

# Masked Mirror Validation in Graphon Estimation

Huimin Cheng  
Boston University

Joint work with Yongkai Chen, Ping Ma, Wenxuan Zhong

# Outline

## **1** Background

- Network, Graphon and Graphon Estimation
- Motivation and Challenges

## **2** Proposed Method: Masked Mirror Validation (MMV)

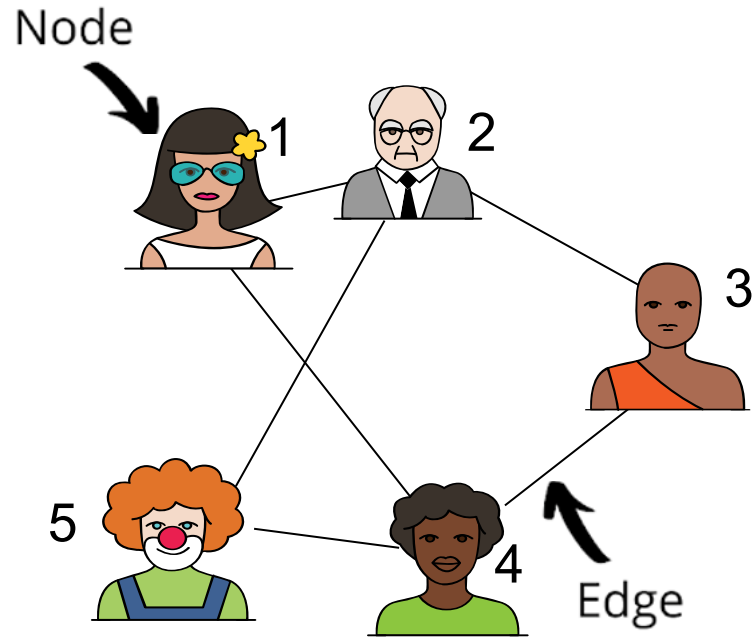
- MMV Procedure
- Theoretical Results
- Simulation Studies

## **3** Application to Drug Repurposing

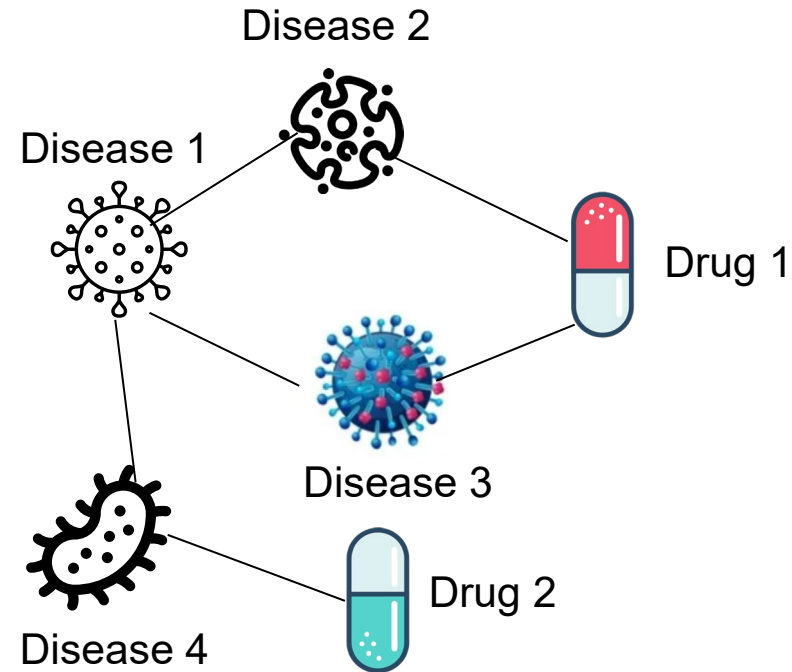
- Drug Repurposing
- Med-Reader AI Tool
- Case Study

# Network (Graph)

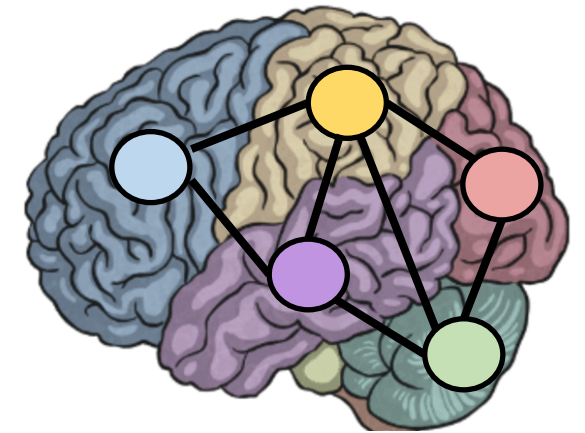
Social network



Drug-disease network

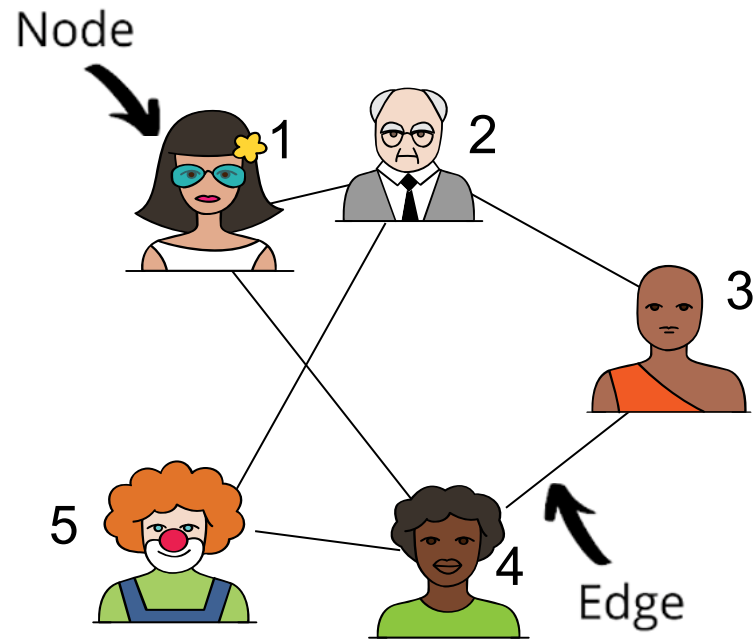


Brain network

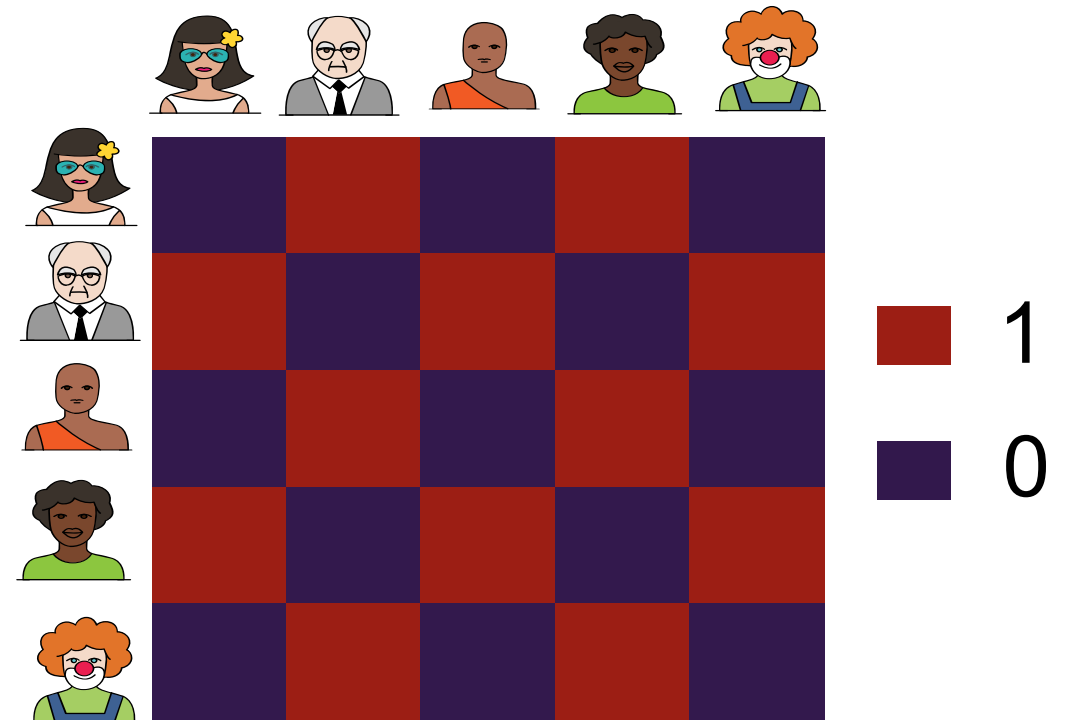


# Mathematical Representation of a Network

Social network



Adjacency matrix  $A = (a_{ij})$

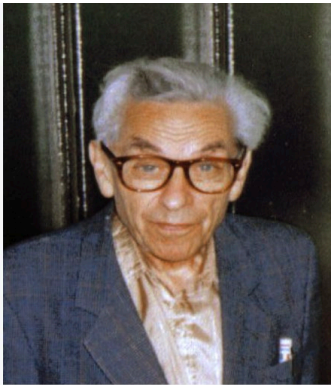


# Generating Model for Network

1960

## Erdős–Rényi model

Paul Erdős and Alfréd Rényi



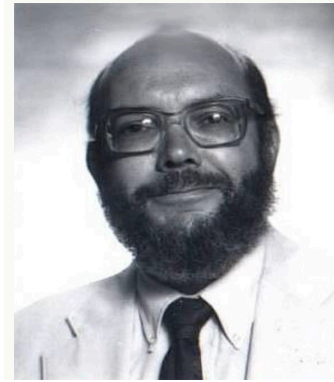
$$a_{ij} \sim \text{Ber}(p)$$

Erdos, and Rényi. *Publ. Math. Inst. Hung. Acad. Sci.* 1960

1983

## Stochastic Block Model (SBM)

Paul W. Holland



$$a_{ij} \sim \text{Ber}(p_{in}), \text{ if } i, j \text{ are in a same group}$$
$$a_{ij} \sim \text{Ber}(p_{out}), \text{ otherwise}$$

Holland, Laskey and Leinhardt. *Social networks* 1983.

2006

## Graphon model

László Lovász

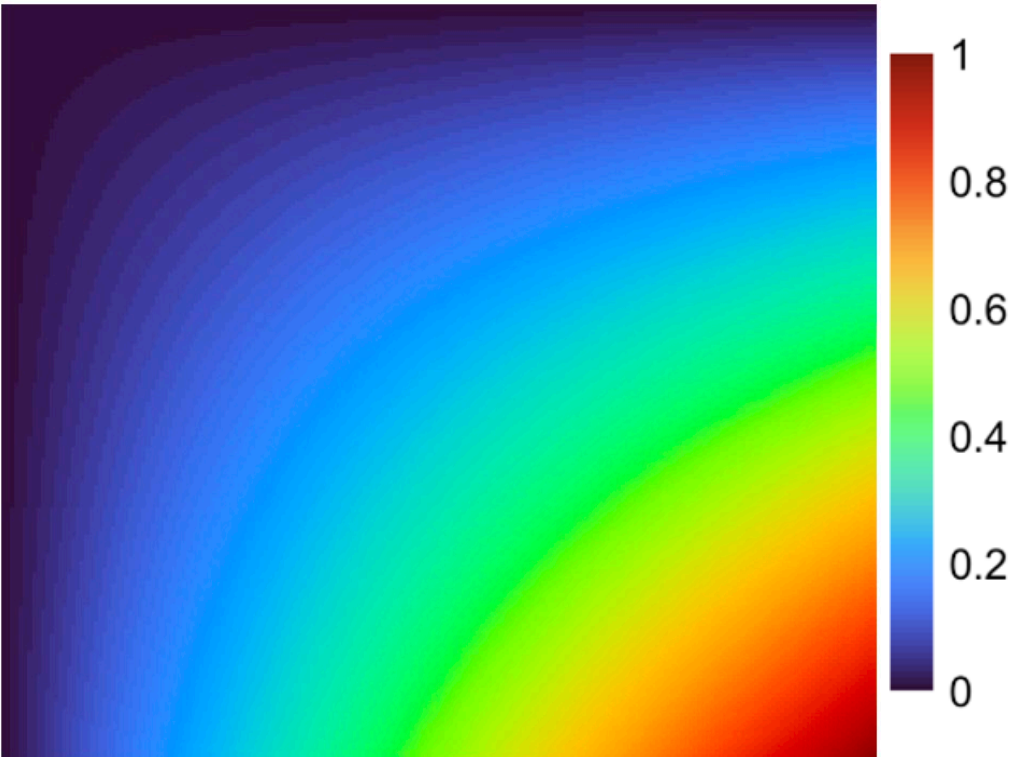


$$a_{ij} \sim \text{Ber}(p_{ij}),$$
$$p_{ij} = f(u_i, u_j)$$
$$u_i \in [0,1], u_j \in [0,1]$$

