



# How Do Hyperedges Overlap in Real-World Hypergraphs? Patterns, Measures, and Generators

---



Geon Lee\*



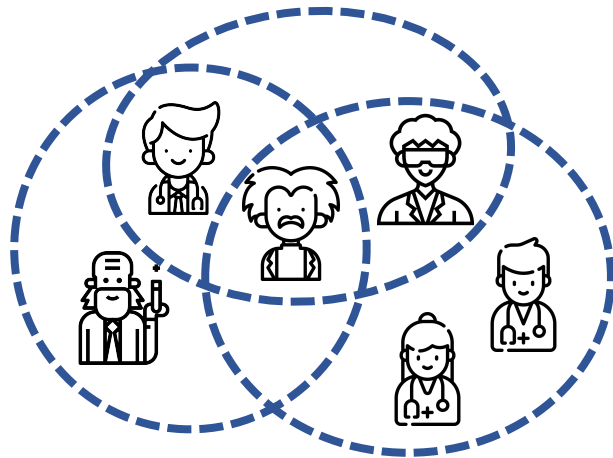
Minyoung Choe\*



Kijung Shin

# Hypergraphs are Everywhere

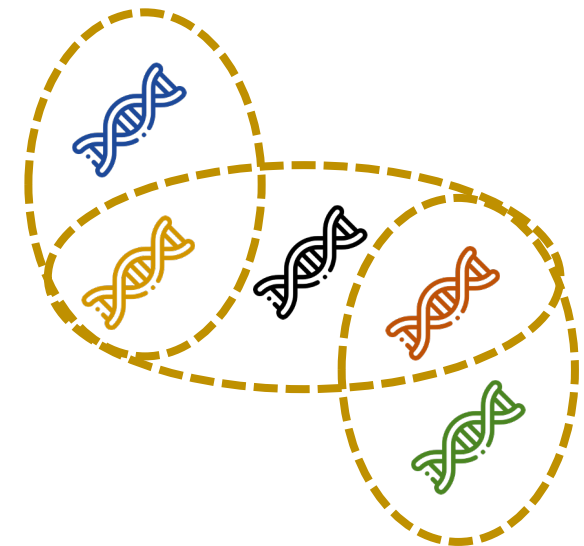
- **Hypergraphs** consist of nodes and hyperedges.
- Each **hyperedge** is a subset of any number of nodes.
- Hyperedges can **overlap** in infinitely many different ways.



**Collaborations of Researchers**



**Co-purchases of Items**

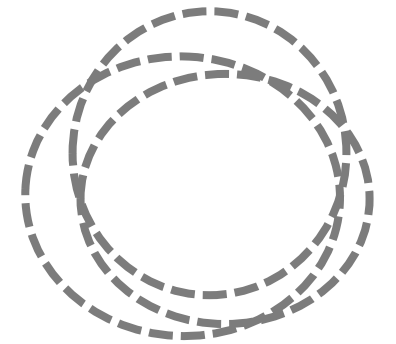
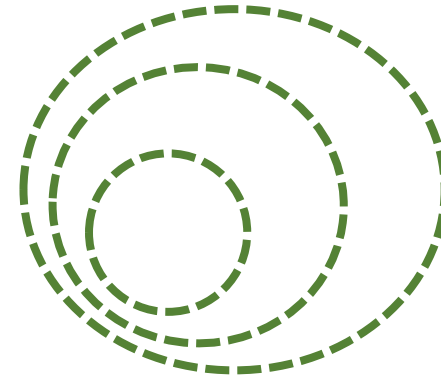
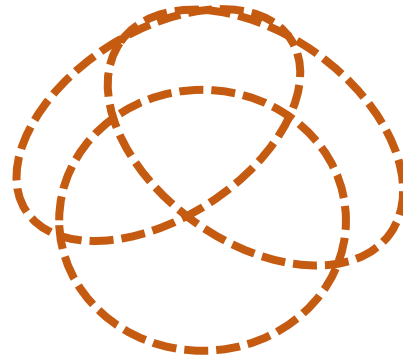
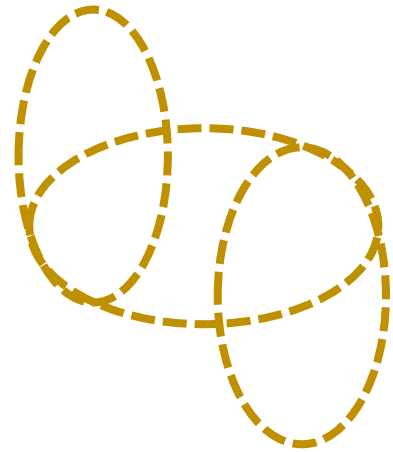
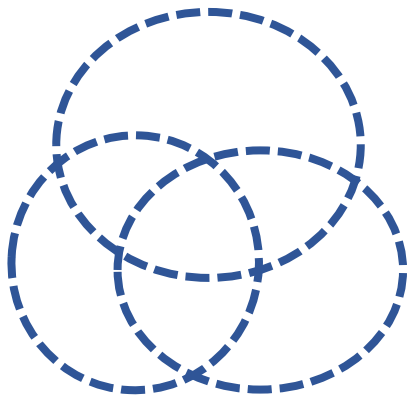


**Joint Interactions of Proteins**

# Our Questions

---

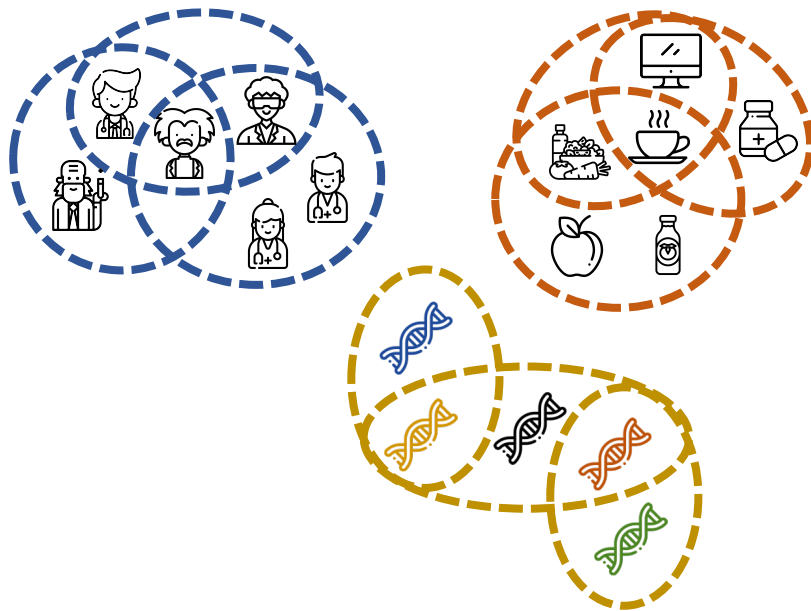
**Q1** How do hyperedges overlap in real-world hypergraphs?



