# Explaining Identity-aware Graph Classifiers through the Language of Motifs

**Alan Perotti**
Paolo Bajardi
Francesco Bonchi
André Panisson
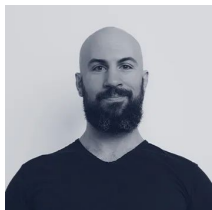
CENTAI

- Private company, working w/ universities

- Fundamental research

- Applied/industrial research projects

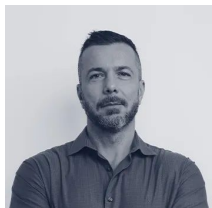- Research projects (e.g. Horizon Europe)

Disclaimer:
In this talk, the focus is heavily on the XAI side.

Happy to put you in touch w/ colleagues working on:

- Graph Machine Learning

- (explainable) Link Prediction

- Social Networks (opinion dynamics etc.)

- Graph Counterfactuals

- Hypergraphs, higher-order data

- Complex systems

- Old-school graph algorithms (MST, search, etc.)
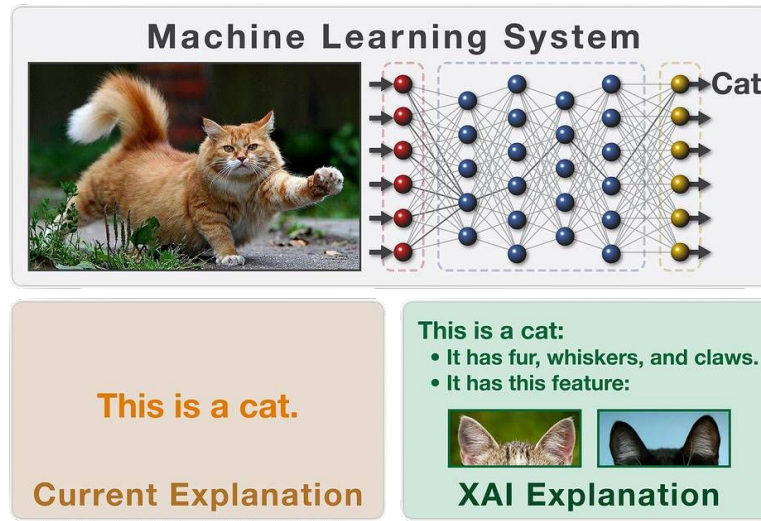
# The Black-Box Problem

Many modern ML models are hard to interpret and it is difficult to understand why they make a certain decision or recommendation. This might cause several problems:

- No trust from experts.
- Biased systems.
- Right for the wrong reasons.
- GDPR non-compliance.
- Adversarial Vulnerability.

Intuition: decorate model's output with additional information.
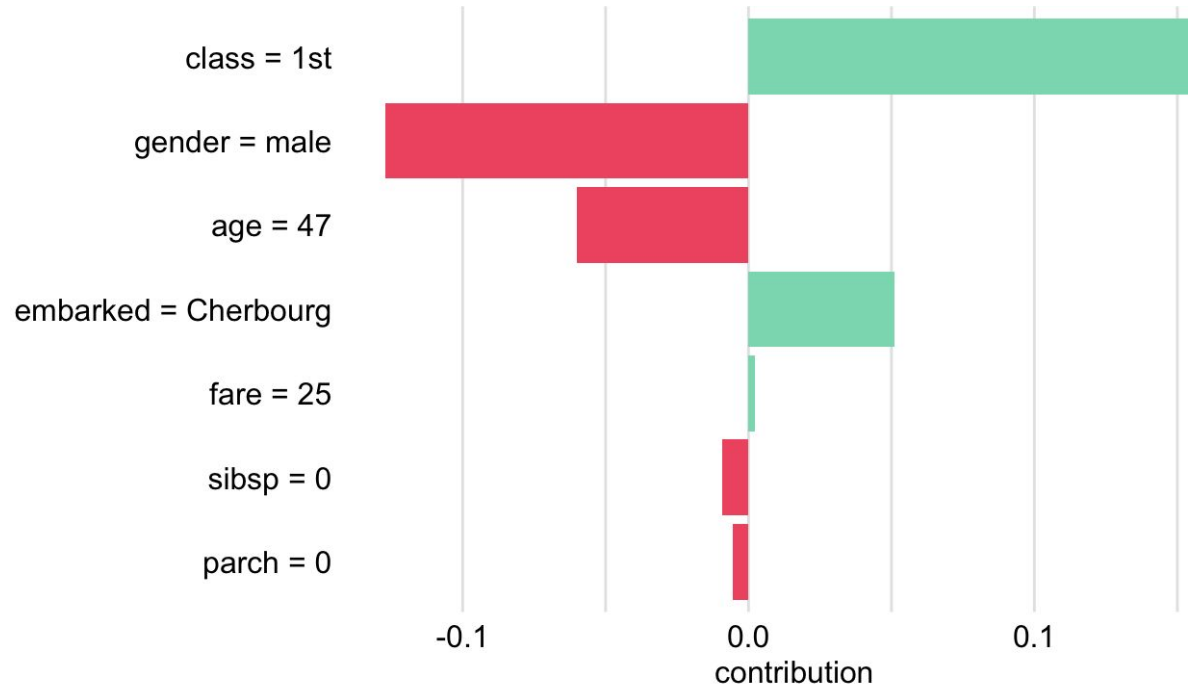


