$
Kernel	WL	$49.33 \pm 4.75$	<b><math>73.4 \pm 4.63</math></b>	<b><math>79.02 \pm 1.77</math></b>
	GK	$43.89 \pm 0.38$	$65.87 \pm 0.98$	$72.84 \pm 0.28$
	ER	OOM	$64.00 \pm 4.93$	OOM
	kR	$34.47 \pm 2.42$	$45.8 \pm 3.45$	OOM
	AWE (NN)	<b><math>51.54 \pm 3.61</math></b>	<b><math>74.45 \pm 5.83</math></b>	<b><math>73.93 \pm 1.94</math></b>
Ours	AWE (GK)	<b><math>51.58 \pm 4.66</math></b>	$73.13 \pm 3.28$	$70.99 \pm 1.49$

[1] Deep Graph Kernels, Yanardag et al. 2012

[2] Weisfeiler-Lehman Graph Kernels, Shervashidze et al. 2011

[3] Efficient graphlet kernels for large graph comparison, Shervashidze et al. 2009

[4] Graph Kernels, Vishwanathan et al. 2010

# Evaluation accuracy

	Algorithm	RE-B	RE-M5K	RE-M12K
Neural	DGK	$78.04 \pm 0.39$	$41.27 \pm 0.18$	$32.22 \pm 0.10$
Kernel	WL	$81.1 \pm 1.9$	$49.44 \pm 2.36$	$38.18 \pm 1.3$
	GK	$65.87 \pm 0.98$	$41.01 \pm 0.17$	$31.82 \pm 0.08$
	ER	OOM	OOM	OOM
	kR	OOM	OOM	OOM
	AWE (NN)	<b><math>87.89 \pm 2.53</math></b>	<b><math>50.46 \pm 1.91</math></b>	<b><math>39.20 \pm 2.09</math></b>
Ours	AWE (GK)	<b><math>82.97 \pm 2.86</math></b>	<b><math>54.74 \pm 2.93</math></b>	<b><math>41.51 \pm 1.98</math></b>

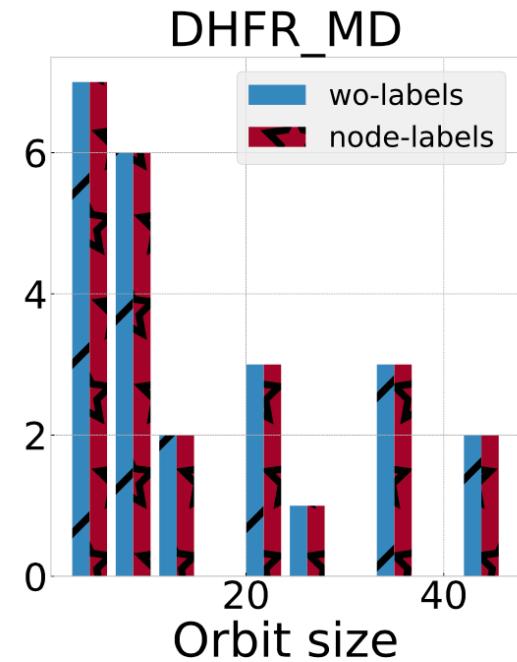
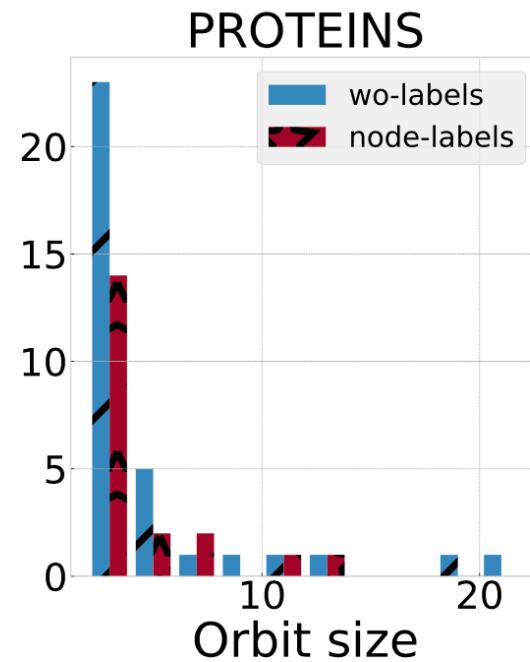
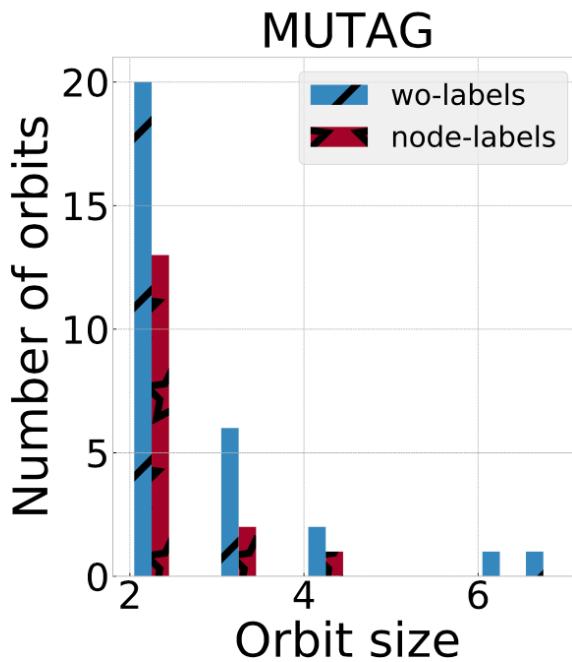
# Evaluation accuracy