Lovász, and Szegedy. *Journal of Combinatorial Theory, Series B* 2006

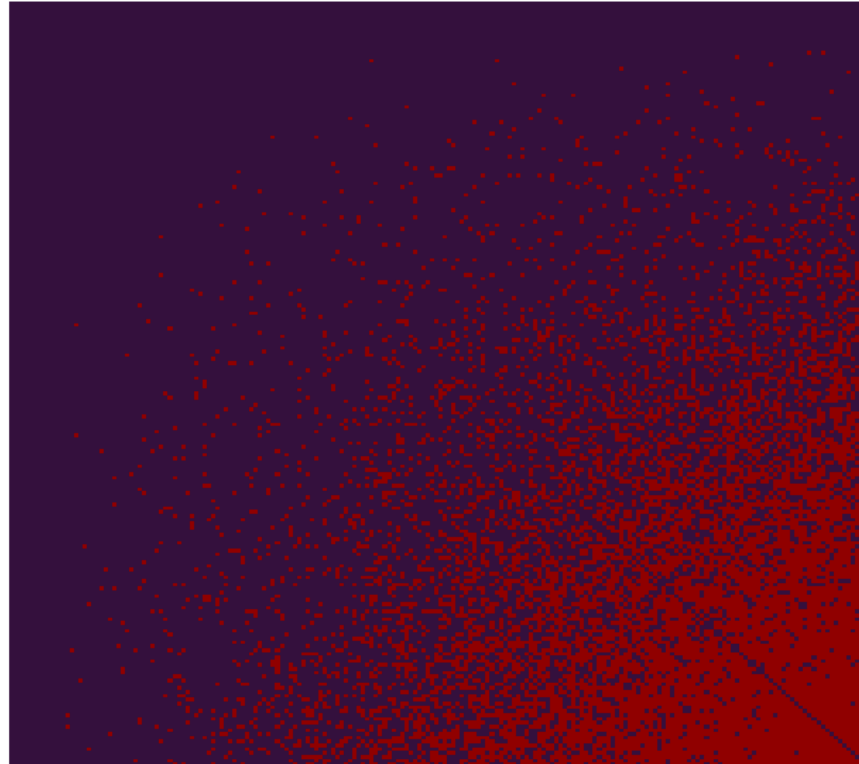
# Heuristic of Probability Matrix

Graphon Probability matrix  $\mathbf{P} = (p_{ij})$



Generate  
→  
←  
Estimate

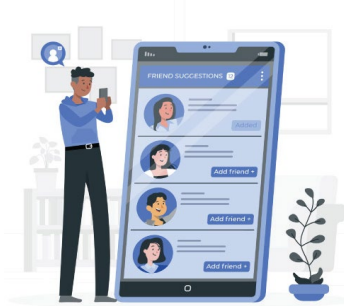
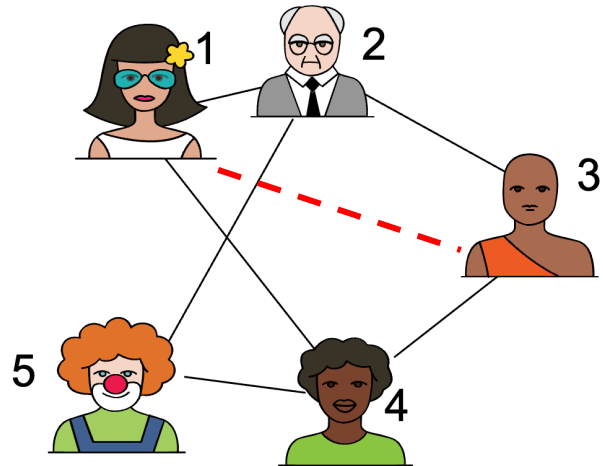
Adjacency matrix



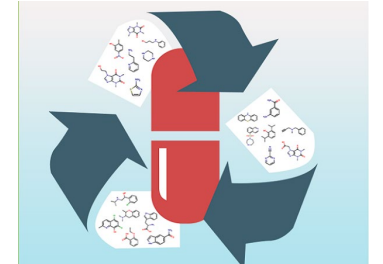
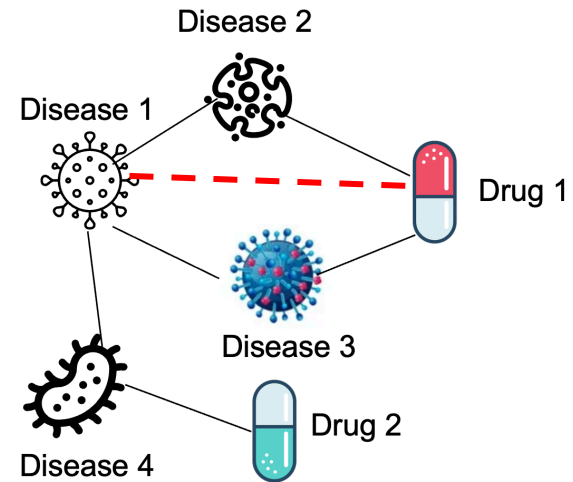
■ 1 ■ 0

# Why Care Graphon Estimation?

## Friend Recommendation



## Drug Repurposing

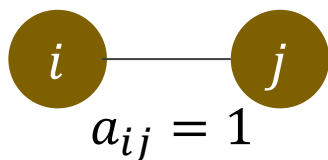


# Graphon Estimation Methods

Sort-and-smoothing method (Chan and Edoardo. ICML 2014) (SAS) and

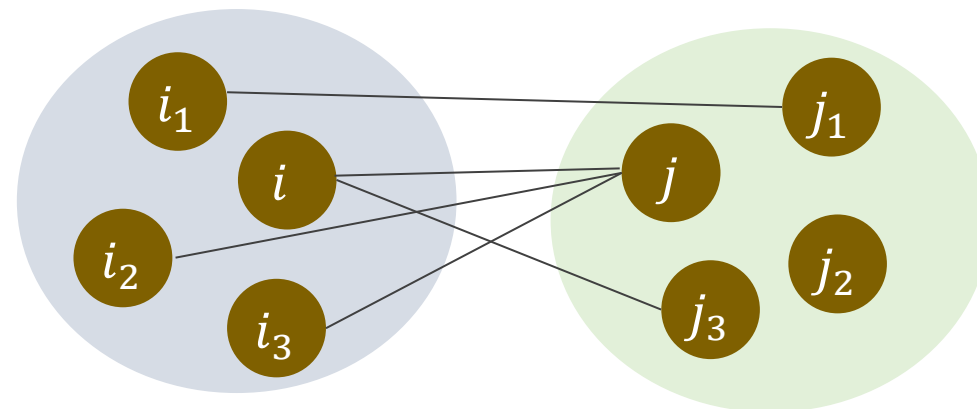
Depends on a **hyperparameter: Neighborhood size  $h$**

## Challenge



Only one observation for  $p_{ij}$ !

“Neighbors” of  $i$       “Neighbors” of  $j$



“Replicated observations” for  $p_{ij}$

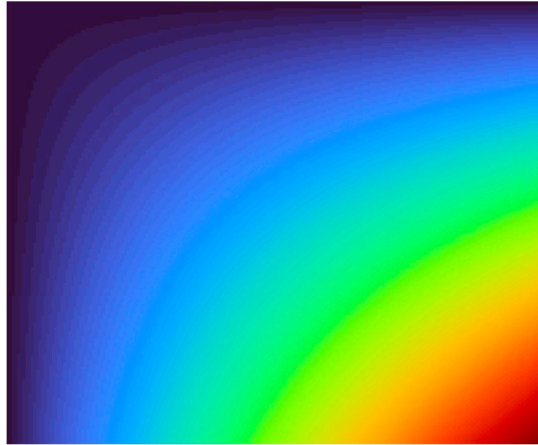
Neighborhood smoothing (NS) method (Zhang, Levina and Zhu. *Biometrika* 2017)

Depends on a **hyperparameter: Neighborhood size  $m\sqrt{n \log n}$**



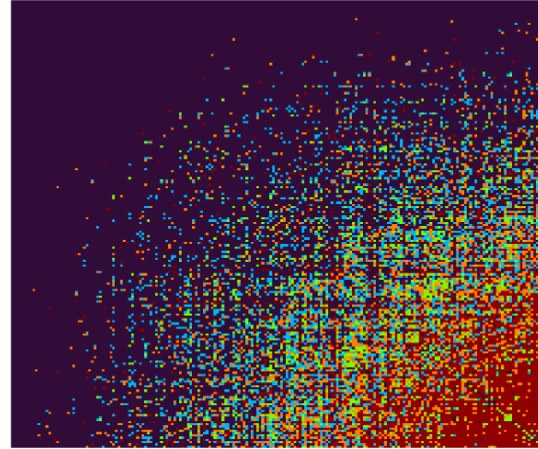
# Hyperparameter Influences Graphon Estimation

True  $\mathbf{P} = (p_{ij})$

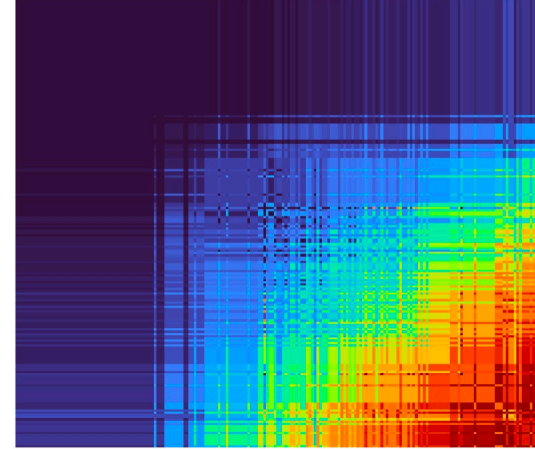


Estimator  $\hat{\mathbf{P}} = (\hat{p}_{ij})$

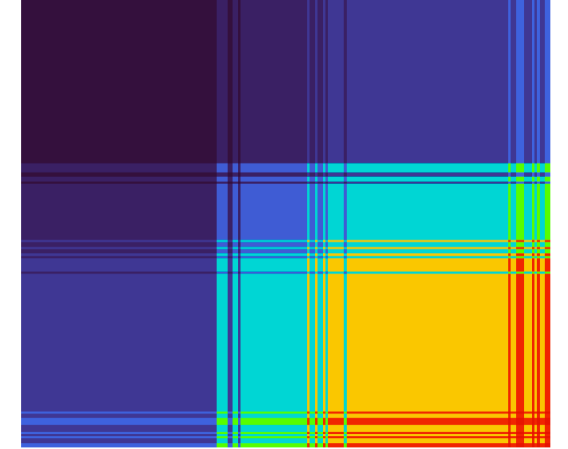
SAS  $h = 2$



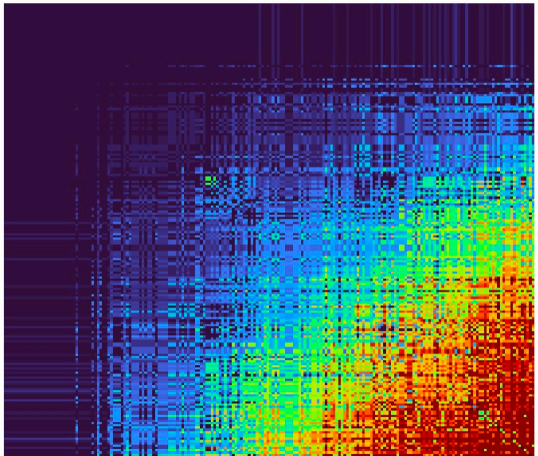
$h = 10$



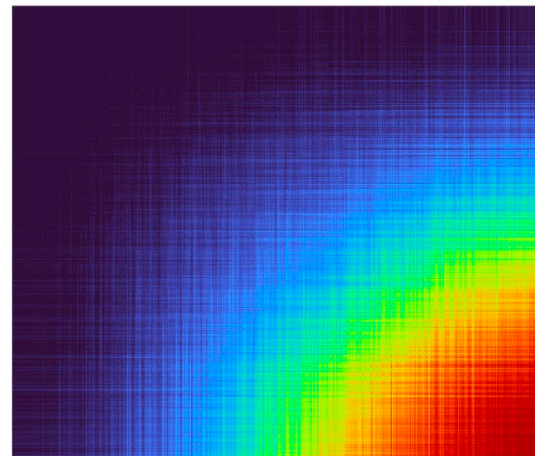
$h = 50$



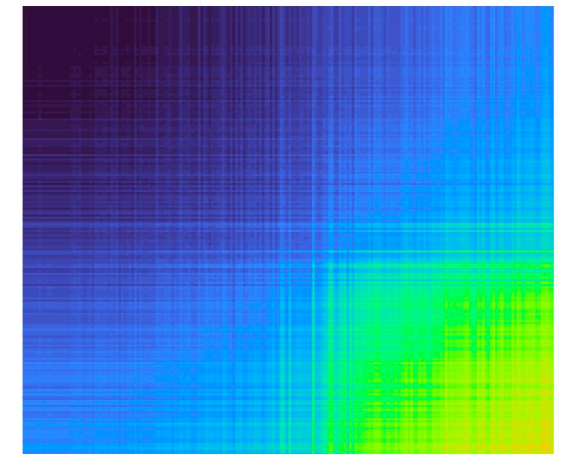
NS  $m = 0.5$



$m = 2$



$m = 5$



# Hyperparameter Tuning

$$\text{MSE}(m) = \frac{1}{n^2} \|\mathbf{P} - \hat{\mathbf{P}}_m\|_F^2$$

$$m^* = \operatorname{argmin}_{m \in M} \text{MSE}(m)$$

## Approaches

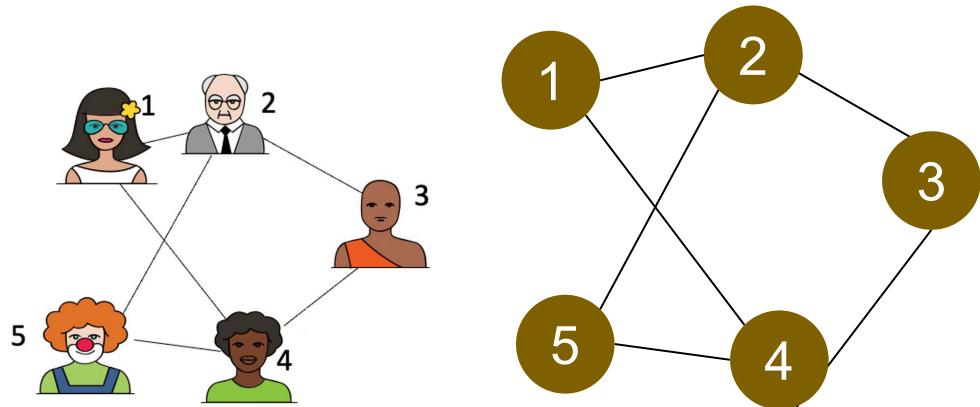
- Empirical experience →
- Analytic solution →
- .....
- Cross-validation →