Very hot area: much fundamental research to be done, strong interest from private companies, European research calls, etc.
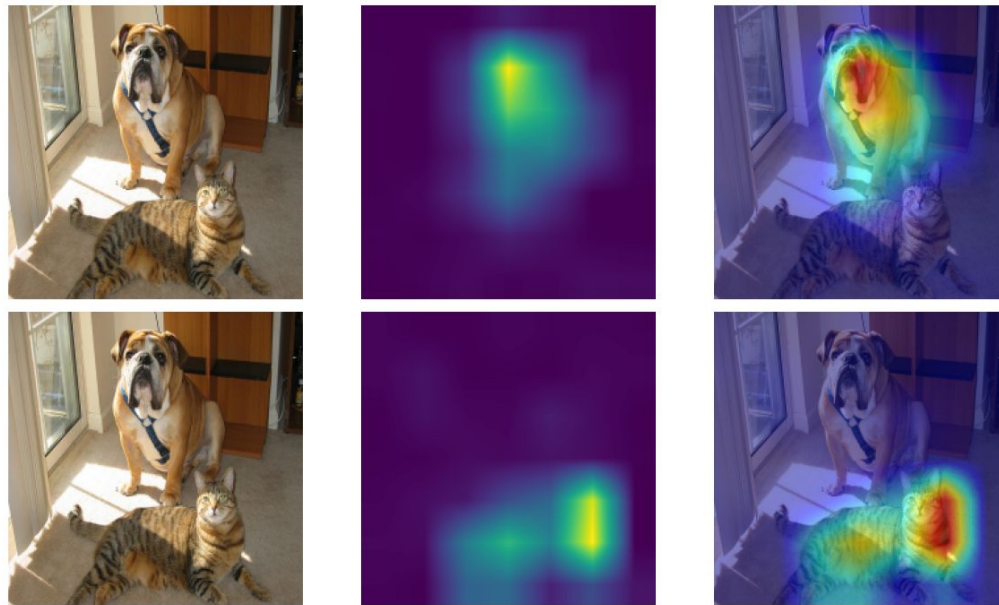
ML input language
VS
explanation language
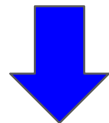
ML input language VS XAI explanation language