	Algorithm	Enzymes	DD	Mutag
Neural	DGK	$27.08 \pm 0.79$	—	$82.66 \pm 1.45$
Kernel	WL	<b><math>53.15 \pm 1.14</math></b>	<b><math>77.95 \pm 0.70</math></b>	$80.72 \pm 3.00$
	GK	$32.70 \pm 1.20$	<b><math>78.45 \pm 0.26</math></b>	$81.58 \pm 2.11$
	ER	$14.97 \pm 0.28$	OOM	$71.89 \pm 0.66$
	kR	$30.01 \pm 1.01$	OOM	$80.05 \pm 1.64$
Ours	AWE (GK)	<b><math>35.77 \pm 5.93</math></b>	$71.51 \pm 4.02$	<b><math>87.87 \pm 9.76</math></b>

### 3. Learning to correct predictions for Graph Classification with noisy labels

# Graph classification and graph isomorphism

- Even theoretically expressive graph classification models do not fully exploit the isomorphism between the graphs in the data set at hand
- Most of the data sets used in graph classification have isomorphic instances
- In some datasets the fraction of the unique non-repeating graphs is as low as 20% of the total size, while those graphs that are isomorphic often have non-unique target label



data set	Size, $N$	Isomorphic, %	Noise, %
Letter-low	2250	99.78	8.72
IMDB-MULTI	1500	80.8	6.39
AIDS	2000	62.95	0.13
COX2	467	60.6	0.6
IMDB-BINARY	1000	57.9	0.67
MUTAG	188	42.02	0.49
COLLAB	5000	21.54	0.11
PROTEINS	1113	13.57	0.1
NCI1	4110	12.73	0.01
NCI109	4127	12.58	0.01
ENZYMES	600	1.67	0
REDDIT-BINARY	2000	0.2	0
REDDIT-MULTI-12K	11929	0.14	0

Table 1: The number of isomorphic graphs and noisy target labels in popular graph data sets. The last column indicates the proportion of graphs that contain at least one isomorphic copy in the same data set, which has a different target label.

---

## Algorithm 1 Error-Correcting Graph Classification

---

```
1: Get predictions of the model on  $Y_{test}$ ;  
2: # Error-Correcting Phase  
3: for  $G_i \in Y_{test}$  do  
4:   Initialize orbit  $O_i = [G_i]$   
5:   for  $G_j \in Y_{train}$  do  
6:     if  $G_i \cong G_j$  then  
7:       Add  $G_j$  to  $O_i$   
8:     end if  
9:   end for  
10:  Predict  $y_i = \arg \max_y \mathbb{P}(y|G_i, O_i)$ .  
11: end for
```

---

# Error-correction performance

Graph Isomorphism Network (GIN)

Weisfeiler-Lehman graph kernel (WL)

Vertex histogram kernel (V)

- Note that GIN and WL kernels have sufficient expressiveness to recognize most of the isomorphic graphs in real data sets and hence one may hypothesize that the error-correcting phase is unnecessary for these methods.
- This hypothesis is false and they *do benefit* from additional correction of the predictions
- We add a prefix EC (error correction) for each model that follows Algorithm 1 and a prefix ECU that has an Error-Correcting phase only for the orbits with a Unique target label

Table 2: Mean classification accuracy for test sets  $Y_{test}$  and  $Y_{iso}$  (in brackets) in 10-fold cross-validation. Top-1 result for each type of algorithm in bold. Error correction consistently achieves increase of accuracy for each model.