## Pros and Cons

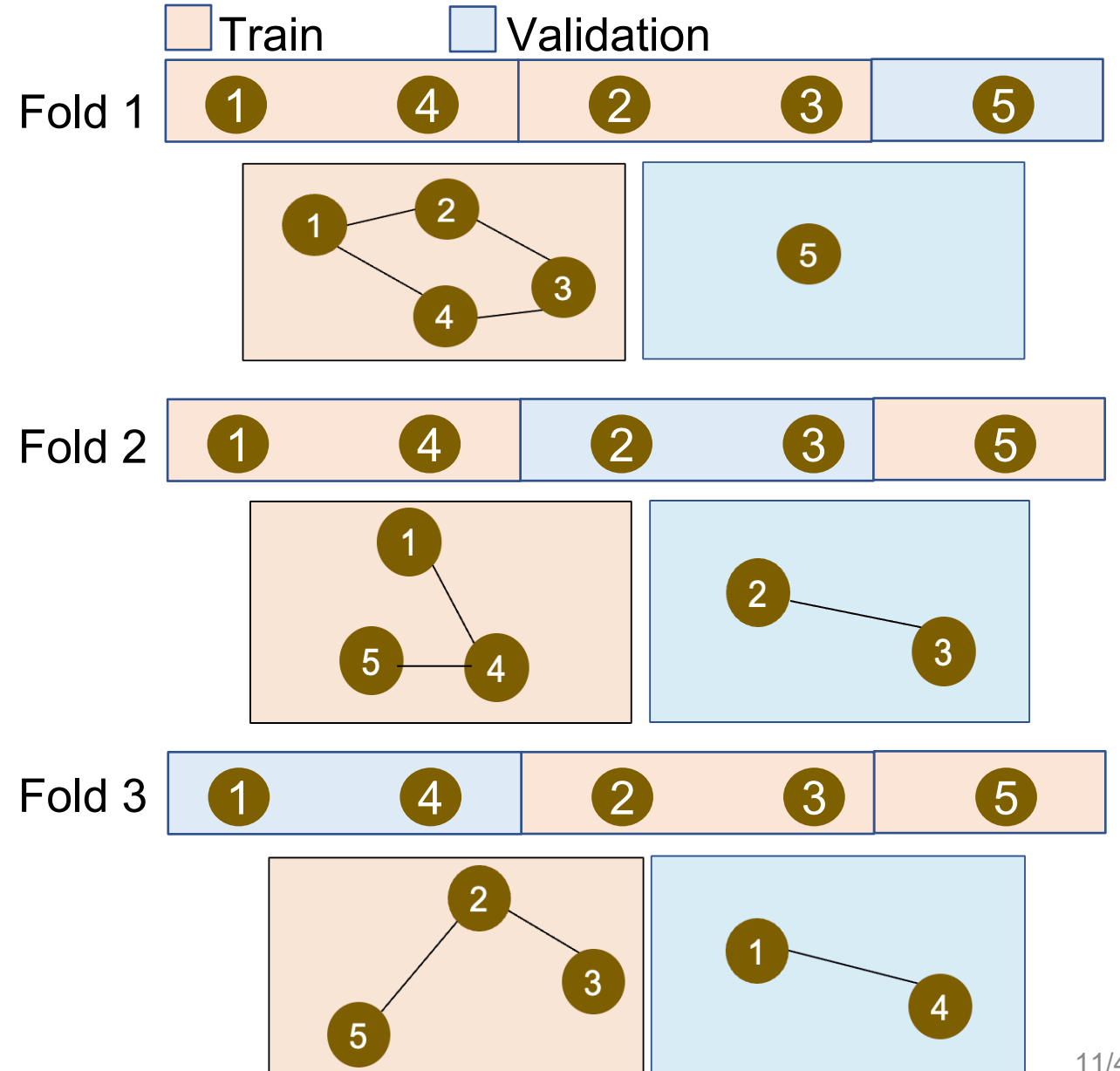
- Lack of theoretical guarantee
- Explicit form between  $\hat{\mathbf{P}}_m$  and  $m$  is usually elusive
- Theoretical guarantee; Effective**

# Problems of Node Splitting

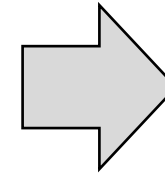
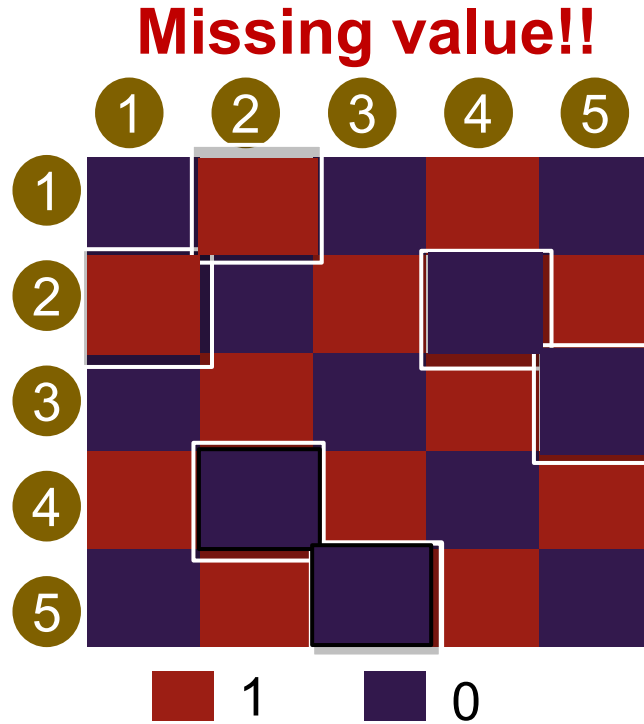
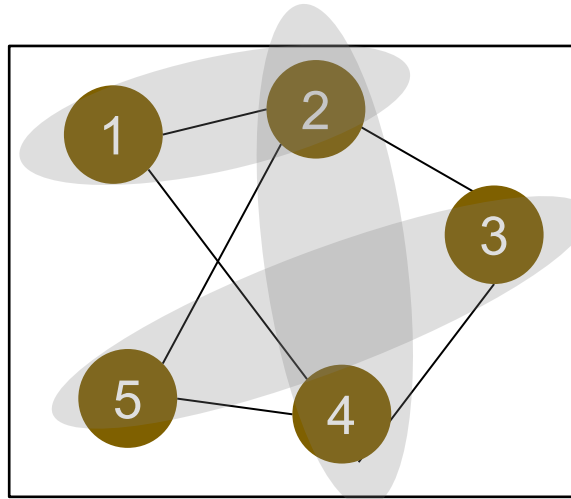
Node splitting will **destroy**  
**the network structure!**