Attribution-based explanations on tabular data

# Heatmaps for Computer Vision ML models

Graph classification

## Graph classification and motifs (connected subgraphs)

G

M
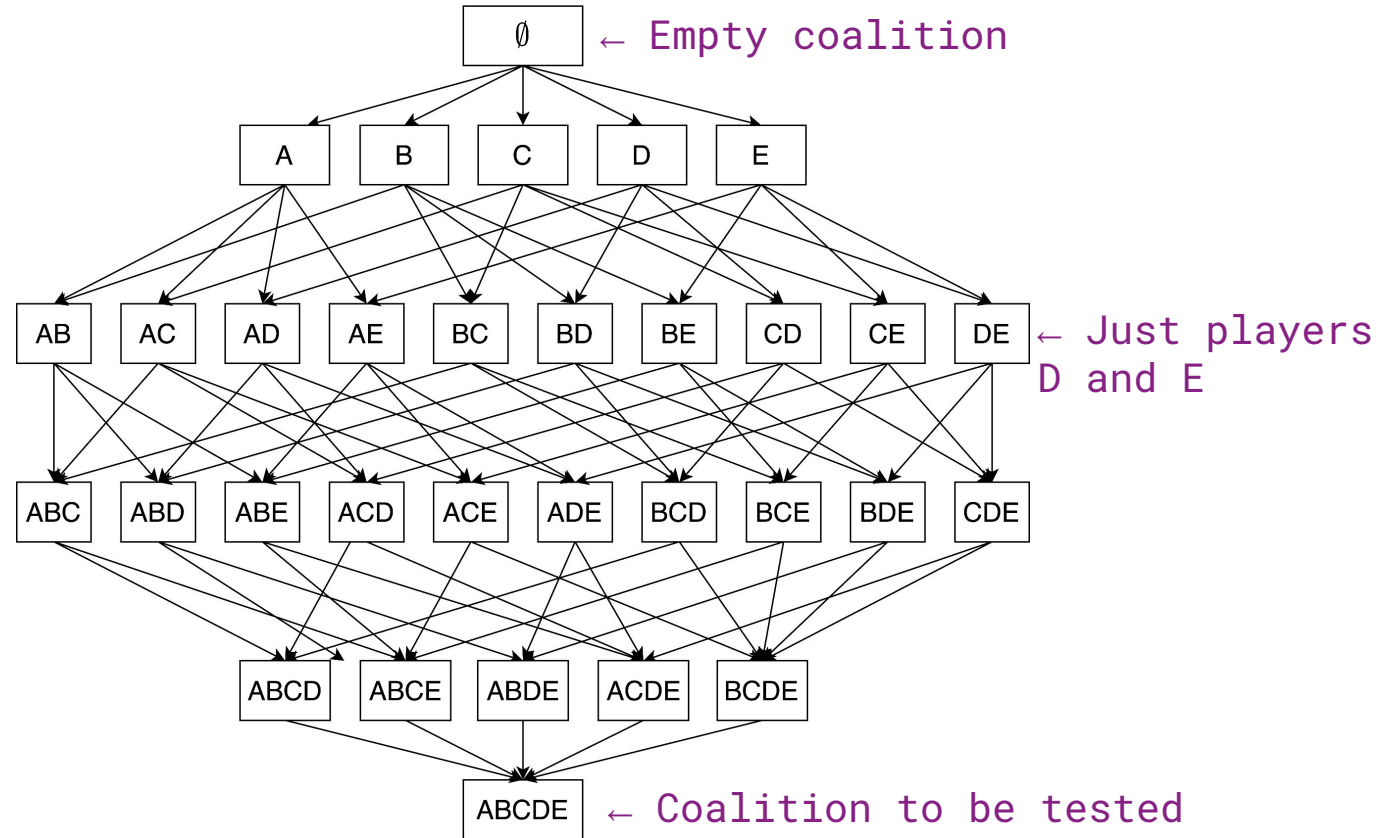
- Sub-graph of the induced complete graph, so might not occur.
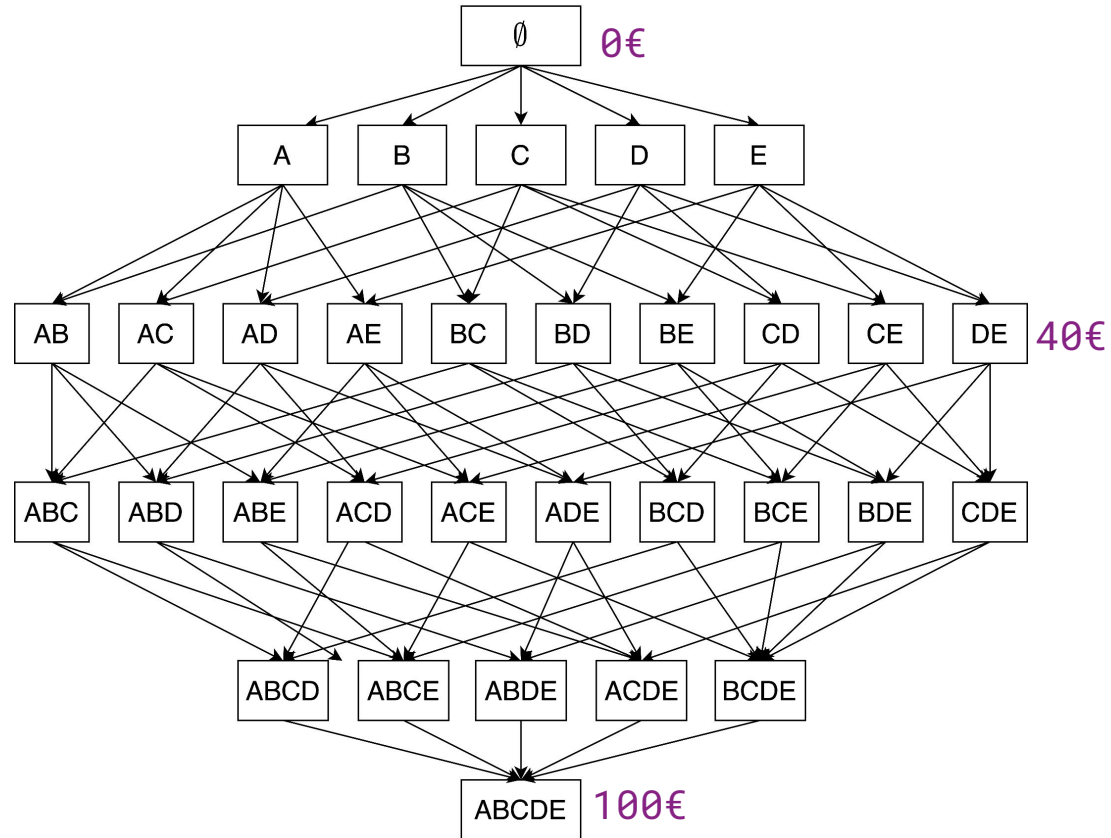
- Links unique nodes.

- Connected, but this can be relaxed.

Given a graph $G \in \mathcal{G}$, a black-box $B : \mathcal{G} \to [0, 1]$ and a set of motifs $\mathcal{M}$, the problem tackled in this paper is that of assigning an *explanation score* $\xi(G, B, M_i) \in [-1, 1]$ to each motif $M_i \in \mathcal{M}$, quantifying the impact of the motif in explaining the label $B(G)$: a value close to -1 means that $M_i$ is important in explaining $B(G) = 0$, a value close to 1 means that $M_i$ is important for $B(G) = 1$.