	MUTAG	IMDB-B	IMDB-M	COX2	AIDS	PROTEINS
GIN	0.829 (0.840)	0.737 (0.733)	0.501 (0.488)	0.82 (0.872)	0.996 (0.998)	0.737 (0.834)
GIN-EC	0.856 (0.847)	0.737 (0.731)	0.499 (0.486)	0.795 (0.83)	0.996 (0.999)	0.729 (0.709)
GIN-ECU	<b>0.867</b> (1.000)	<b>0.756</b> (1.000)	<b>0.522</b> (1.000)	<b>0.838</b> (1.000)	<b>0.996</b> (1.000)	<b>0.742</b> (1.000)
WL	0.862 (0.867)	0.734 (0.990)	0.502 (0.953)	0.800 (0.974)	0.993 (0.999)	0.747 (0.950)
WL-EC	0.870 (0.838)	0.724 (0.715)	0.495 (0.487)	0.794 (0.844)	0.994 (0.999)	0.740 (0.742)
WL-ECU	<b>0.907</b> (1.000)	<b>0.736</b> (1.000)	<b>0.504</b> (1.000)	<b>0.810</b> (1.000)	<b>0.994</b> (1.000)	<b>0.749</b> (1.000)
V	0.836 (0.902)	0.707 (0.820)	0.503 (0.732)	0.781 (0.966)	0.994 (0.997)	0.726 (0.946)
V-EC	0.827 (0.844)	0.724 (0.728)	0.496 (0.481)	0.768 (0.852)	0.996 (0.999)	0.719 (0.741)
V-ECU	<b>0.859</b> (1.000)	<b>0.750</b> (1.000)	<b>0.517</b> (1.000)	<b>0.794</b> (1.000)	<b>0.996</b> (1.000)	<b>0.729</b> (1.000)

Table 3: Mean classification accuracy for test sets  $Y_{test}$  in 10-fold cross-validation on data sets with no noise. Each cell is  $a \pm b(c)$ , where  $a$  is mean accuracy,  $b$  is standard deviation, and  $c$  is absolute difference of accuracy between original (noisy) and new data sets.

Data	DGCNN	DiffPool	GIN	Baseline
MUTAG	80.59±11 (-1.60)	81.48±8 (-1.75)	<b>84.93±9</b> (-0.55)	71.08±6 (5.84)
IMDB-B	70.53±6 (-1.23)	72.96±5 (-2.43)	<b>74.84±3</b> (-4.84)	52.47±5 (0.83)
IMDB-M	46.11±4 (-0.17)	45.90±6 (-0.28)	<b>46.82±6</b> (-0.93)	41.11±6 (-7.07)
COX2	69.62±7 (9.05)	<b>75.11±9</b> (7.19)	73.83±5 (6.92)	68.30±11 (8.31)
BZR	79.33±6 (0.72)	<b>81.64±7</b> (1.48)	80.69±5 (3.28)	78.43±7 (4.12)
PTC_FM	<b>63.96±5</b> (-3.41)	59.46±6 (0.37)	57.76±6 (2.50)	60.89±6 (-2.43)
PTC_MM	<b>66.83±3</b> (-4.23)	59.97±9 (5.00)	60.36±6 (4.71)	58.66±7 (2.66)
PTC_MR	<b>60.72±4</b> (-4.48)	57.69±5 (-0.98)	56.01±8 (-0.44)	59.72±7 (-3.17)
PTC_FR	<b>66.01±2</b> (-2.76)	63.51±2 (0.30)	61.39±6 (4.61)	60.86±9 (1.73)

## 4. fMRI-based Depression Diagnostics

# fMRI-based Depression Diagnostics

- 50 million people worldwide have epilepsy, making it one of the most common neurological diseases;
- 300 million people globally affected with depression as common psychiatric disorder;
- 37% of **epilepsy** patients suffer from **depression**.

The biomarkers of major depressive disorder as the **concomitant disease** are still unclear as well as their relation to epilepsy.

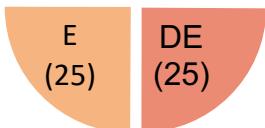
# Data set

Medical partner Research and Clinical Center of Neuropsychiatry named after Solovyov.

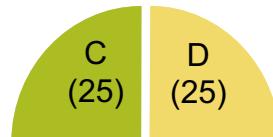
Dataset of **100 functional 1.5 T MRI sets** in four groups (mean age:  $33 \pm 9$  years).