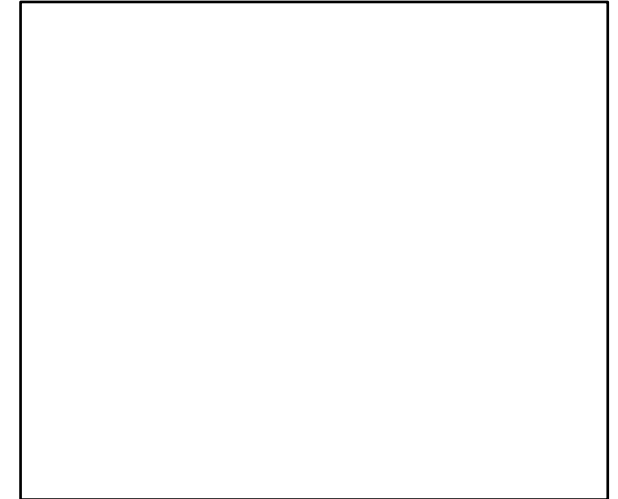
Three folds



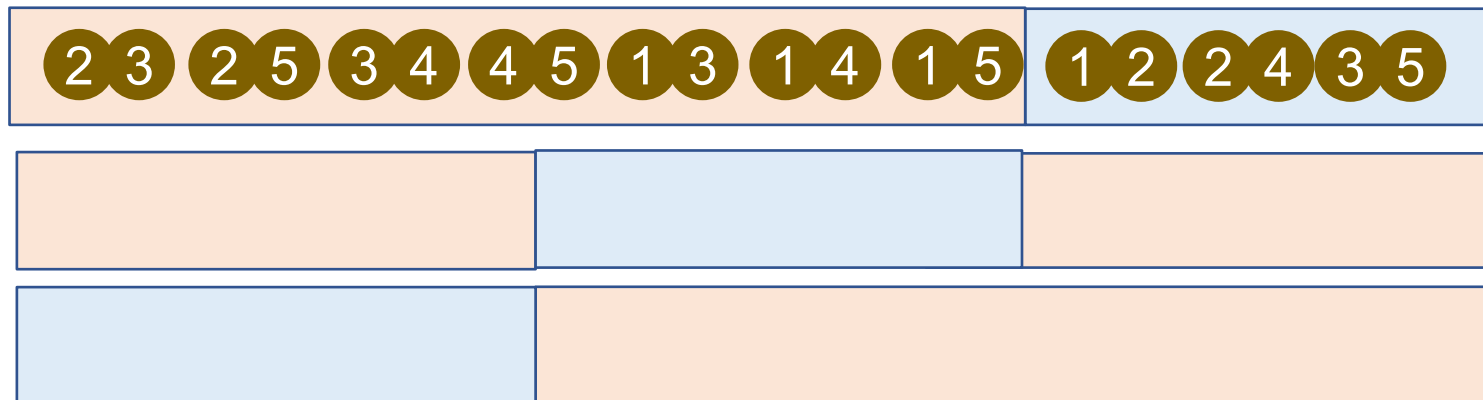
# Challenges of Edge Splitting



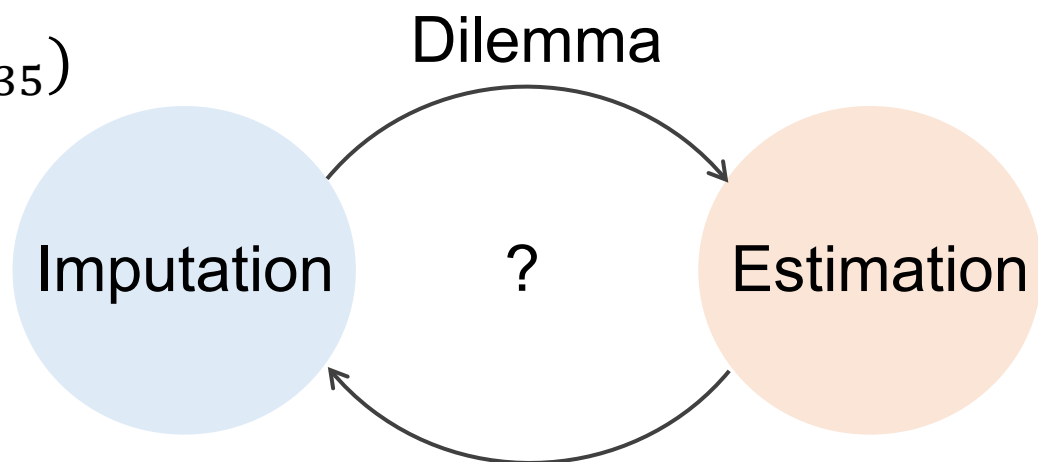
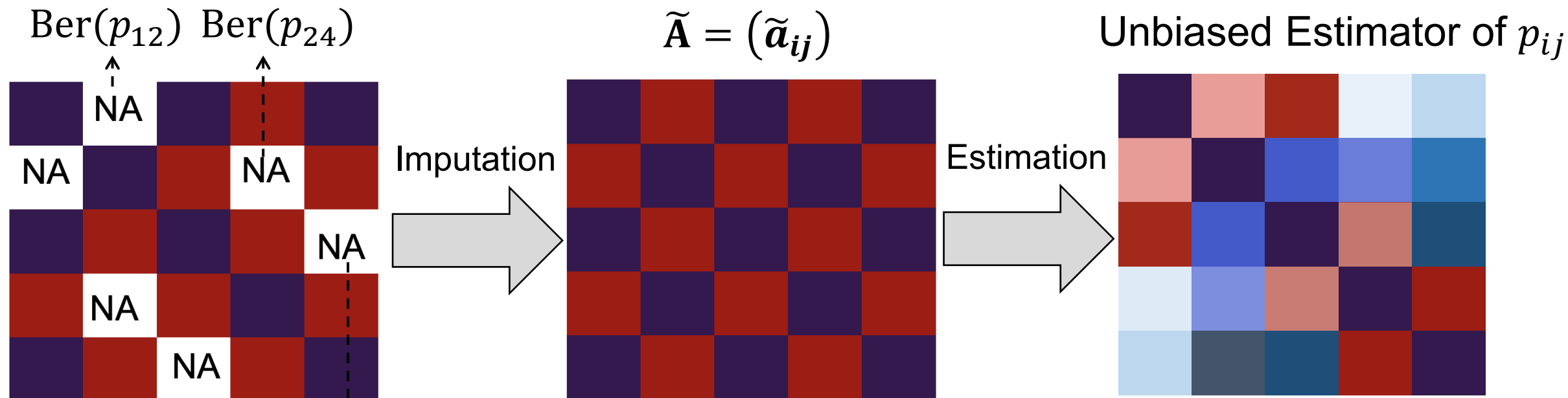
Validation data



Three-fold split    ■ Train    ■ Validation



# Challenges of Missing Value Imputation

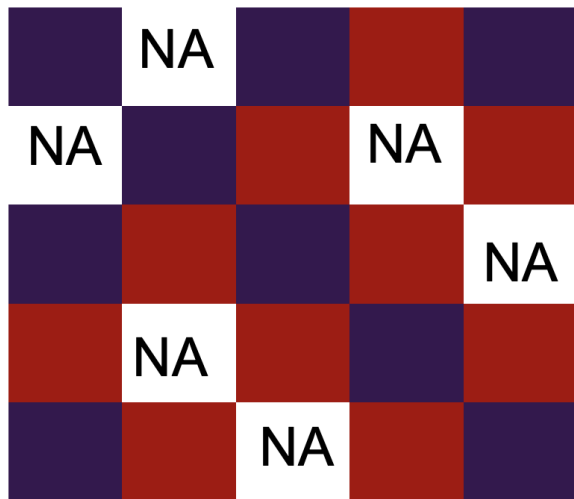


# Literature Review of Network Cross-Validation

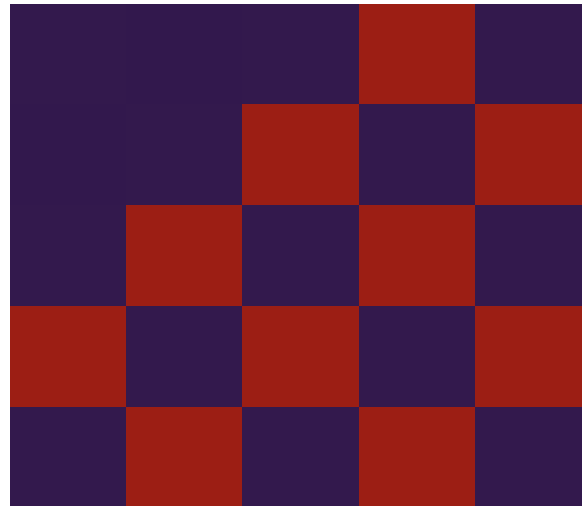
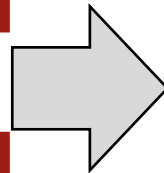
## Edge Cross-Validation (ECV)

( Li, Levina and Zhu. Biometrika 2020 )

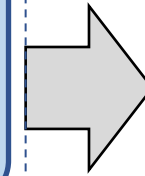
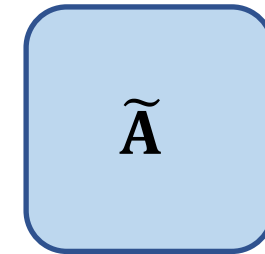
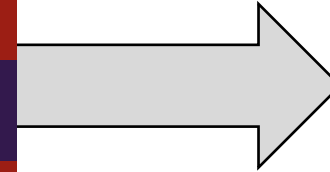
### Imputation



Fill  
zeros



SVD  
Thresholding



### Estimation

Estimator  
of  $p_{ij}$

- Restrictive assumption:  $\text{Rank}(P) \leq \frac{n}{\delta}$ , where  $\delta$  is average degree
- Introduce additional hyperparameter: Threshold
- Expensive computational cost:  $O(n^3)$

# Outline

## 1 Background

- Network, Graphon and Graphon Estimation
- Motivation and Challenges of Network Cross-Validation

## 2 Proposed Method: Masked Mirror Validation (MMV)

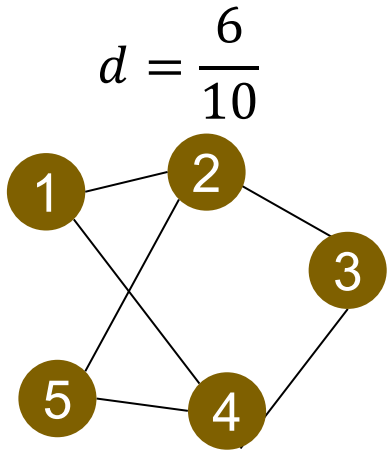
- MMV Procedure
- Theoretical Results
- Simulation Studies

## 3 Application to Drug Repurposing

- Drug Repurposing
- Med-Reader AI Tool
- Case Study

# MMV Procedure

## Step 1: Imputation



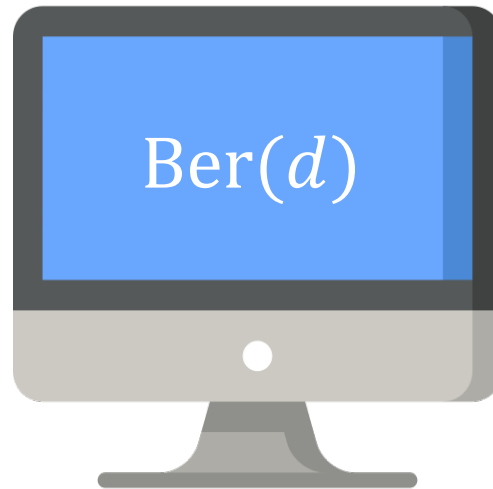
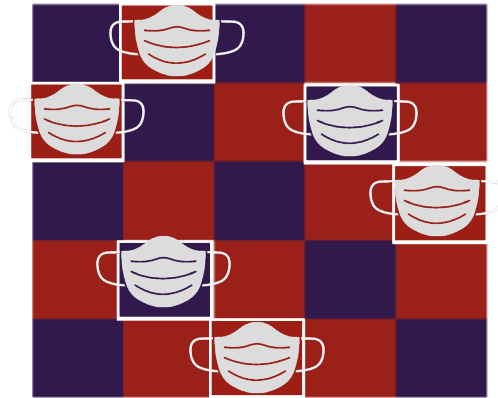
$$\text{Pooled mean } p_{..} = \frac{\sum_{i < j} p_{ij}}{\binom{n}{2}}$$

Network density

$$d = \frac{\sum_{i < j} a_{ij}}{\binom{n}{2}}$$

is an unbiased estimator for  $p_{..}$

Training data  $\tilde{\mathbf{A}} = (\tilde{a}_{ij})$

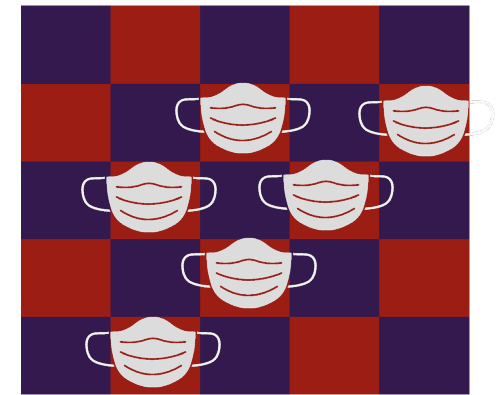
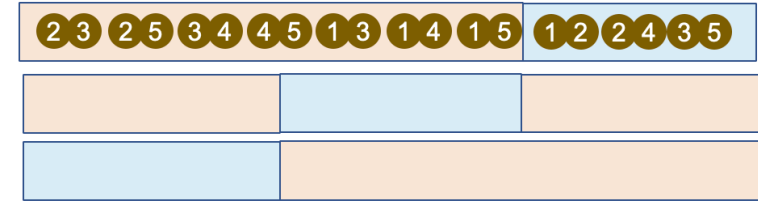


Fold 1

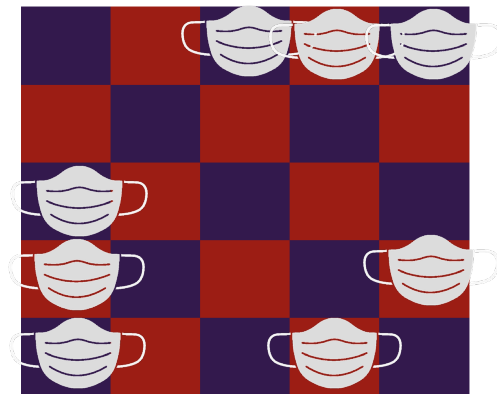
$$\tilde{a}_{12} = 1$$

$$\tilde{a}_{24} = 0$$

$$\tilde{a}_{35} = 1$$



Fold 2



Fold 3



# Distribution of $\tilde{\mathbf{A}}$

$$\tilde{\mathbf{A}} = (\tilde{a}_{ij}), \tilde{a}_{ij} \sim \text{Ber}(\tilde{p}_{ij})$$

$K$  is the number of folds

If  $i$  and  $j$  are in validation

$$\tilde{a}_{ij} \sim \text{Ber}(d)$$

$$\mathbb{P}(i \text{ and } j \text{ are in validation}) = \frac{1}{K}$$

$$\mathbb{P}(\tilde{a}_{ij} = 1 | i \text{ and } j \text{ are in validation}) = d$$

If  $i$  and  $j$  are not in validation

$$\tilde{a}_{ij} = a_{ij} \sim \text{Ber}(p_{ij})$$

$$\mathbb{P}(i \text{ and } j \text{ are not in validation}) = \frac{K-1}{K}$$

$$\mathbb{P}(\tilde{a}_{ij} = 1 | i \text{ and } j \text{ are not in validation}) = p_{ij}$$

By law of total probability

$$\tilde{p}_{ij} := \mathbb{P}(\tilde{a}_{ij} = 1) = \frac{1}{K}d + \frac{K-1}{K}p_{ij}$$

# Derivation of Bias Correction Formula

$$\tilde{p}_{ij} := \mathbb{P}(\tilde{a}_{ij} = 1) = \frac{1}{K}d + \frac{K-1}{K}p_{ij}$$