# Shapley value, 1951: a lattice of coalitions

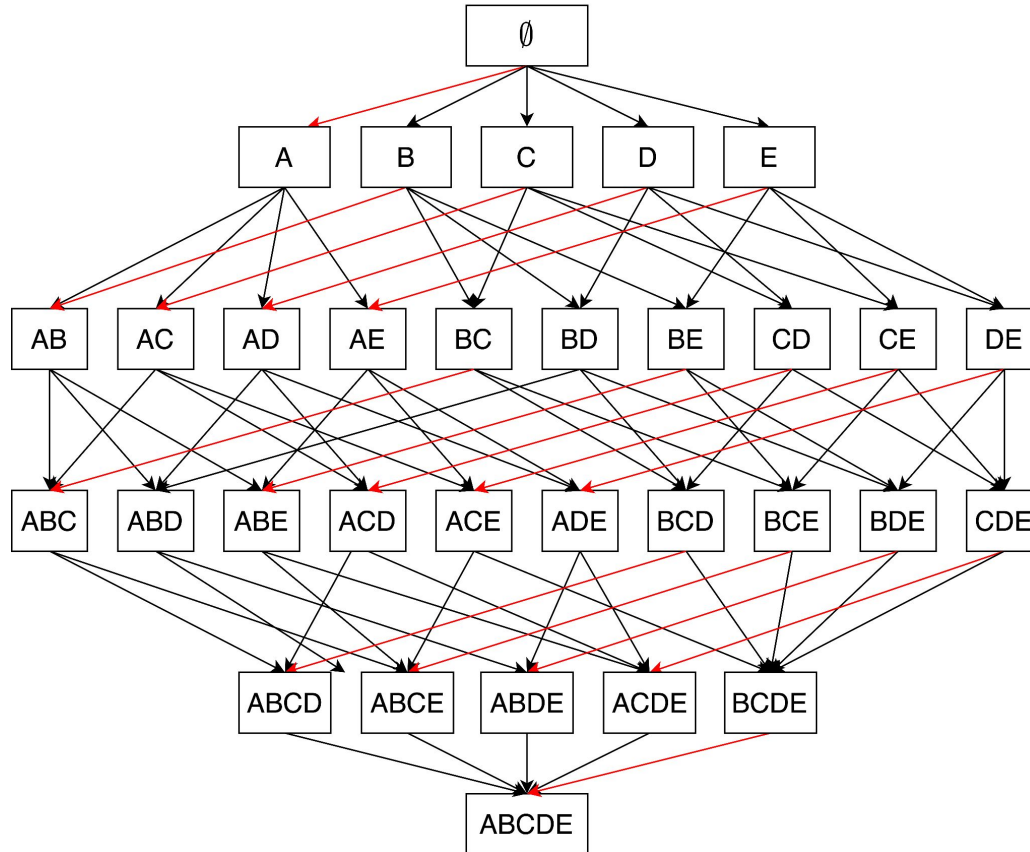Team ABCDE wins some money.
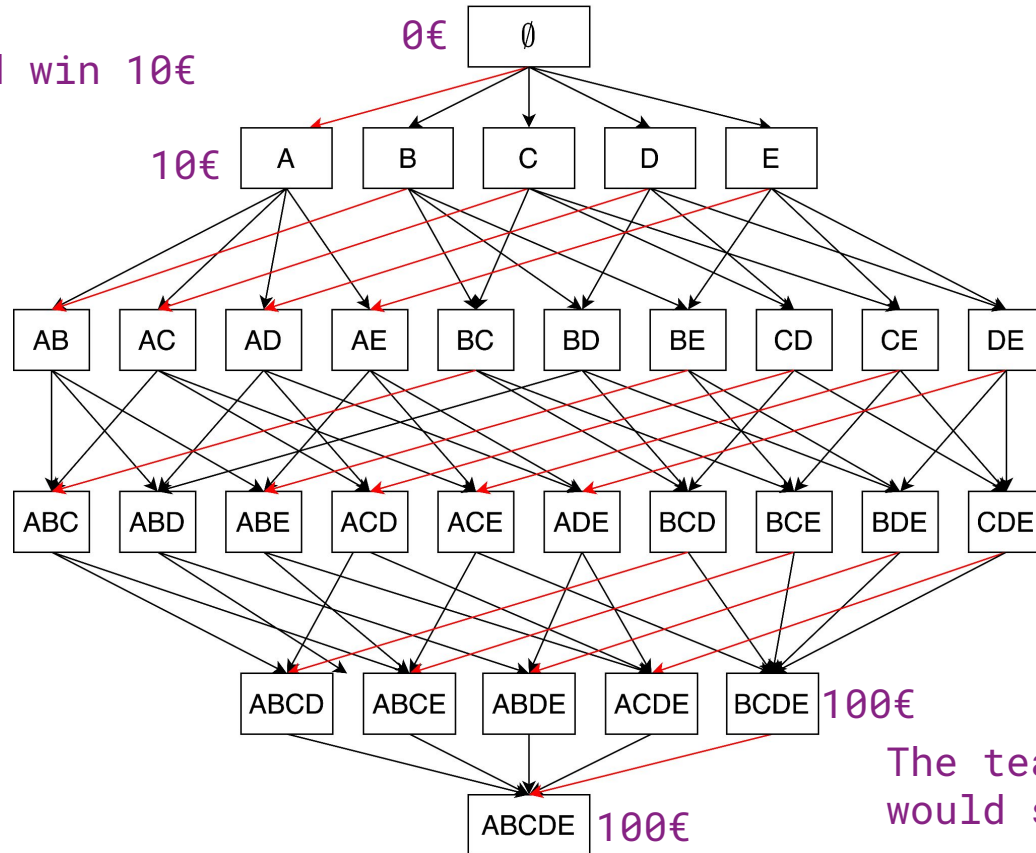How to distribute among players?