# Our Questions (cont.)

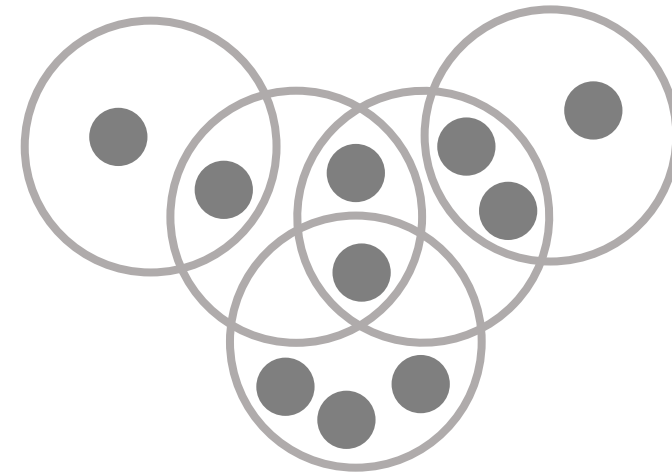
---

**Q2** Are there any non-trivial patterns that distinguish real-world hypergraphs from random hypergraphs?



**Real-world Hypergraphs**

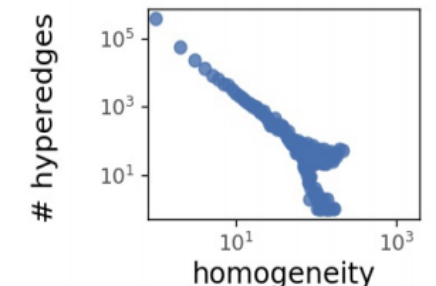
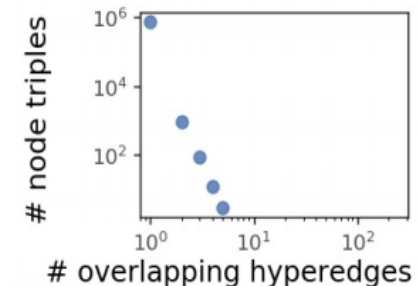
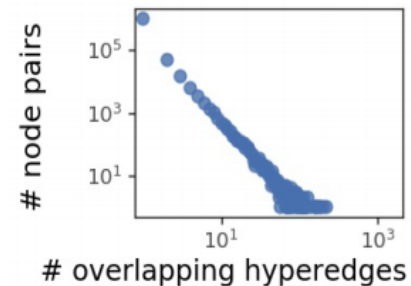
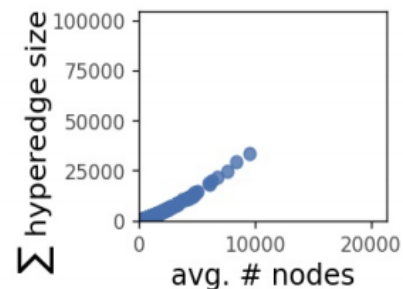
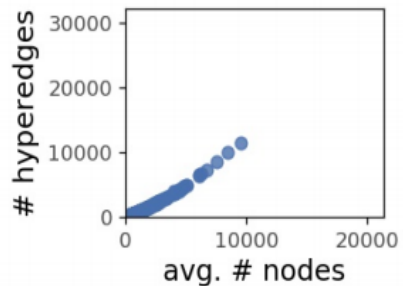
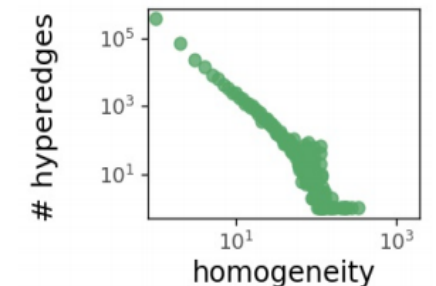
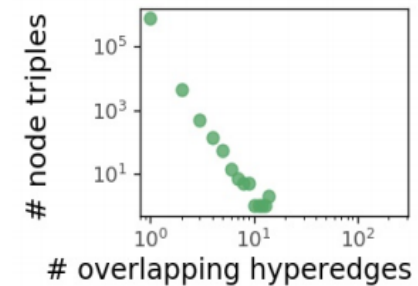
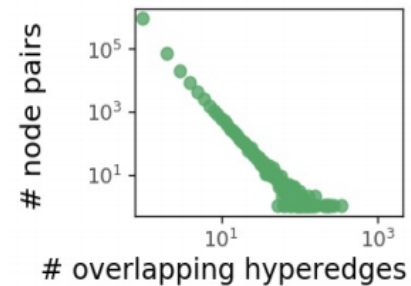
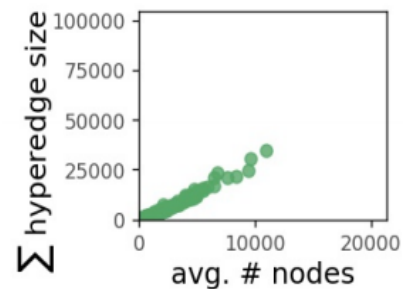
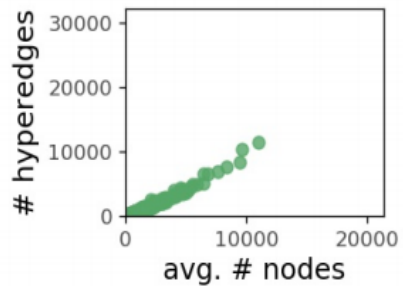
**VS**



**Random Hypergraph**

# Our Questions (cont.)

**Q3** How can we reproduce the patterns through simple mechanisms?

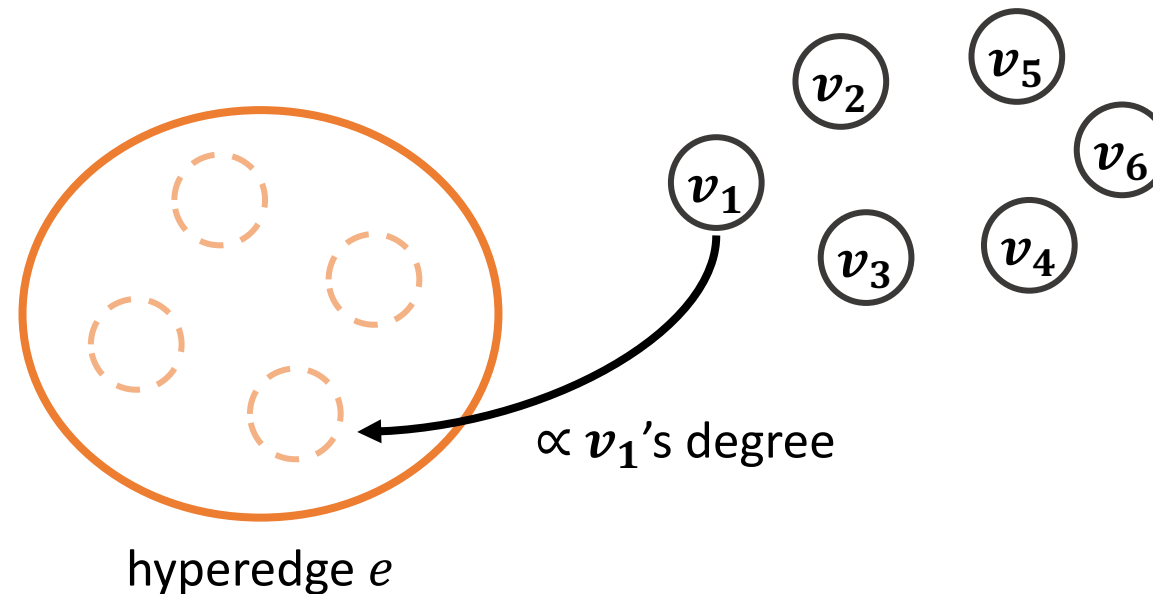


# Null Model

---

## HyperCL: Random Hypergraph Generator (Null Model)

- Nodes are sampled with probability proportional to the degree of each node.
- The degree distribution of nodes is empirically preserved.