Reverse map

$$p_{ij} = \frac{K\tilde{p}_{ij} - d}{K-1}$$

Used for bias correction

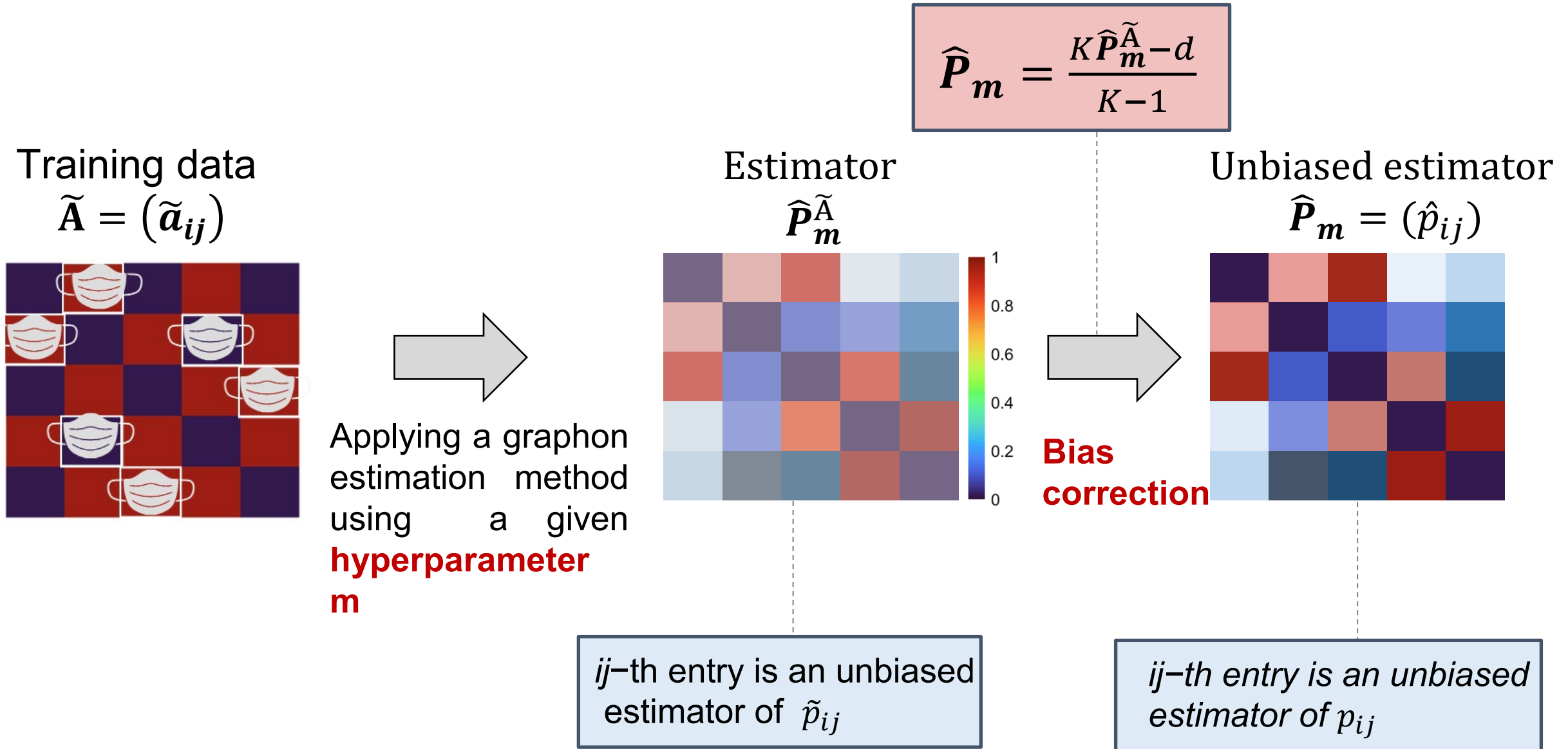
$$\hat{p}_{ij} = \frac{K\hat{\tilde{p}}_{ij} - d}{K-1}$$

Unbiased estimator of  $p_{ij}$

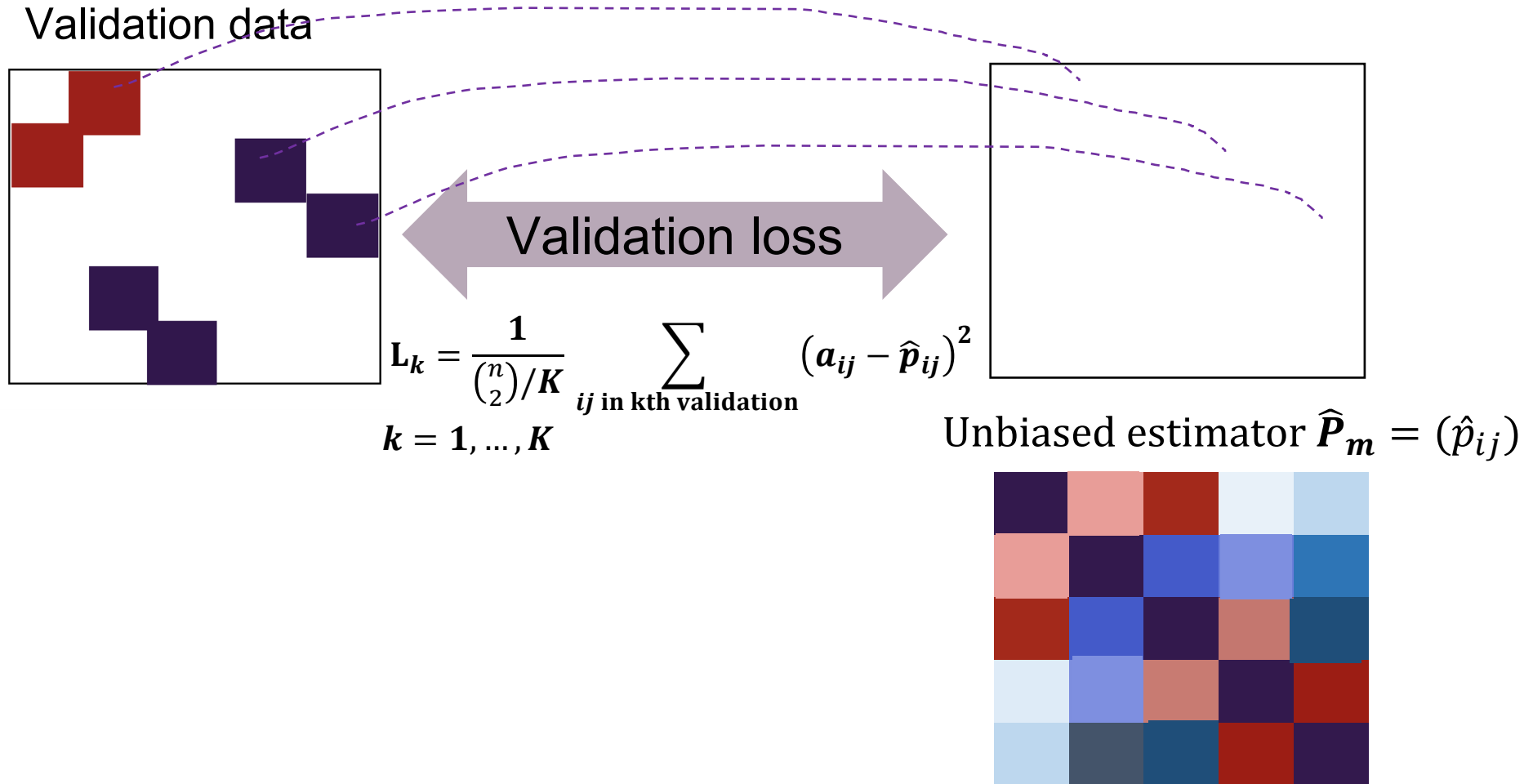
Unbiased estimator of  $\tilde{p}_{ij}$

By applying graphon estimation method to  $\tilde{\mathbf{A}}$ , we can get an unbiased estimator of  $\tilde{p}_{ij}$ .

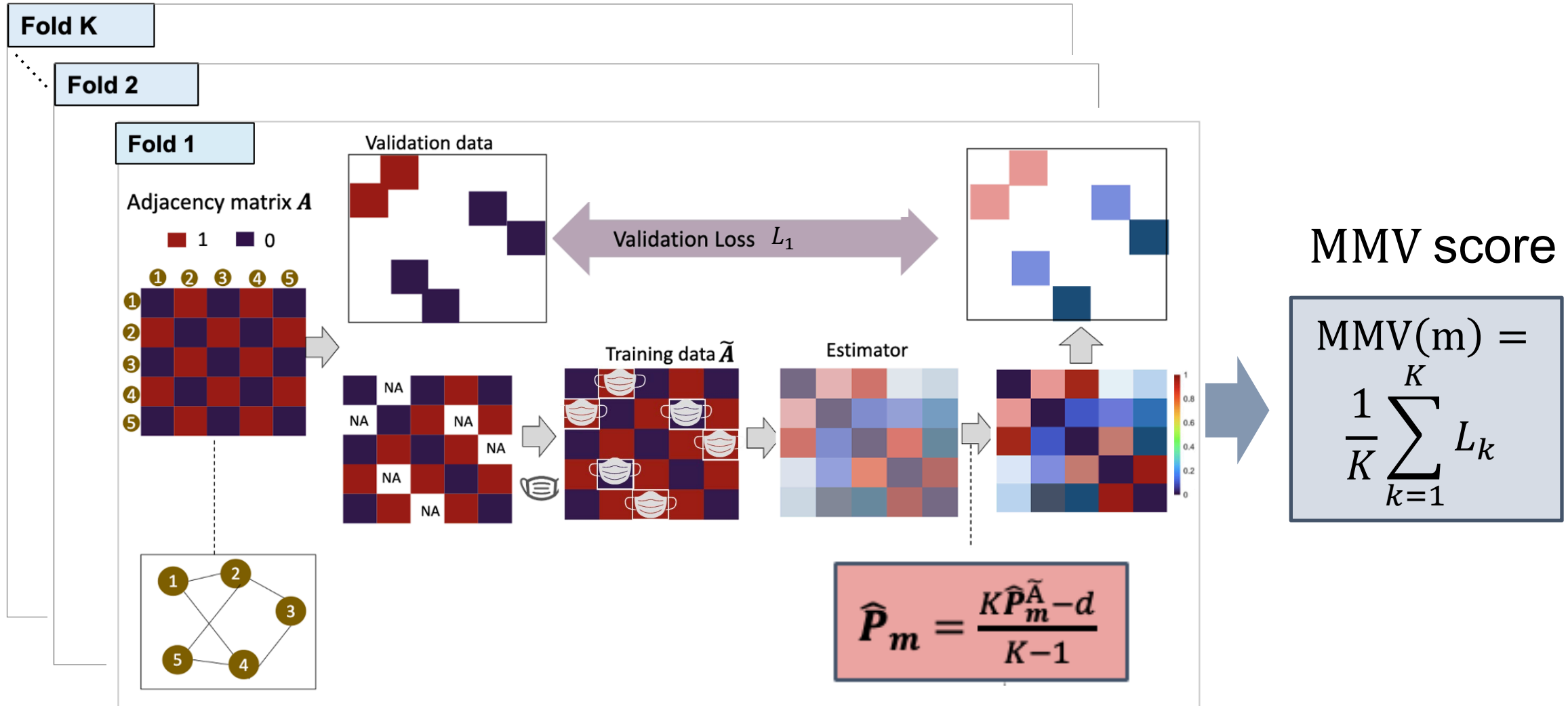
# Step 2: Estimation and Bias Correction



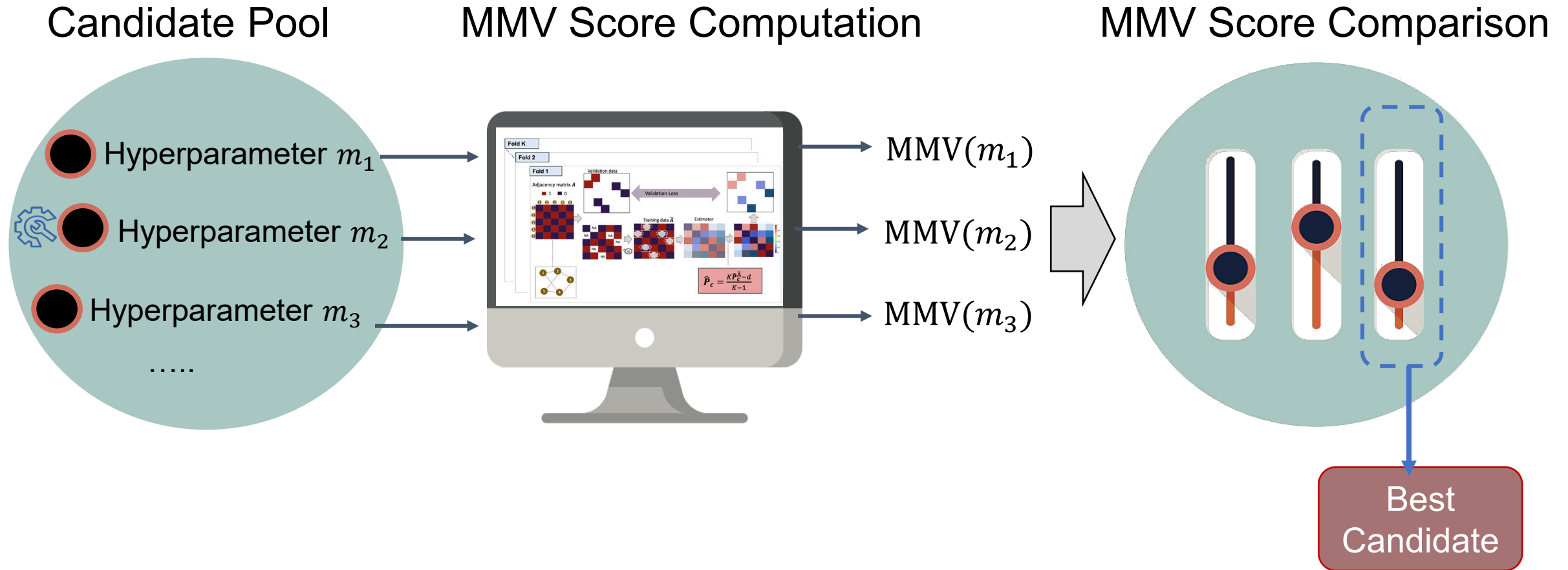
# Step 3: Validation Loss



# MMV Score of a Hyperparameter



# Hyperparameter Tuning Using MMV score



# Theoretical Results

Assumption 1:  $K = O(n)$ .