∅ ← Empty coalition

A  B  C  D  E

AB  AC  AD  AE  BC  BD  BE  CD  CE  DE ← Just players D and E

ABC  ABD  ABE  ACD  ACE  ADE  BCD  BCE  BDE  CDE

ABCD  ABCE  ABDE  ACDE  BCDE

ABCDE ← Coalition to be tested

Impact of player A

Shapley: every edge is adding a new player

A, alone, would win 10€

0€  ∅

10€  A  B  C  D  E

AB  AC  AD  AE  BC  BD  BE  CD  CE  DE

ABC  ABD  ABE  ACD  ACE  ADE  BCD  BCE  BDE  CDE

ABCD  ABCE  ABDE  ACDE  BCDE  100€

ABCDE  100€
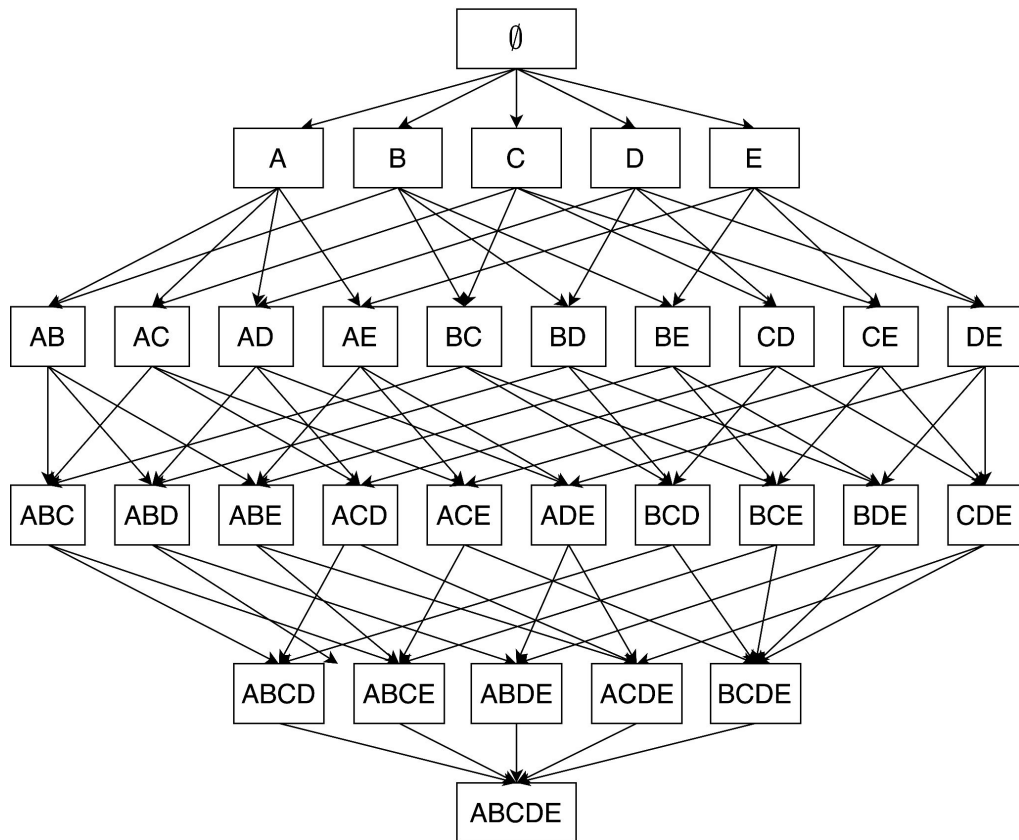
The team, without A, would still win 100€

Players → Features
Coalitions/teams → Data Points
Team value → ML output (continuous)