# Datasets

---

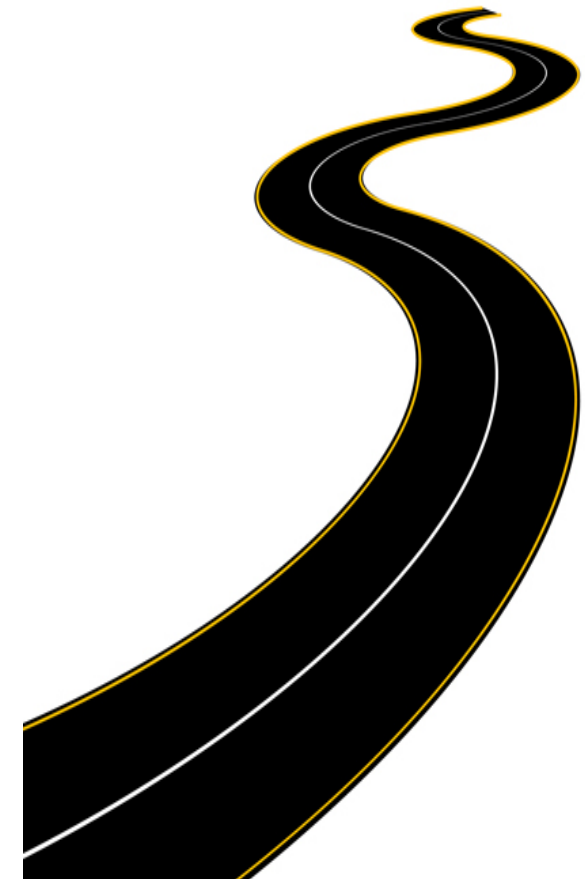
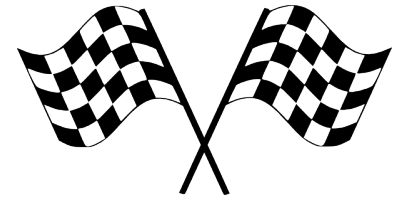
Thirteen real-world hypergraphs from six domains

Domain	Datasets
Email	email-Enron, email-Eu
Contact	contact-primary, contact-high
Drugs	NDC-classes, NDC-substances
Tags	tags-ubnutu, tags-math
Threads	threads-ubuntu, threads-math
Co-authorship	coauth-DBLP, coauth-geology, coauth-history

# Roadmap

---

1. **Observation: Egonet Level**
2. Observation: Pair/Triple of Nodes Level
3. Observation: Hyperedge Level
4. Generators
5. Conclusions

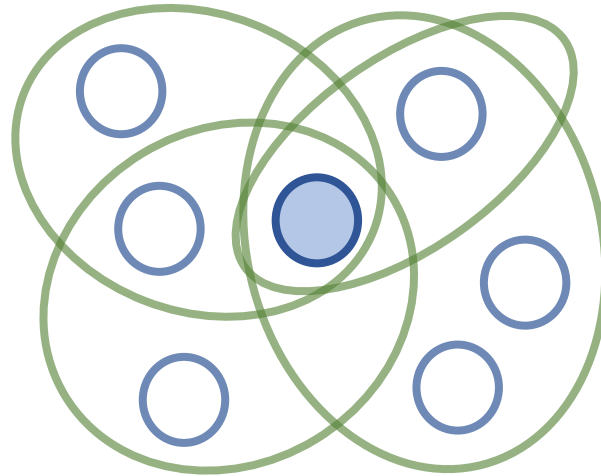




# Observation: Egonet Level

---

How substantially do the hyperedges around a node overlap with each other in the real-world hypergraphs?



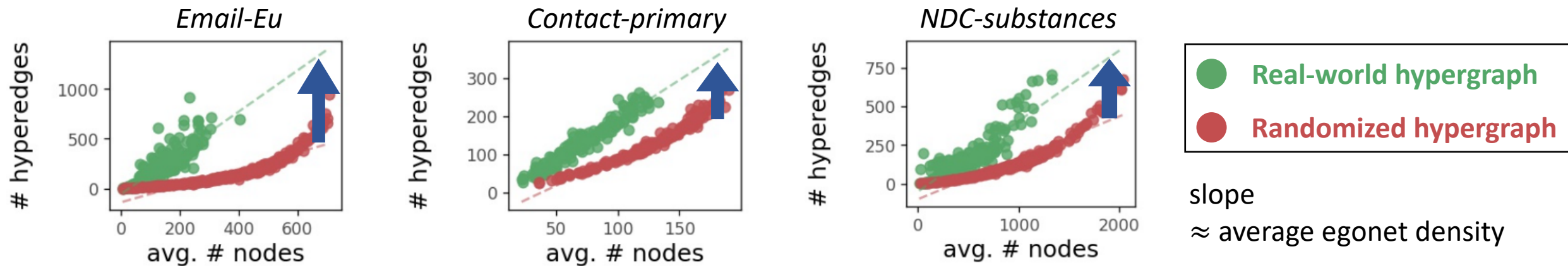
We quantitatively measure this by using **density** and **overlapness**.

# Density of Egonets

**Egonet of a node ( $\mathcal{E}$ ):** set of hyperedges that contains the node

**Density:**  $\rho(\mathcal{E}) := \frac{|\mathcal{E}|}{|\cup_{e \in \mathcal{E}} e|}$

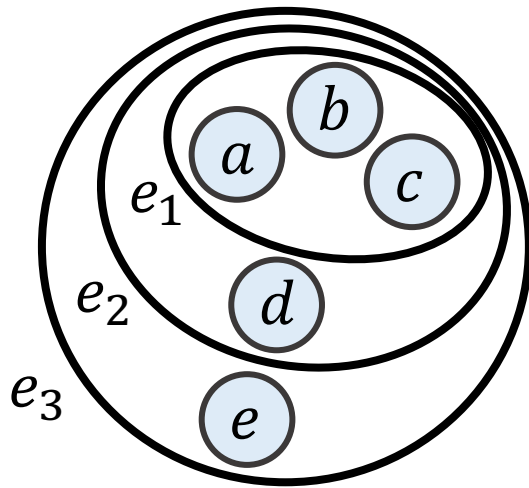
← number of hyperedges  
← number of nodes



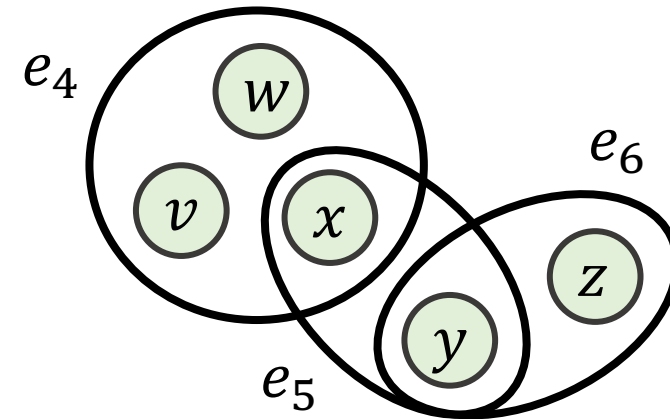
Egonets in real-world hypergraphs tend to have **higher density** than those in randomized ones.

# Density of Egonets (cont.)

Does **density** fully capture the degree of overlaps of a set of hyperedges?



$$\mathcal{E}_1 = \{e_1, e_2, e_3\}$$



$$\mathcal{E}_2 = \{e_4, e_5, e_6\}$$

Our intuition:  $\mathcal{E}_1$  is more overlapped than  $\mathcal{E}_2$ .

Density:  $\rho(\mathcal{E}_1) = \rho(\mathcal{E}_2) = \frac{3}{5}$

**What is the principled measure for evaluating the degree of overlaps of a set of hyperedges?**

# Degree of Hyperedge Overlaps

Any reasonable measure  $f$  of the hyperedge overlaps should satisfy the following axioms.

## Axiom 1: Number of Hyperedges

Consider two sets of hyperedges  $\mathcal{E}$  and  $\mathcal{E}'$ .

If  $\mathcal{E}$  and  $\mathcal{E}'$  have the same (1) hyperedge sizes and (2) number of distinct nodes, but  $\mathcal{E}$  have more hyperedges than  $\mathcal{E}'$ , then  $f(\mathcal{E}) > f(\mathcal{E}')$ .