Assumption 2:  $\frac{\|\hat{\mathbf{P}}_m^{\tilde{\mathbf{A}}} - \hat{\mathbf{P}}_m^{\mathbf{A}}\|_F}{\|\tilde{\mathbf{A}} - \mathbf{A}\|_F} = o_p(1)$ .

## Theorem (Selection consistency of MMV)

Given a set of hyperparameters, under Assumptions 1 and 2, as  $n \rightarrow \infty$ , the probability of MMV selecting the optimal hyperparameter converges to one.

↓  
Minimize the MSE

# Simulation Setting

Generate  $\mu_i$

Generate  $p_{ij}$

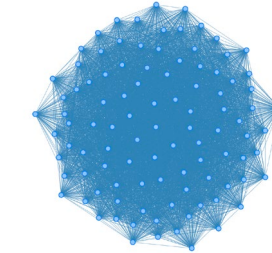
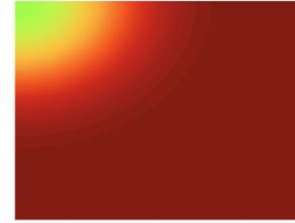
Probability matrix  
 $P = (p_{ij}) \in \mathbb{R}_{n \times n}$

Generate network  
using  $a_{ij} \sim \text{Ber}(p_{ij})$

$$\mu_i \sim U(0,1), \\ i = 1, \dots, n$$

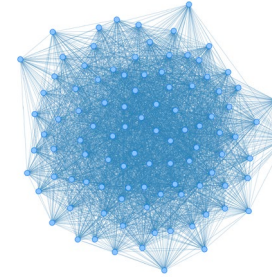
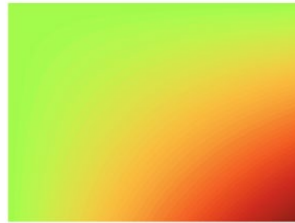
Graphon Setting 1

$$p_{ij} = \frac{1}{1 + \exp(-10(\mu_i^2 + \mu_j^2))}$$



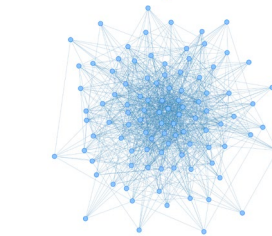
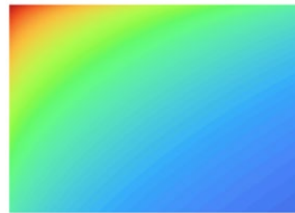
Graphon Setting 2

$$p_{ij} = 0.5 + \frac{\mu_i \mu_j}{3}$$



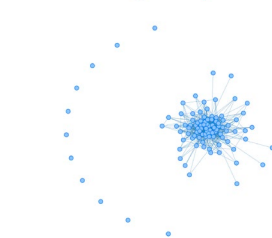
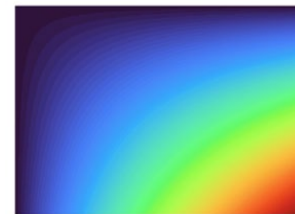
Graphon Setting 3

$$p_{ij} = \mu_i \mu_j$$



Graphon Setting 4

$$p_{ij} = \exp(-(\mu_i^{0.7} + \mu_j^{0.7}))$$



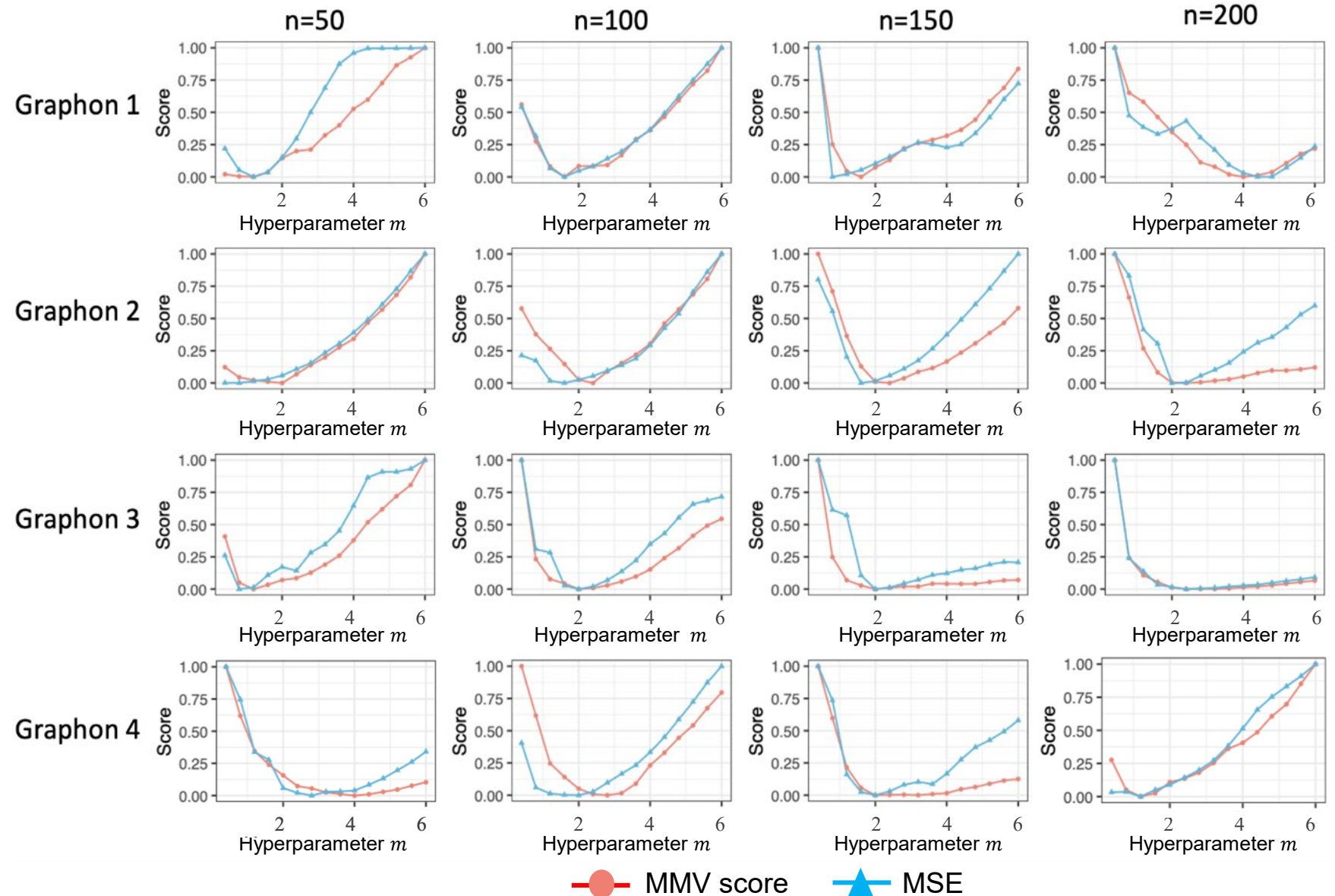
$n = 50, 100,$   
 $150, 200$

$K = 0.1n$

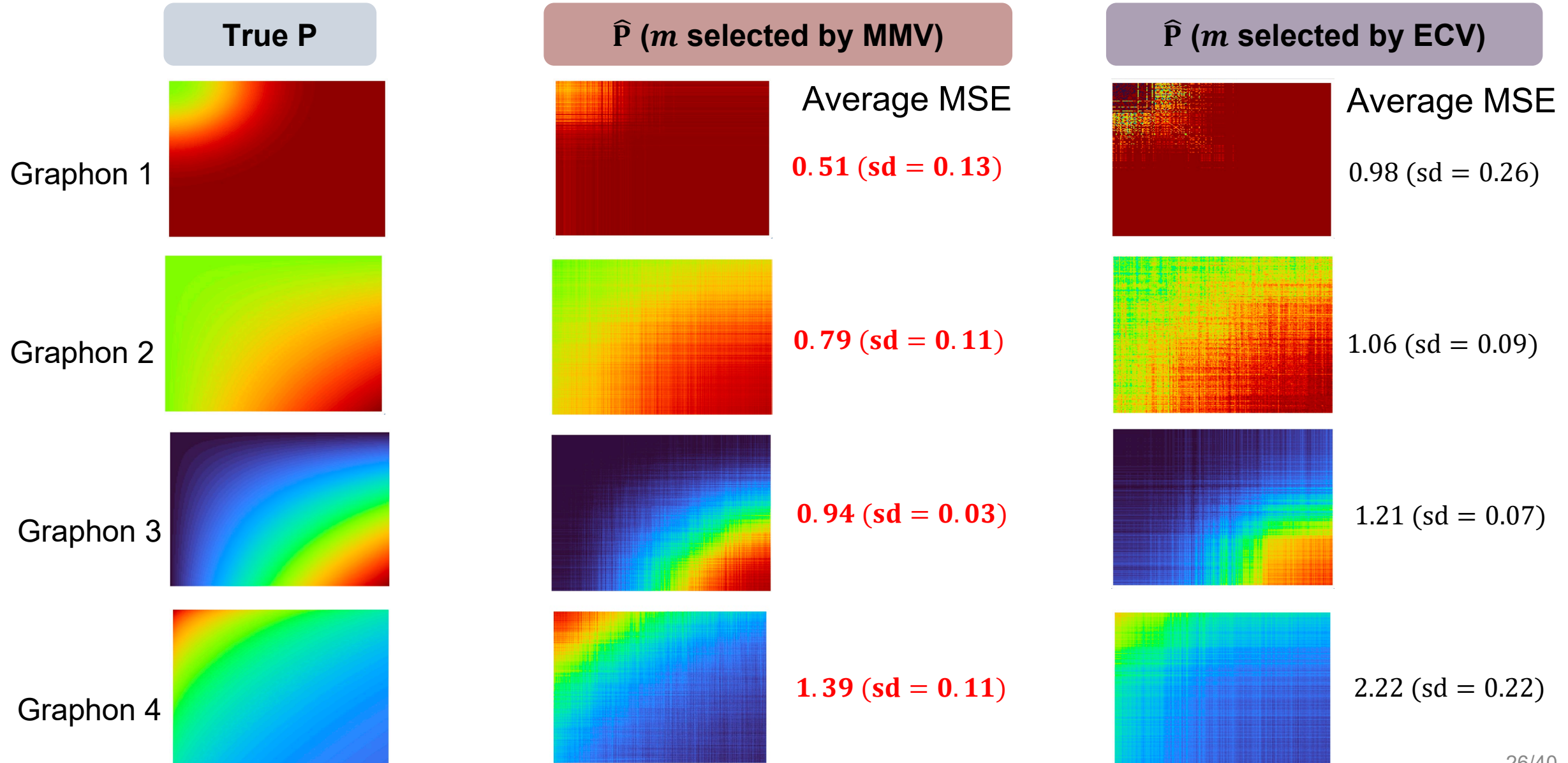


# Simulation Results: Validation of Our Theorem

Tuning hyperparameter  $m$  in NS method  
(Zhang, Levina and Zhu. Biometrika 2017)

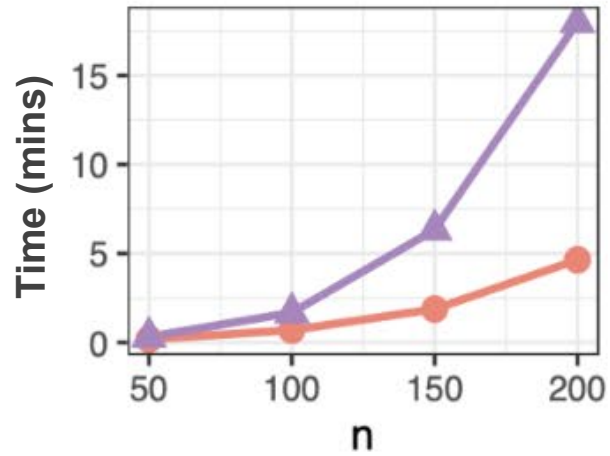


# Simulation Results: Graphon Estimation

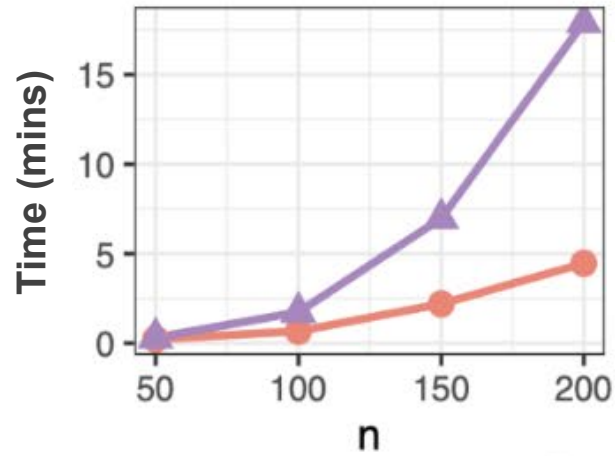


# Simulation Results: Computation Time

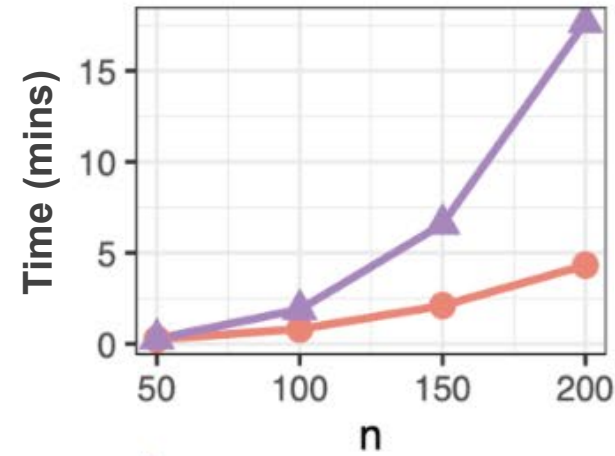
Graphon 1



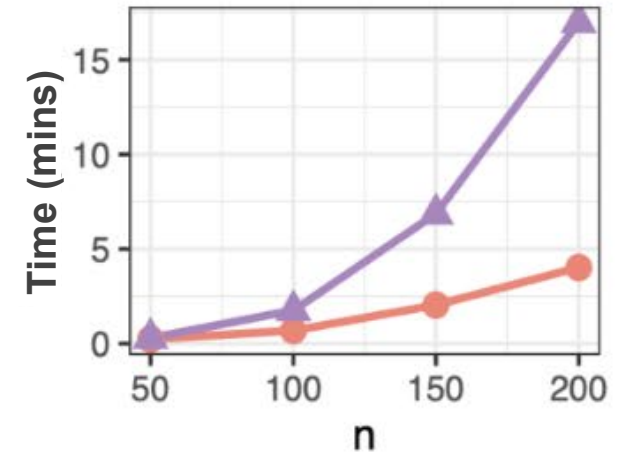
Graphon 2



Graphon 3



Graphon 4



● MMV ▲ ECV

# Outline

## 1 Background

- Network, Graphon and Graphon Estimation
- Motivation and Challenges of Network Cross-Validation