Classification tasks:

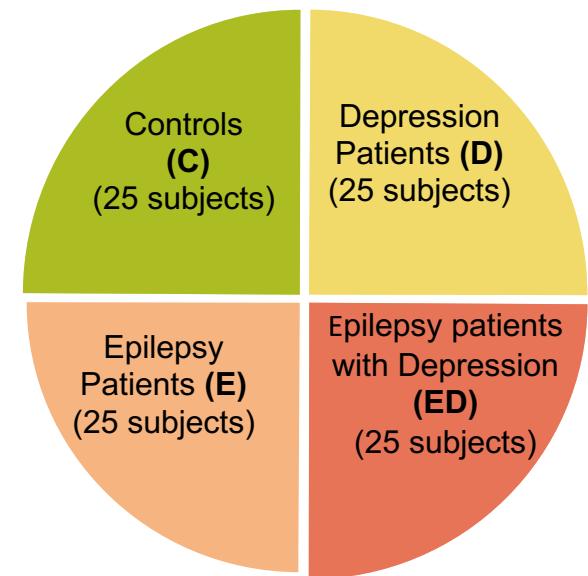
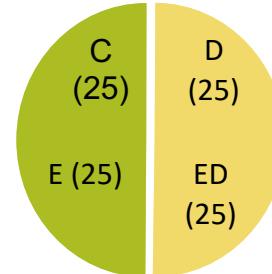
E vs DE:



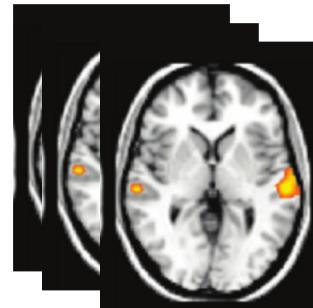
D vs C:



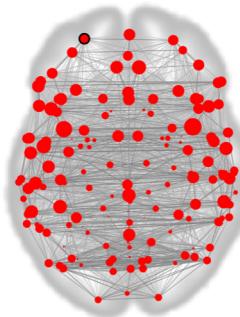
D vs NoD:



# Data content



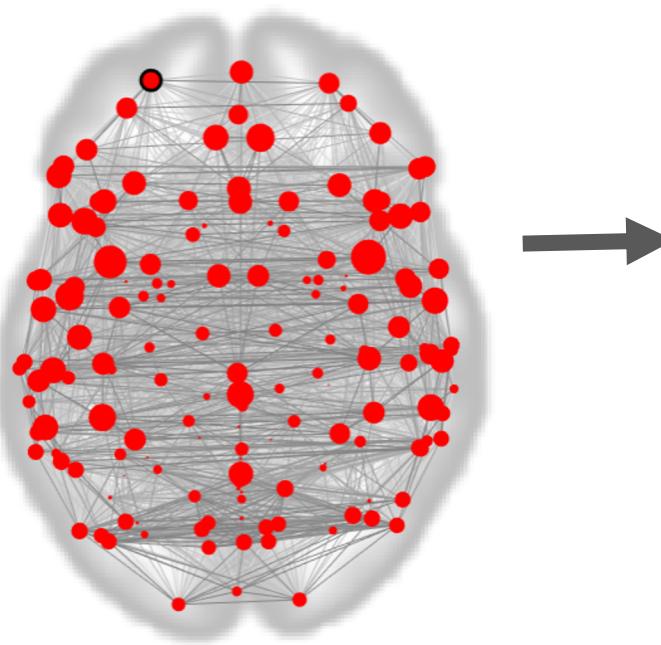
Preprocessing:  
Time series extraction



Resting state T2\* EPI fMRI  
[64, 64, 30] voxels  
133 time steps

- 117 regions (Automated Anatomical Labeling atlas)
- Time series for each region
- Correlation matrix

# Connectivity matrix



	<b>ROI_1</b>	<b>ROI_2</b>	...	<b>ROI_116</b>	<b>ROI_117</b>
<b>ROI_1</b>	1	0.2418	...	0.034	0.086
<b>ROI_2</b>	0.2418	1	...	0.053	0.247
...	...	...	1	...	...
<b>ROI_116</b>	0.034	0.053	...	1	0.017
<b>ROI_117</b>	0.086	0.247	...	0.017	1

# Feature-based graph classification

## 1. Feature based methods

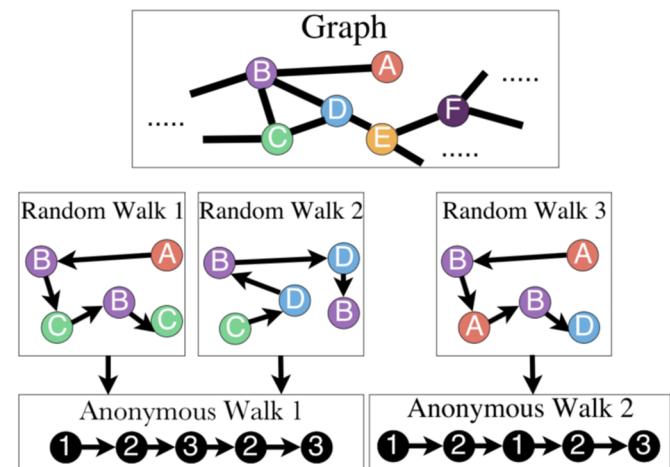
- Degree
- Average neighbor degree
- Degree centrality
- Betweenness/ closeness centrality
- Clustering coefficient
- Average path length
- Local/ global efficiency