# Degree of Hyperedge Overlaps (cont.)

Any reasonable measure  $f$  of the hyperedge overlaps should satisfy the following axioms.

## Axiom 2: Number of Distinct Nodes

Consider two sets of hyperedges  $\mathcal{E}$  and  $\mathcal{E}'$ .

If  $\mathcal{E}$  and  $\mathcal{E}'$  have the same (1) number of hyperedges and (2) size distribution of hyperedges, but  $\mathcal{E}$  have less distinct nodes than  $\mathcal{E}'$ , then  $f(\mathcal{E}) > f(\mathcal{E}')$ .



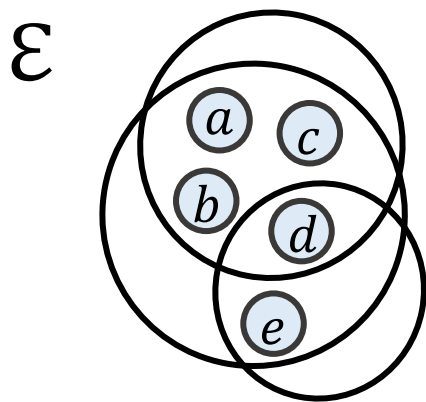
# Degree of Hyperedge Overlaps (cont.)

Any reasonable measure  $f$  of the hyperedge overlaps should satisfy the following axioms.

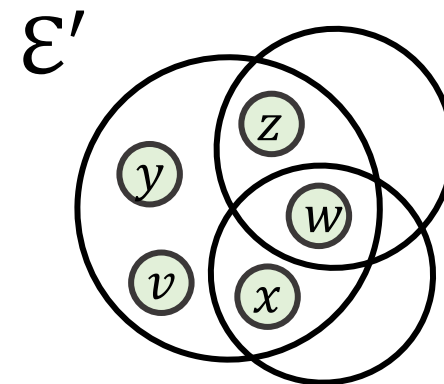
## Axiom 3: Sizes of Hyperedges

Consider two sets of hyperedges  $\mathcal{E}$  and  $\mathcal{E}'$ .

If  $\mathcal{E}$  and  $\mathcal{E}'$  have the same (1) number of distinct nodes and (2) number of hyperedges, but  $\mathcal{E}$  have larger hyperedges than  $\mathcal{E}'$ , then  $f(\mathcal{E}) > f(\mathcal{E}')$ .



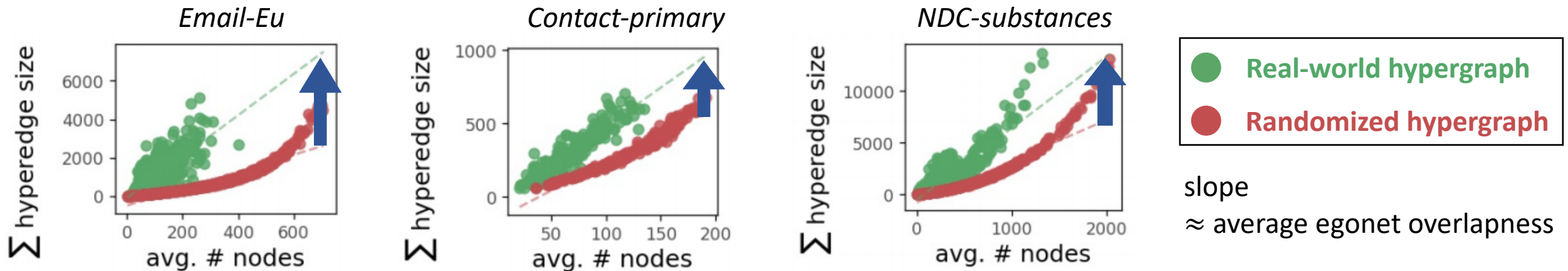
$>$   
more overlap



# Overlapness of Egonets

**Egonet of a node ( $\mathcal{E}$ ):** set of hyperedges that contains the node

**Overlapness:**  $o(\mathcal{E}) := \frac{\sum_{e \in \mathcal{E}} |e|}{|\cup_{e \in \mathcal{E}} e|}$  ← sum of the hyperedge sizes  
← number of nodes



Egonets in real-world hypergraphs tend to have **higher overlapness** than those in randomized ones.

# Overlapness of Egonets (cont.)

---

Overlapness satisfies all the axioms while others does not.

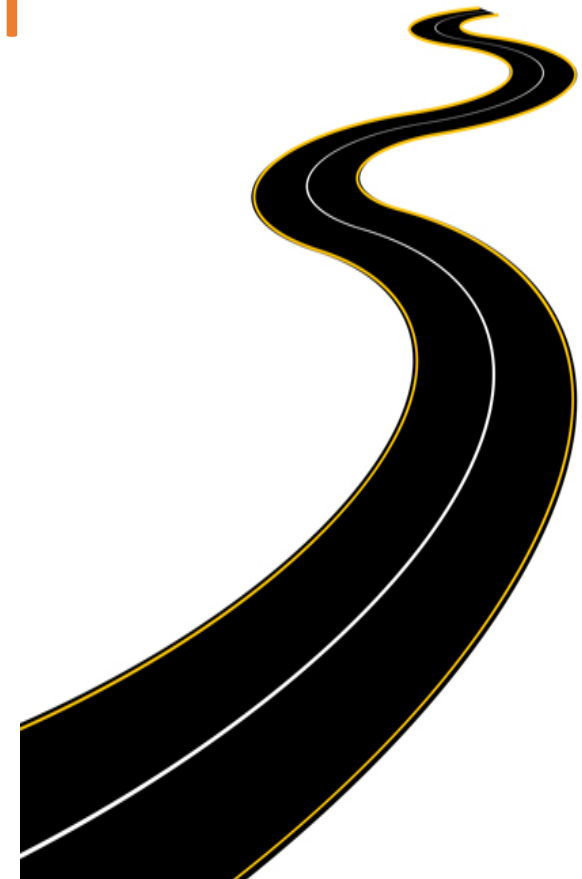
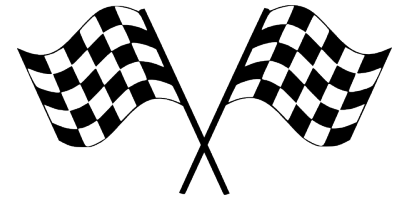
Metric	Axiom 1	Axiom 2	Axiom 3
Intersection	X	X	X
Union Inverse	X	✓	X
Jaccard Index	X	X	X
Overlap Coefficient	X	X	X
Density	✓	✓	X
<b>Overlapness (Proposed)</b>	✓	✓	✓



# Roadmap

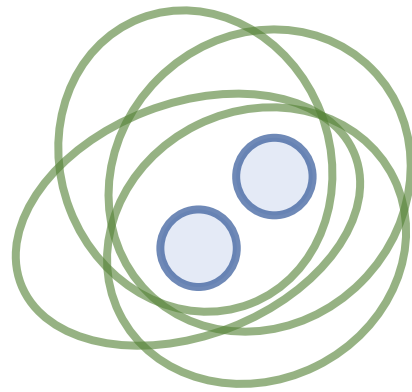
---

1. Observation: Egonet Level
2. **Observation: Pair/Triple of Nodes Level**
3. Observation: Hyperedge Level
4. Generators
5. Conclusions

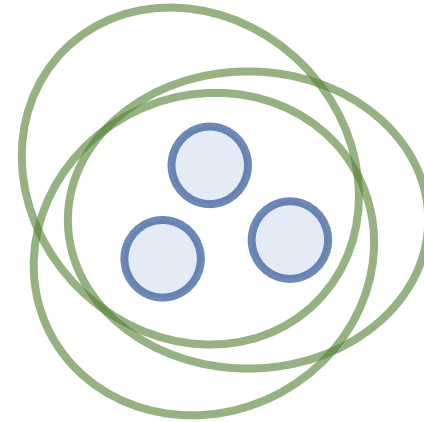


# Observation: Pair/Triple of Nodes Level

How many hyperedges overlap at a pair or triple of nodes in the real-world hypergraphs?