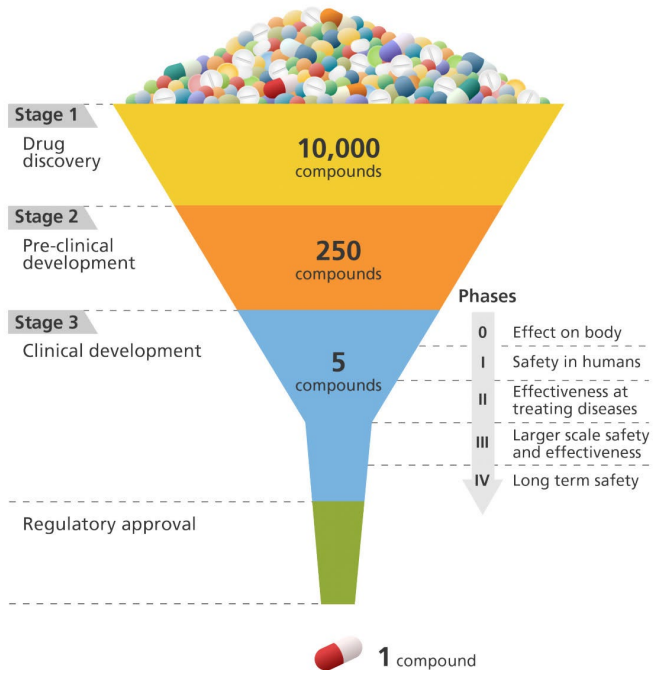
## 2 Proposed Method: Masked Mirror Validation (MMV)

- MMV Procedure
- Theoretical Results
- Simulation Studies

## 3 Application to Drug Repurposing

- Drug Repurposing
- Med-Reader AI Tool
- Case Study

# Drug Repurposing



12-15 years

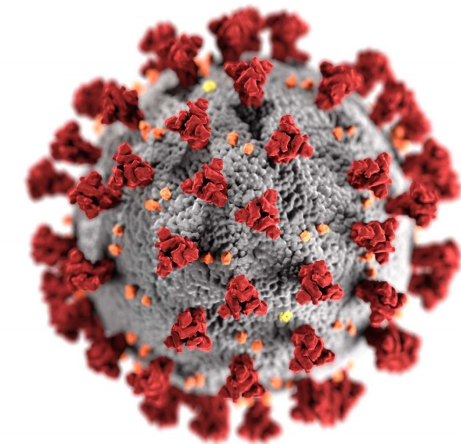
\$2-3 billion

## Drug repurposing example

### Influenza

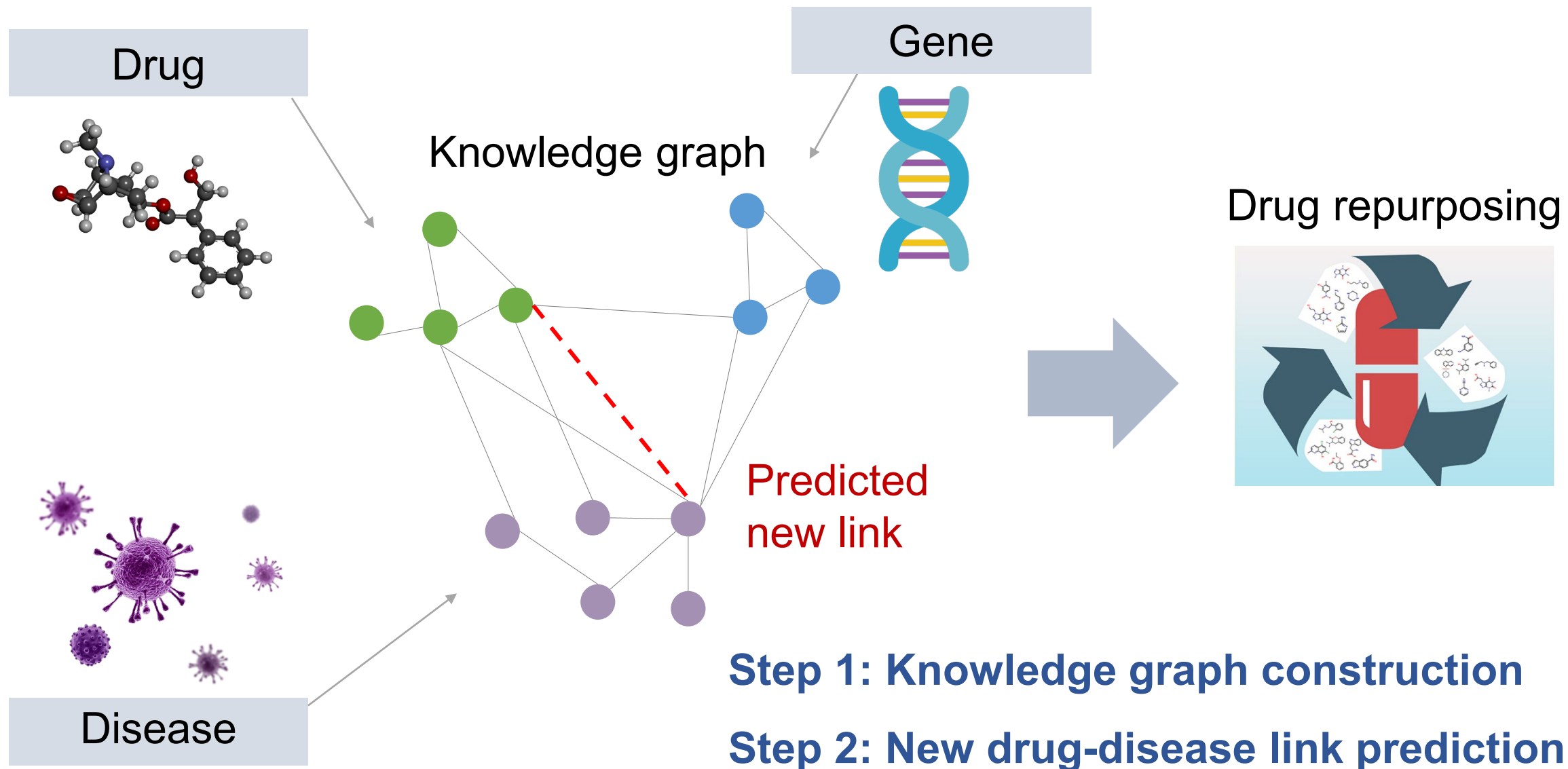


### COVID-19





# Drug Repurposing Using Knowledge Graph



# Challenge of Step 1

PubMed.gov COVID-19 Search

Advanced Create alert Create RSS User Guide

Save Email Send to Sorted by: Best match Display options

MY NCBI FILTERS 326,069 results Page 1 of 1,631

If the reading speed is 30 mins one paper

It takes **19 years** to read all papers



# Solution: Our Developed Med-Reader

Med-Reader: Help expedite your research

Home

Liquid Hot Topics

Network with Significance Score

Multi-databases Cross Validation

Hypotheses Generation

Tutorial

About this site



## Medical AI Reader

Publication Start Date

2020-01-01

Publication End Date

2020-04-30

COVID-19

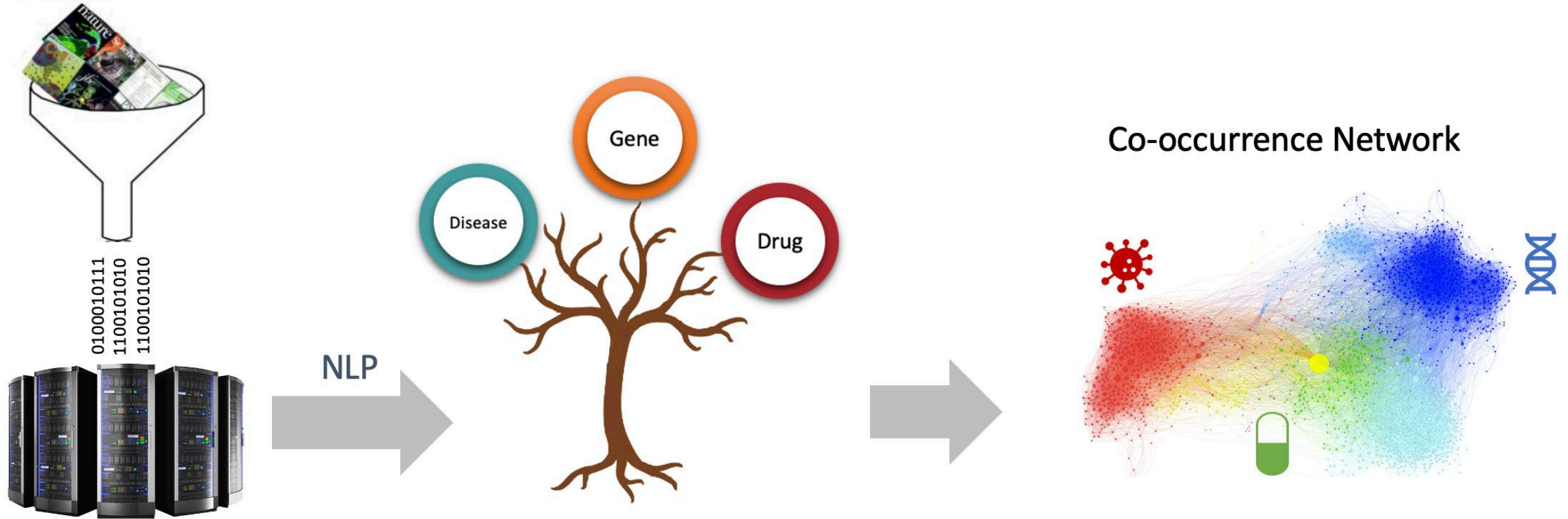
Submit Query

Click 'Submit query' directly to get example results for input query 'COVID-19'.  
Note: This search query will be passed to PubMed to retrieve relevant publications.  
[Click here for PubMed query search help](#)

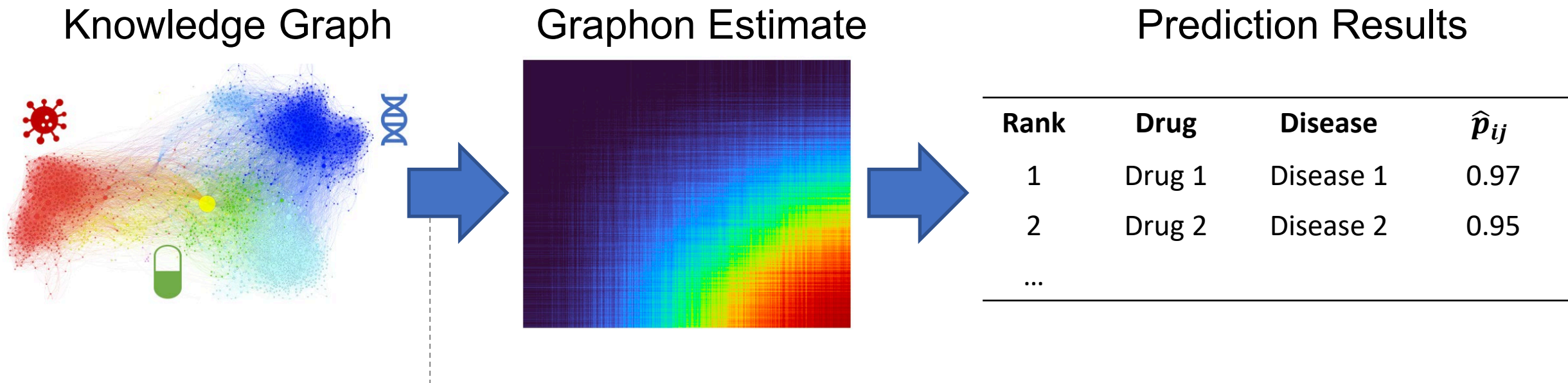




# Med-Reader for Knowledge Graph Construction



# Graphon Estimation in Step 2



Challenge: How to tune hyperparameters?