Removing a player from the team and measuring the change in team value → masking a feature value in the data point and measuring the change in ML output
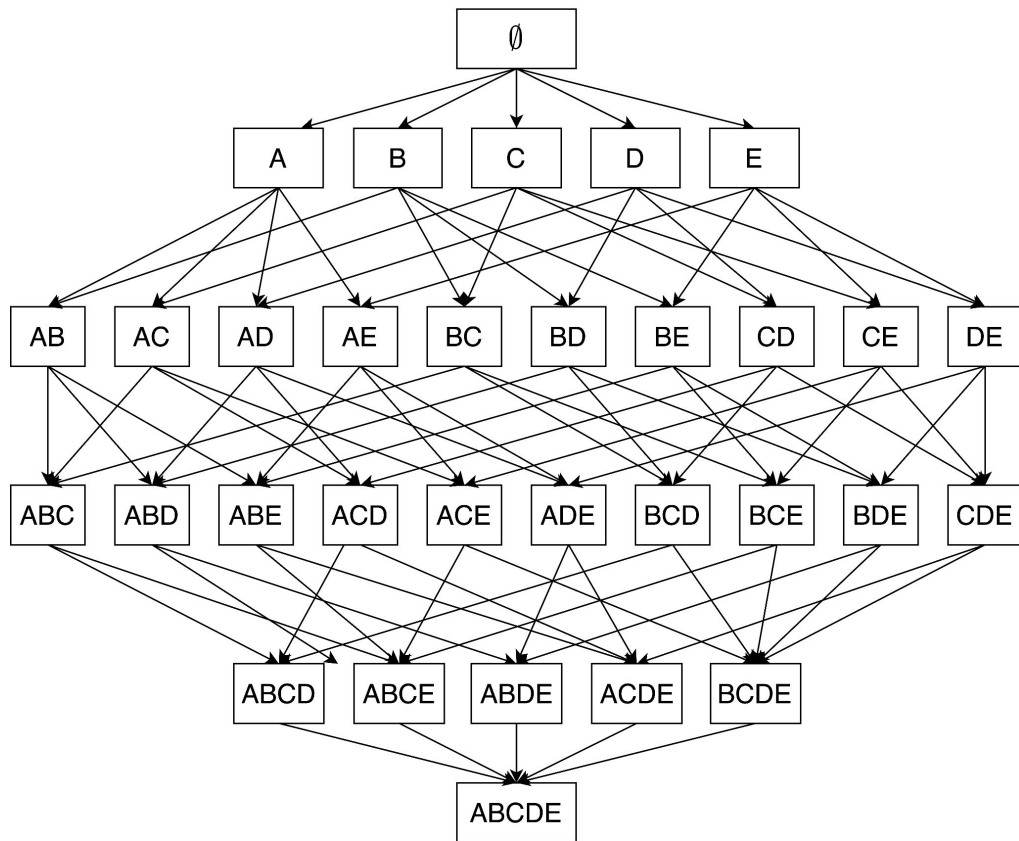
Pro: Arguably, XAI SOTA

Con: TWO devils in the details: scalability/approximation, feature removal

Players → MOTIFS
Coalitions/teams → Data Points
Team value → ML output (continuous)

Removing a player from the team and measuring the change in team value → masking a motif in the graph data point and measuring the change in ML output
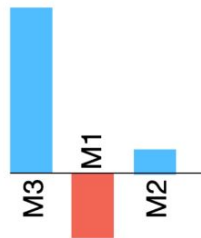
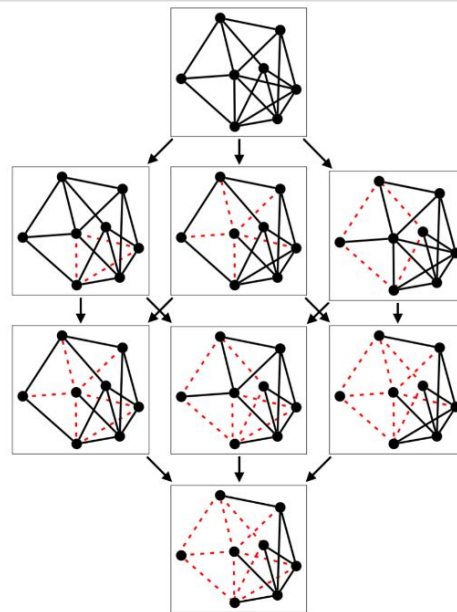(we inherit the same weaknesses)
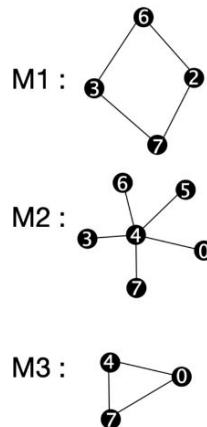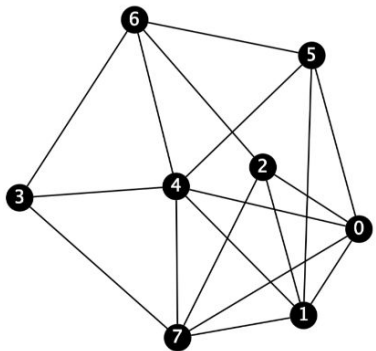
# GraphSHAP pipeline