Pair of nodes



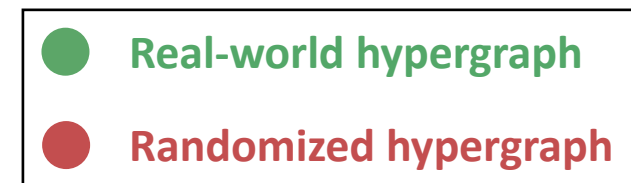
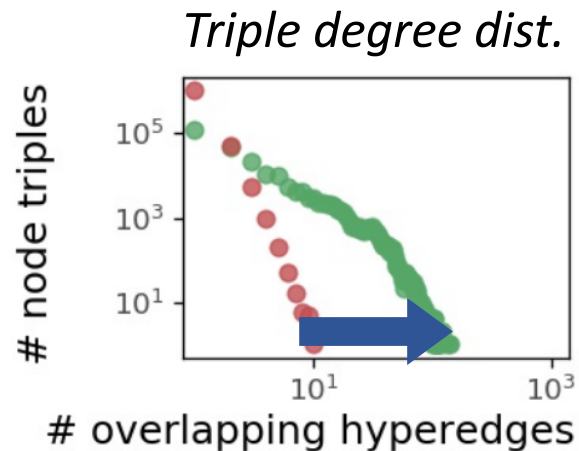
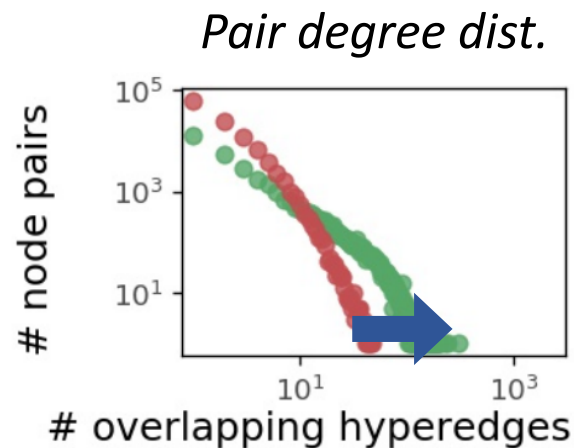
Triple of nodes

We extend the concept of **degree** to pairs and triples of nodes.

# Degree of Node Pair/Triple

$E_S$ : set of hyperedges overlapping at subset  $S$  of nodes

Consider the number of hyperedges overlapping at each **pair or triple of nodes**:  $|E_{\{i,j\}}|$  and  $|E_{\{i,j,k\}}|$ .

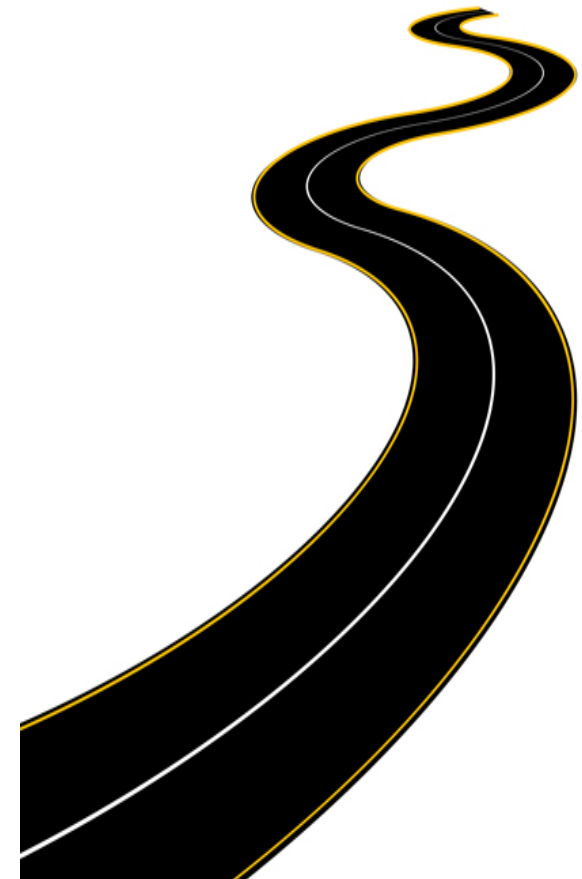
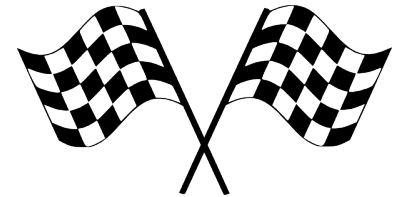


The distribution of the number of hyperedges overlapping at each node pair & triple is **more skewed with a heavier tail** in real-world hypergraphs than in randomized ones.

# Roadmap

---

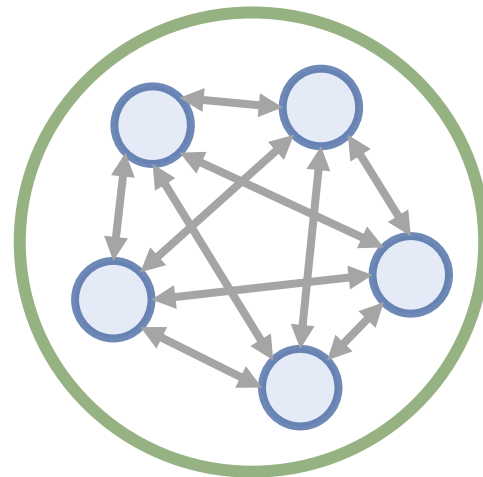
1. Observation: Egonet Level
2. Observation: Pair/Triple of Nodes Level
3. **Observation: Hyperedge Level**
4. Generators
5. Conclusions



# Observation: Hyperedge Level

---

How structurally similar are nodes that form hyperedges together related to each other in the real-world hypergraphs?



We define a new measure to investigate the similarity.

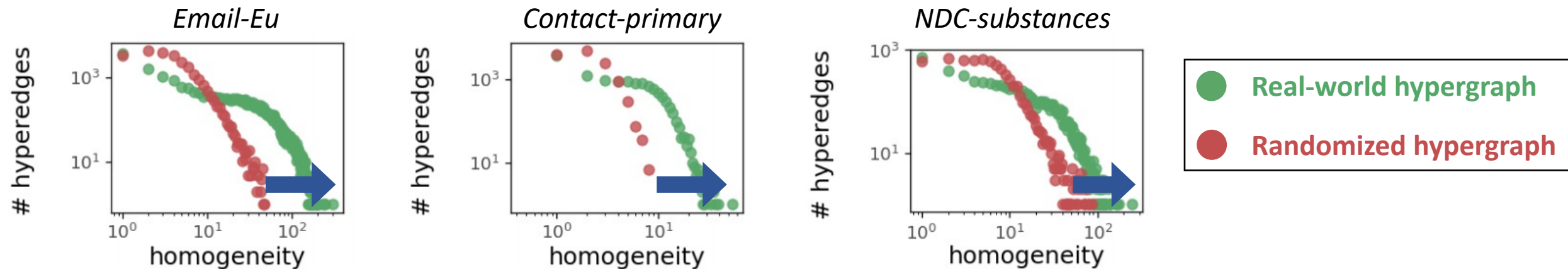
# Homogeneity of a Hyperedge

How to measure the similarity among the nodes forming a hyperedge?