**Solution: Masked Mirror Validation!**



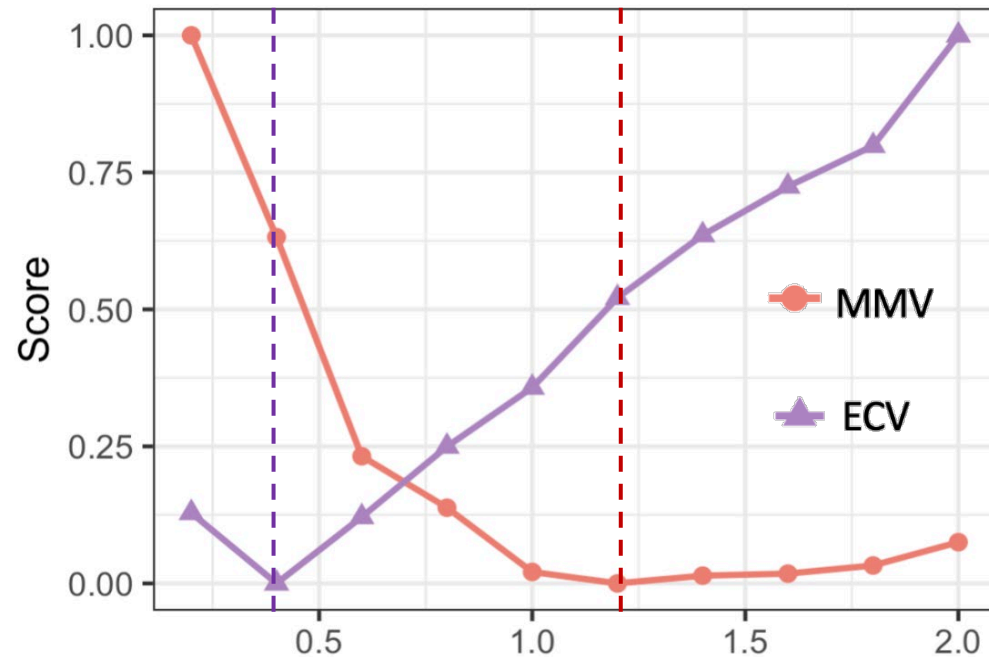
# Case Study of Drug Repurposing

## Step 2: Graphon Estimation

### Hyperparameter Tuning

Tuning hyperparameter  $m$  in NS method  
(Zhang, Levina and Zhu. Biometrika 2017)

MMV selects  $m = 1.2$



ECV selects  $m = 0.2$  Hyperparameter  $m$

### Prediction

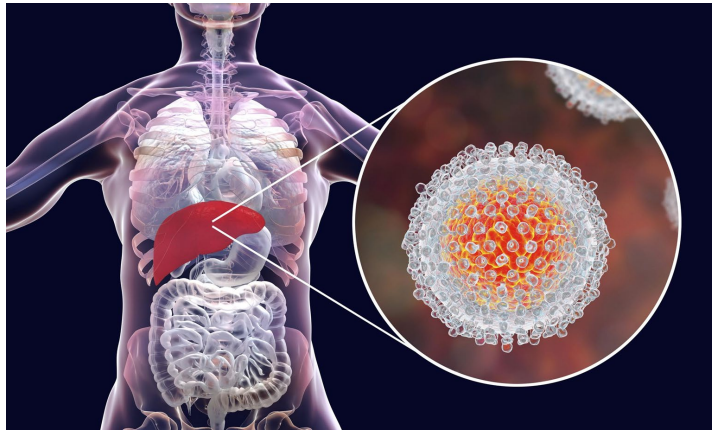
#### Prediction Results using MMV

Rank	Drug	Disease	$\hat{p}_{ij}$
1	Ledipasvir	COVID-19	0.91
2	Budesonide	COVID-19	0.87
...			



# Top Prediction: COVID-19 and Ledipasvir

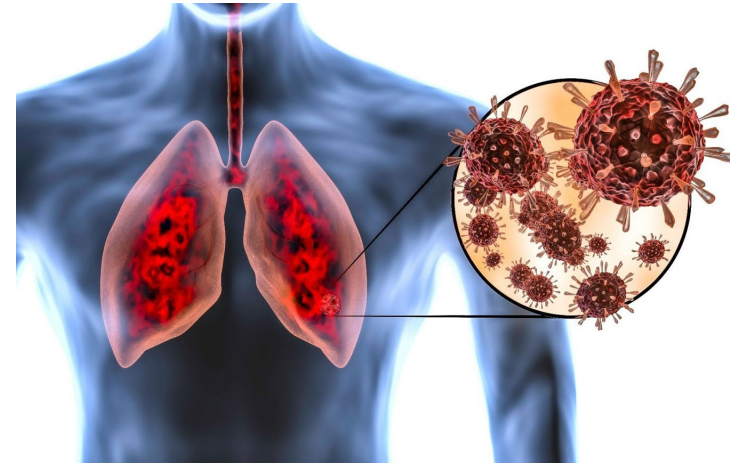
Hepatitis C



Ledipasvir



COVID-19



Scientific way of validating the prediction

1/1/2020 –  
4/30/2020

5/1/2020 – 5/1/2022

Used for **making**  
prediction

Used for **validating** prediction



[Cells](#), 2021 May; 10(5): 1052.

Published online 2021 Apr 29. doi: [10.3390/cells10051052](https://doi.org/10.3390/cells10051052)

PMCID: PMC8146643

PMID: [33946869](https://pubmed.ncbi.nlm.nih.gov/33946869/)

Remdesivir and Ledipasvir among the FDA-Approved Antiviral Drugs Have Potential to Inhibit SARS-CoV-2 Replication

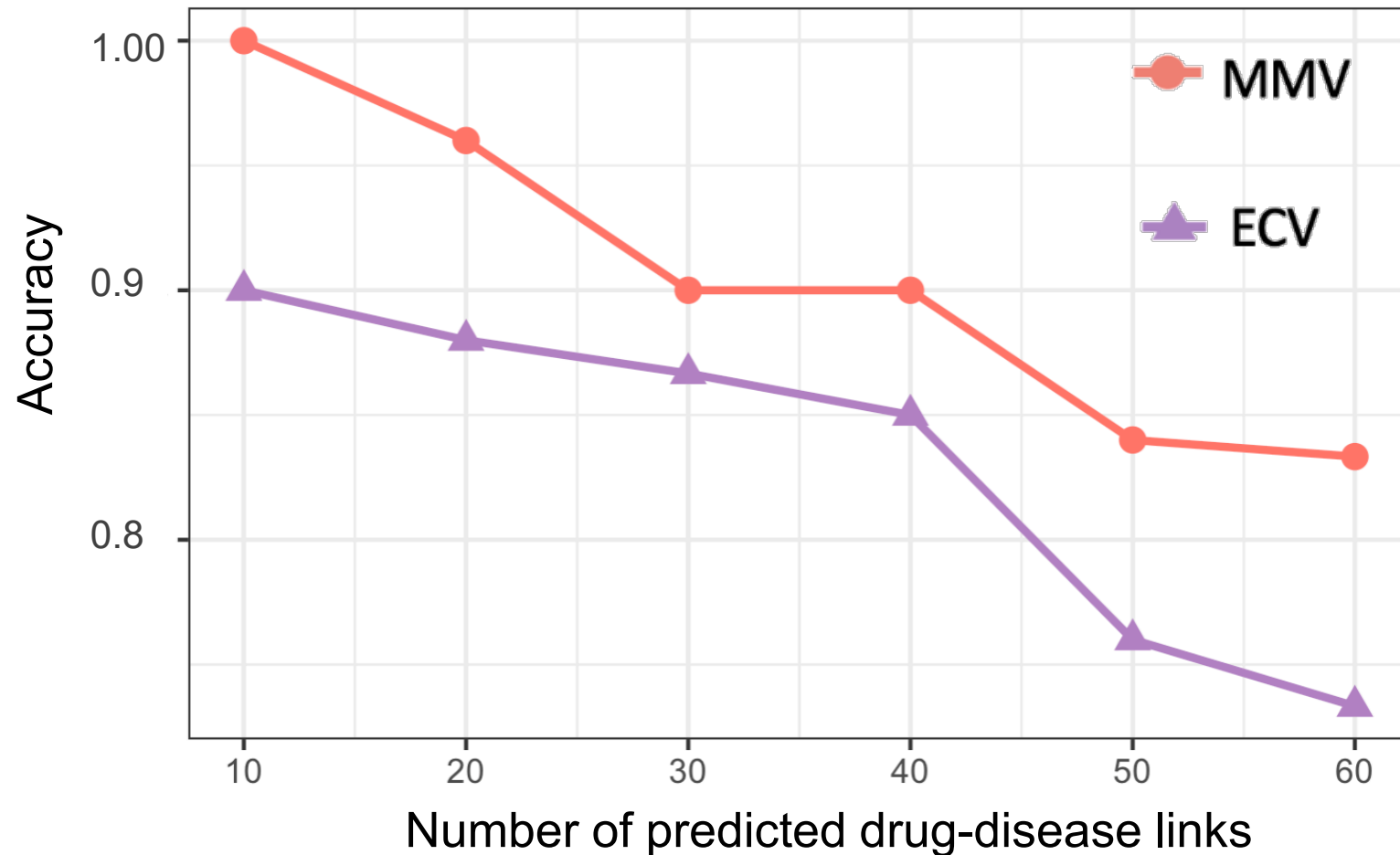
[Rameez Hassan Pirzada](#),<sup>1</sup> [Muhammad Haseeb](#),<sup>1</sup> [Maria Batool](#),<sup>1,2</sup> [MoonSuk Kim](#),<sup>1</sup> and [Sangdun Choi](#)<sup>1,2,\*</sup>

# Comparison Between MMV and ECV

Prediction Accuracy

Time

$$\text{Accuracy} = \frac{\# \text{ validated predictions}}{\# \text{ predictions}}$$



Computational Time (mins)

- **MMV: 28.90 ± 1.30**
- **ECV: 250.65 ± 2.11**

# Take Home Message

## Med-Reader

Med-Reader: Help expedite your research

Home Liquid Hot Topics Network with Significance Score Multi-databases Cross Validation Hypotheses Generation Tutorial About this site


**Medical AI Reader**

Publication Start Date: 2019-12-01

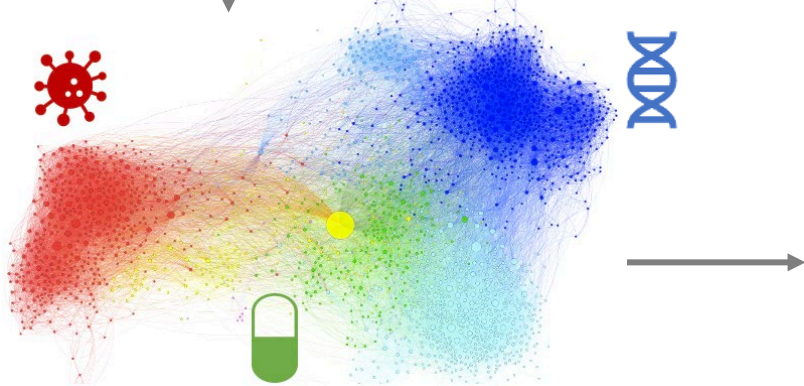
Publication End Date: 2020-03-15

COVID-19

Click 'Submit query' directly to get example results for input query 'COVID-19'.  
Note: This search query will be passed to PubMed to retrieve relevant publications.  
[Click here for PubMed query search help](#)

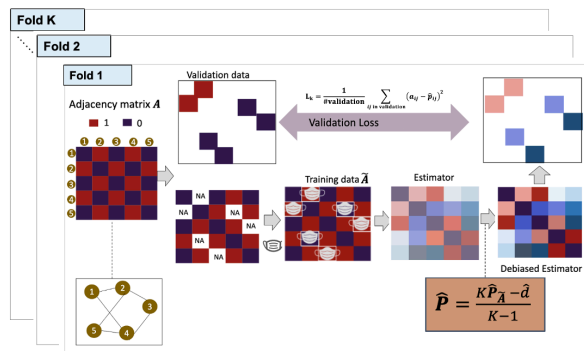


Automatic



Knowledge graph

## Masked mirror validation (MMV)



- ✓ Theoretical guarantee
- ✓ Fast
- ✓ Effective

Better estimation

Graphon Estimation

Promotes

Drug repurposing



Thank you!