$$E(G) = \frac{1}{n(n-1)} \sum_{i < j \in G} \frac{1}{d(i,j)}$$

where n is number of nodes - regions in a graph and  $d(i,j)$  is the length of the shortest path between nodes  $i$  and  $j$ .

## 2. Kernel methods

- Weisfeiler-Lehman kernel (WL)
- **Random walks (Anonymous walk embedding (AWE))**



# Evaluation

- 4 groups of patients: healthy (H), depression (D), epilepsy (E), depression + epilepsy (DE)
- Each group has 25 graphs
- Two classification tasks: DvsH and DvsDE

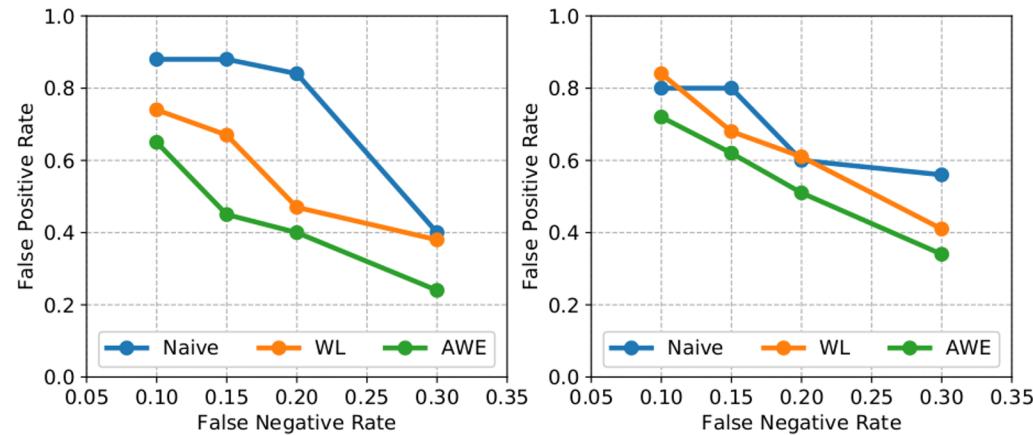


Fig.1. Classification results on DvsH and EvsDE tasks in terms of False Negative Rate vs. False Positive Rate performance curves.  
Left image: DvsH task. Right image: EvsDE task.

Task	Naïve	WL	AWE
DvsH	$73 \pm 15\%$	$78 \pm 15\%$	$80 \pm 12\%$
EvsDE	$67 \pm 15\%$	$75 \pm 14\%$	$76 \pm 16\%$

# Main results

- Proposed and justified a new graph representation that **strictly provides isomorphism property**;
- Designed two approaches for **approximate efficient computation of embeddings**;
- Demonstrated **superior quality of embeddings** in graph classification problem.

# Published Papers

- S. IVANOV, E. BURNAEV, “*ANONYMOUS WALK EMBEDDINGS*”, **PROCEEDINGS OF INTERNATIONAL CONFERENCE ON MACHINE LEARNING (ICML)**, 2018
- S. IVANOV, N. DURASOV, E. BURNAEV, “*LEARNING NODE EMBEDDINGS FOR INFLUENCE SET COMPLETION*”, **IEEE INTERNATIONAL CONFERENCE IN DATA MINING (ICDM) 2018 WORKSHOPS PROCEEDINGS**, 2018
- M. SHARAEV, A. ARTEMOV, A. BERNSTEIN, E. KONDRATYEVA, S. SUSHCHINSKAYA, E. BURNAEV, S. IVANOV, “*LEARNING CONNECTIVITY PATTERNS VIA GRAPH KERNELS FOR FMRI-BASED DEPRESSION DIAGNOSTICS*”, **IEEE INTERNATIONAL CONFERENCE IN DATA MINING (ICDM) 2018 WORKSHOPS PROCEEDINGS**, 2018
- S. IVANOV, S. SVIRIDOV, E. BURNAEV, “*UNDERSTANDING ISOMORPHISM BIAS IN GRAPH DATA SETS*”, **SUBMITTED, ARXIV 1910.12091**, 2019