**Homogeneity of a hyperedge:**

$$\text{homogeneity}(e) := \begin{cases} \frac{\sum_{\{u,v\} \in \binom{e}{2}} |E_{\{u,v\}}|}{\binom{|e|}{2}}, & \text{if } |e| > 1 \\ 0, & \text{otherwise} \end{cases}$$

← Average number of hyperedges overlapping at all the pair of nodes in the hyperedge.

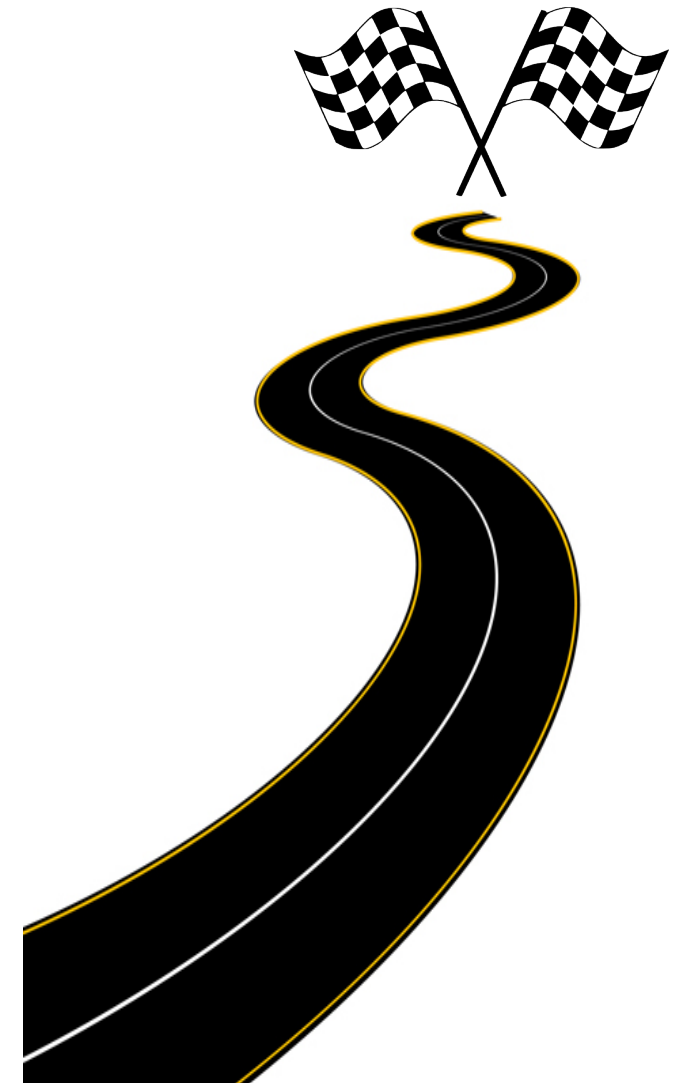


Hyperedges in real-world hypergraphs tend to have **higher homogeneity** than those in randomized ones.

# Roadmap

---

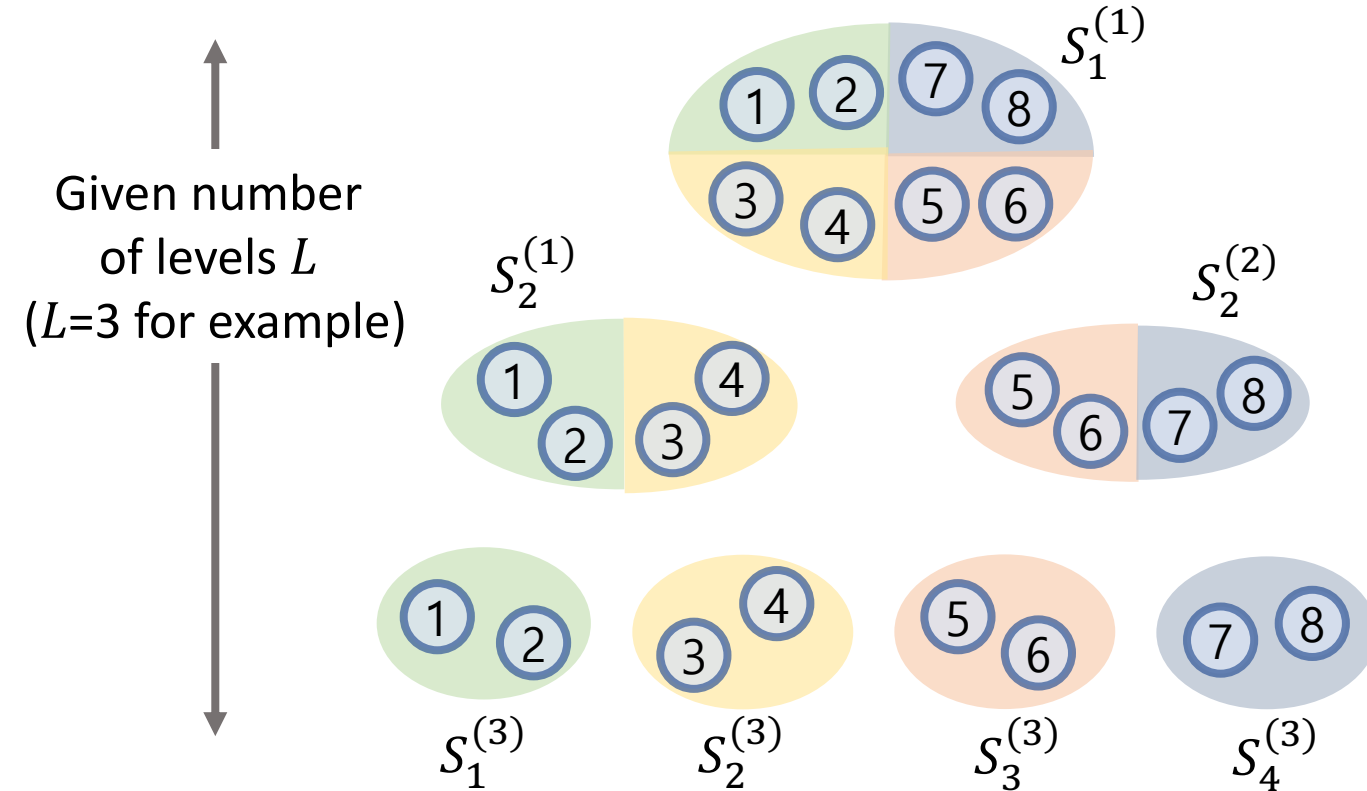
1. Observation: Egonet Level
2. Observation: Pair/Triple of Nodes Level
3. Observation: Hyperedge Level
4. **Generators**
5. Conclusions