Black-Box graph classifier *B*

Graph *G* to be classified and explained

User-defined model-agnostic explanation space
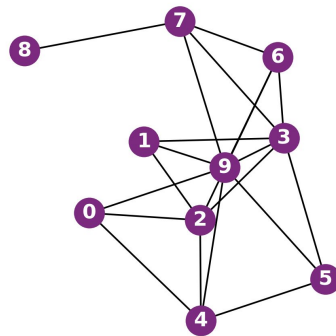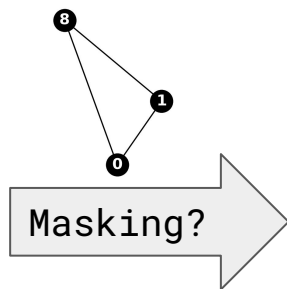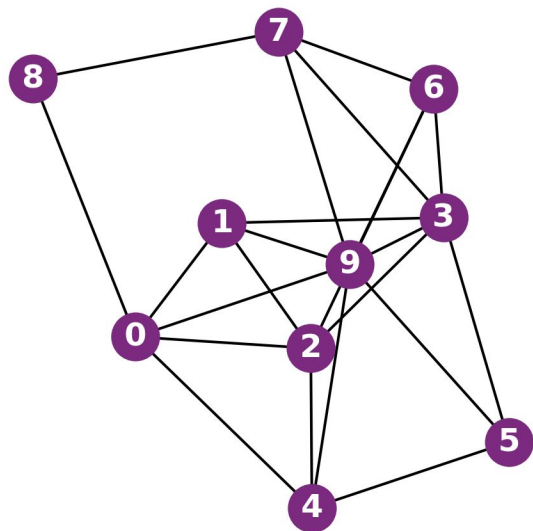
Marginal contribution of explainable features

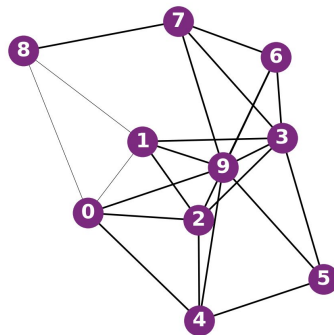Shapley-based explanation on selected motifs

M1 :

M2 :

M3 :

How do we translate the concept of *removing a player* in our graph-ML setting?

SHAP introduces the concept of *background dataset*, and copies values from other data points
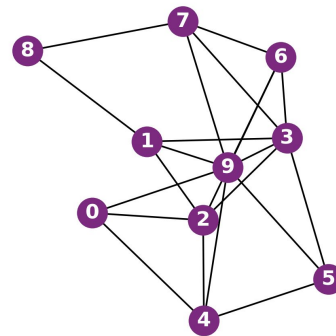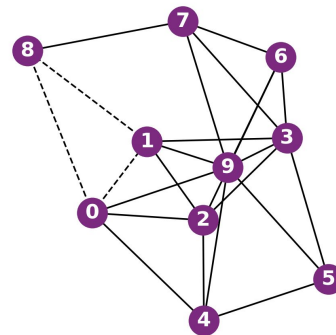
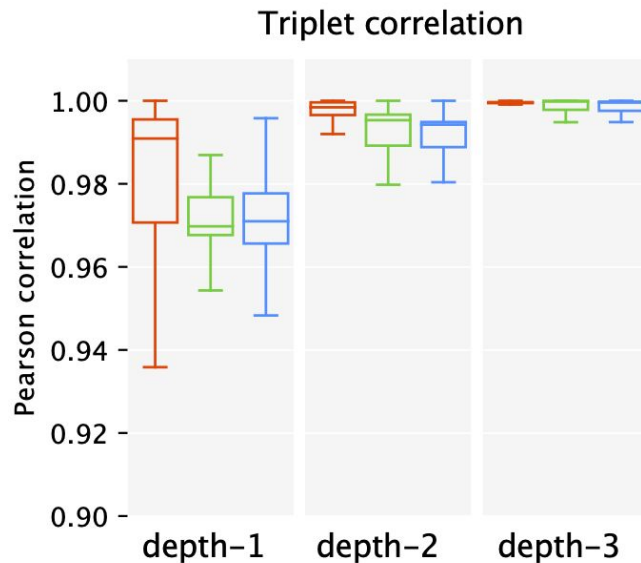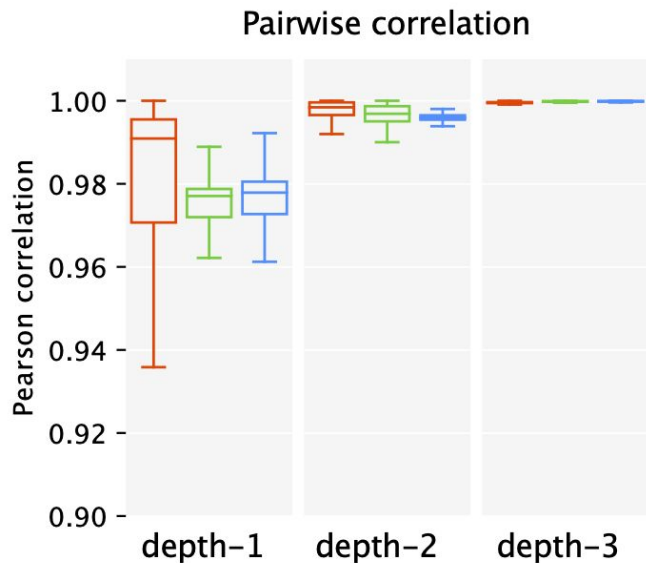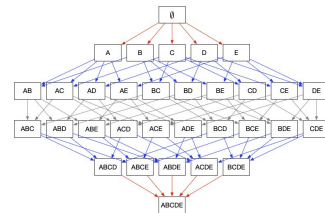Masking?