# Our Model: HyperLap

**Main Idea:** Extension of HyperCL

**Step 1.** Hierarchical Node Partitioning

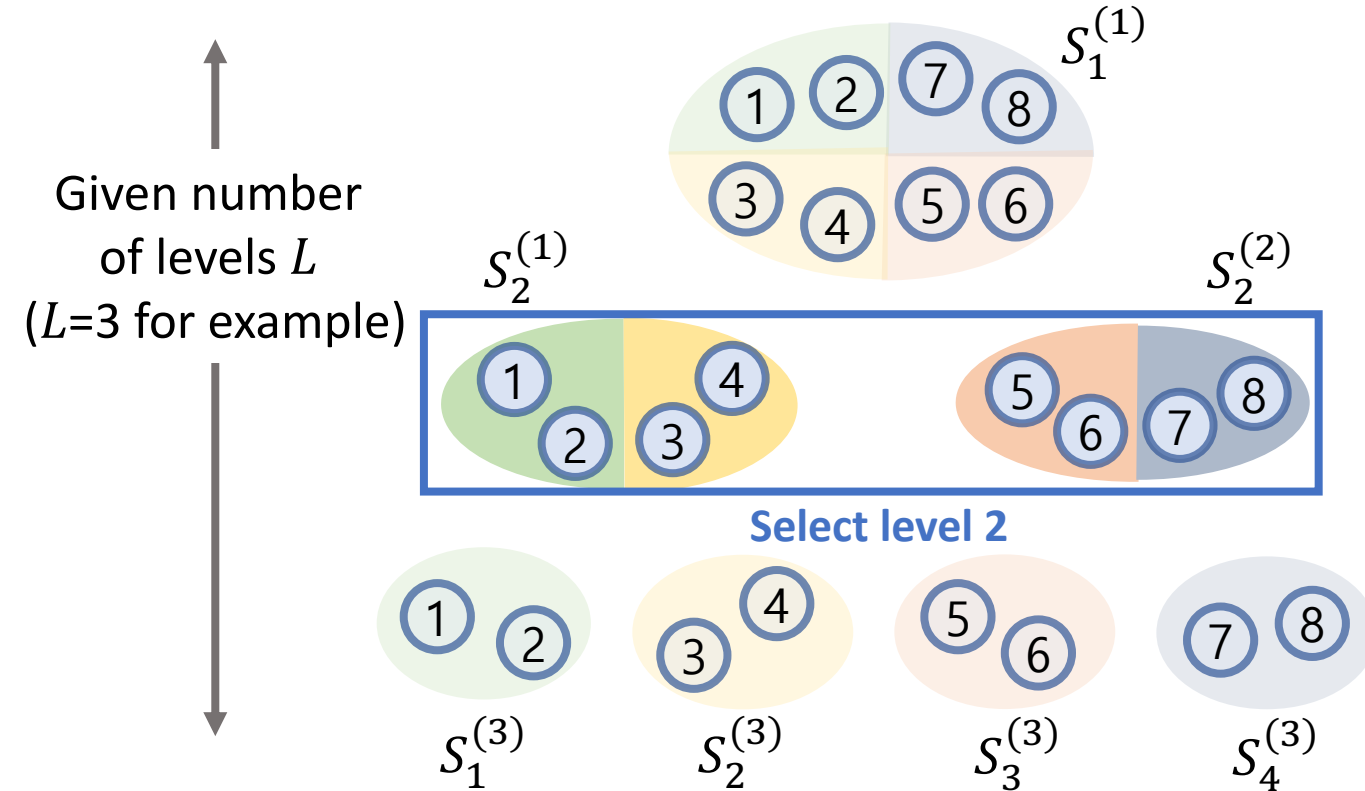




# Our Model: HyperLap (cont.)

**Main Idea:** Extension of HyperCL

**Step 2.** Hyperedge Generation

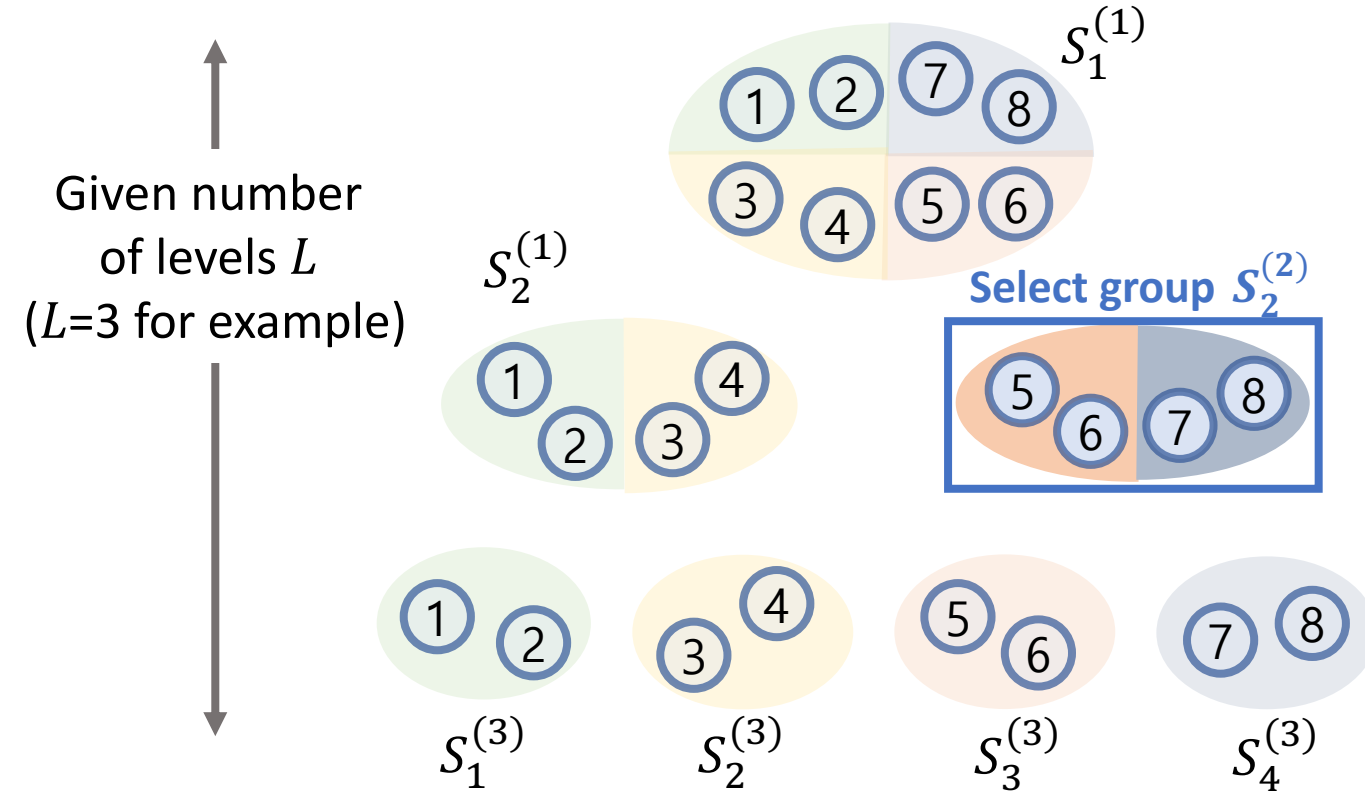


1. Select a level with probability proportional to the given weight of each level  $\{w_1, \dots, w_L\}$ .
2. Select a group uniformly at random.
3. Sample nodes independently with probability proportional to the degree of each node to form a hyperedge.