Remove

Toggle

Weigh

Sample

How do we deal with the lattice's exponential computational complexity (wrt the number of features)?

SHAP introduces the concept of *budget*, and samples the budget according to heavily engineered heuristics

We found a strong approximation (with respect to the full Shapley lattice) in the first Shapley layer.



Pairwise correlation

Triplet correlation

no correlation    mild correlation    strong correlation

## Welcome to the Autism Brain Imaging Data Exchange!

**ABIDE**

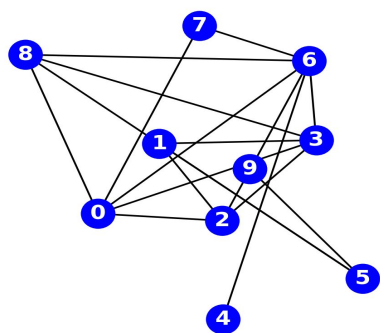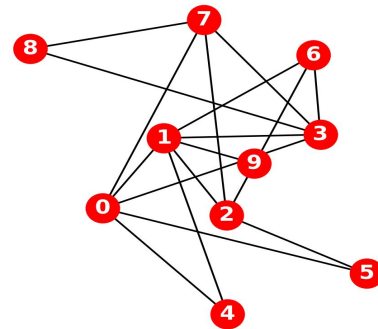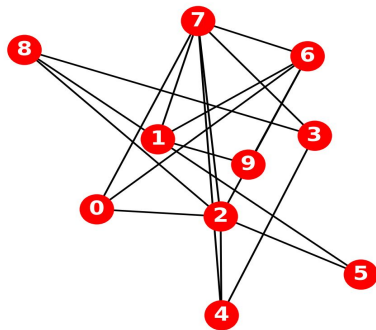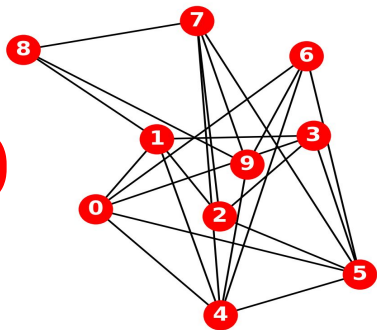Autism Brain Imaging Data Exchange

## Introduction

Autism spectrum disorder (ASD) is characterized by qualitative impairment in social reciprocity, and by repetitive, restricted, and stereotyped behaviors/interests. Previously considered rare, ASD is now recognized to occur in more than 1% of children. Despite continuing research advances, their pace and clinical impact have not kept up with the urgency to identify ways of determining the diagnosis at earlier ages, selecting optimal treatments, and predicting outcomes. For the most part this is due to the complexity and heterogeneity of ASD. To face these challenges, large-scale samples are essential, but single laboratories cannot obtain sufficiently large datasets to reveal the brain mechanisms underlying ASD. In response, the Autism Brain Imaging Data Exchange (ABIDE) initiative has aggregated functional and structural brain imaging data collected from laboratories around the world to accelerate our understanding of the neural bases of autism. With the ultimate goal of facilitating discovery science and comparisons across samples, the ABIDE initiative now includes two large-scale collections: ABIDE I and ABIDE II. Each collection was created through the aggregation of datasets independently collected across more than 24 international brain imaging laboratories and are being made available to investigators throughout the world, consistent with open science principles, such as those at the core of the International Neuroimaging Data-sharing Initiative. For details about these initiatives visit the collection specific pages: **ABIDE I** and **ABIDE II**.

# ABIDE motifs

Patient A (ASD)



Patient B (ASD)

We developed a Shapley-based XAI algorithm for graph classification w/ node identity.

GraphSHAP computes attribution scores (a.k.a. feature importances) for a set of arbitrarily-defined motifs.
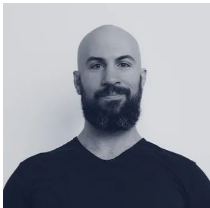
PROs:
- Custom, high-level explanation language
- Scalable algorithm
- Rooted in Shapley's game theory

CONs:
- Requires node identity (so far)
- Requires motifs
- Masking is arbitrary

# Thanks!

Perotti, Bajardi, Bonchi, Panisson, "Explaining Identity-aware Graph Classifiers through the Language of Motifs". International Joint Conference on Neural Networks 2023