# Our Model: HyperLap (cont.)

**Main Idea:** Extension of HyperCL

## Step 2. Hyperedge Generation

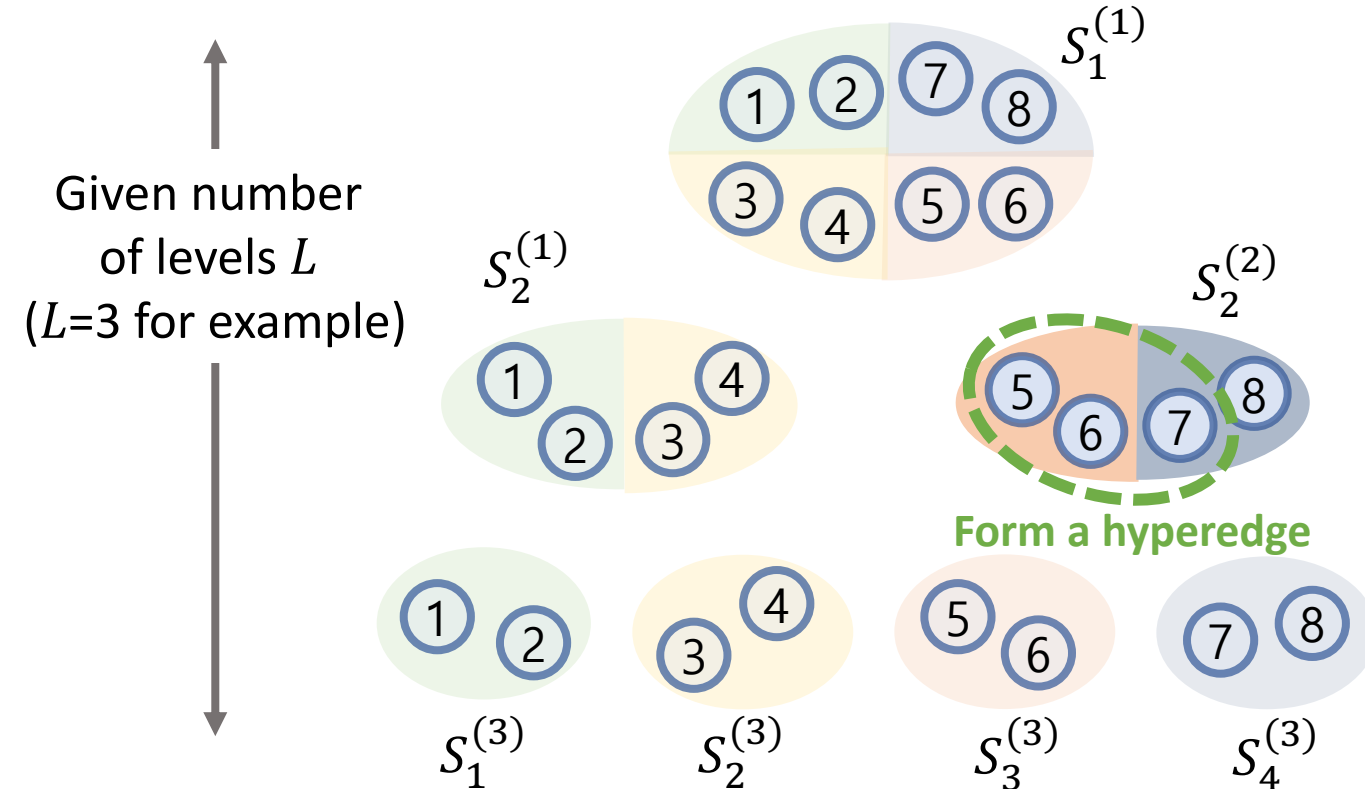


1. Select a level with probability proportional to the given weight of each level  $\{w_1, \dots, w_L\}$ .
2. **Select a group uniformly at random.**
3. Sample nodes independently with probability proportional to the degree of each node to form a hyperedge.

# Our Model: HyperLap (cont.)

**Main Idea:** Extension of HyperCL

## Step 2. Hyperedge Generation



1. Select a level with probability proportional to the given weight of each level  $\{w_1, \dots, w_L\}$ .
2. Select a group uniformly at random.
3. **Sample nodes independently with probability proportional to the degree of each node to form a hyperedge.**

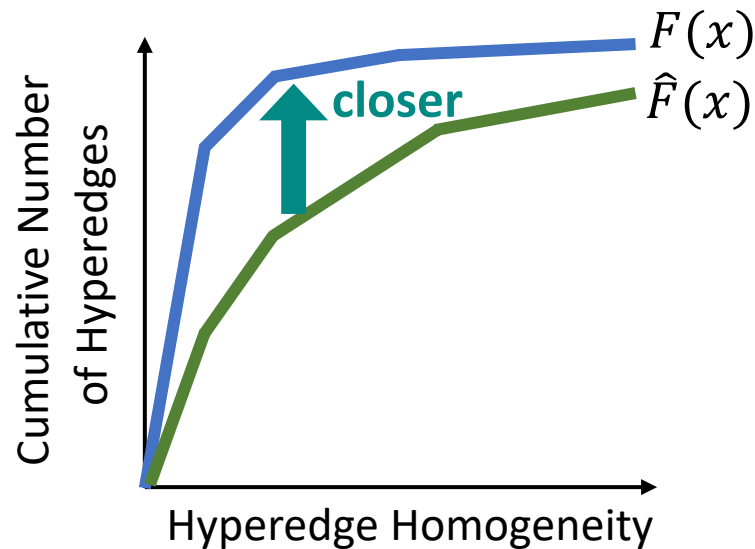
# Our Model: HyperLap<sup>+</sup>

**HyperLap<sup>+</sup>** automatically tunes the parameters of HyperLap.

**Objective:** Minimize the **hyperedge homogeneity distance**  $HHD(G, \hat{G})$ .

$$HHD(G, \hat{G}) = \max_x \{|F(x) - \hat{F}(x)|\}$$

where  $F$  and  $\hat{F}$  are the cumulative hyperedge homogeneity distribution of hypergraphs  $G$  and  $\hat{G}$ , respectively.



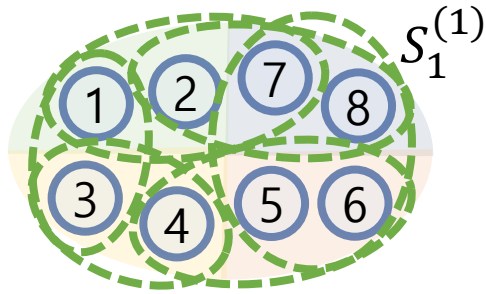
$$\min_{w_1, \dots, w_L} HHD(G, \hat{G}) \text{ where } w_1 + \dots + w_L = 1$$

Learnable parameters

# Our Model: HyperLap<sup>+</sup> (cont.)

---

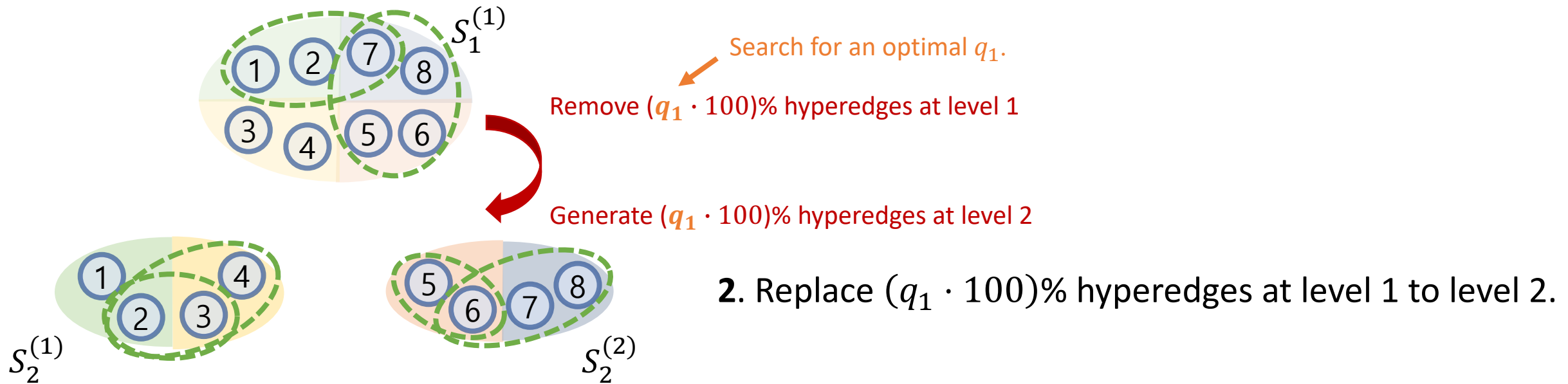
**HyperLap<sup>+</sup>** automatically tunes the parameters of HyperLap.



1. Generate  $|E|$  hyperedges at level 1.

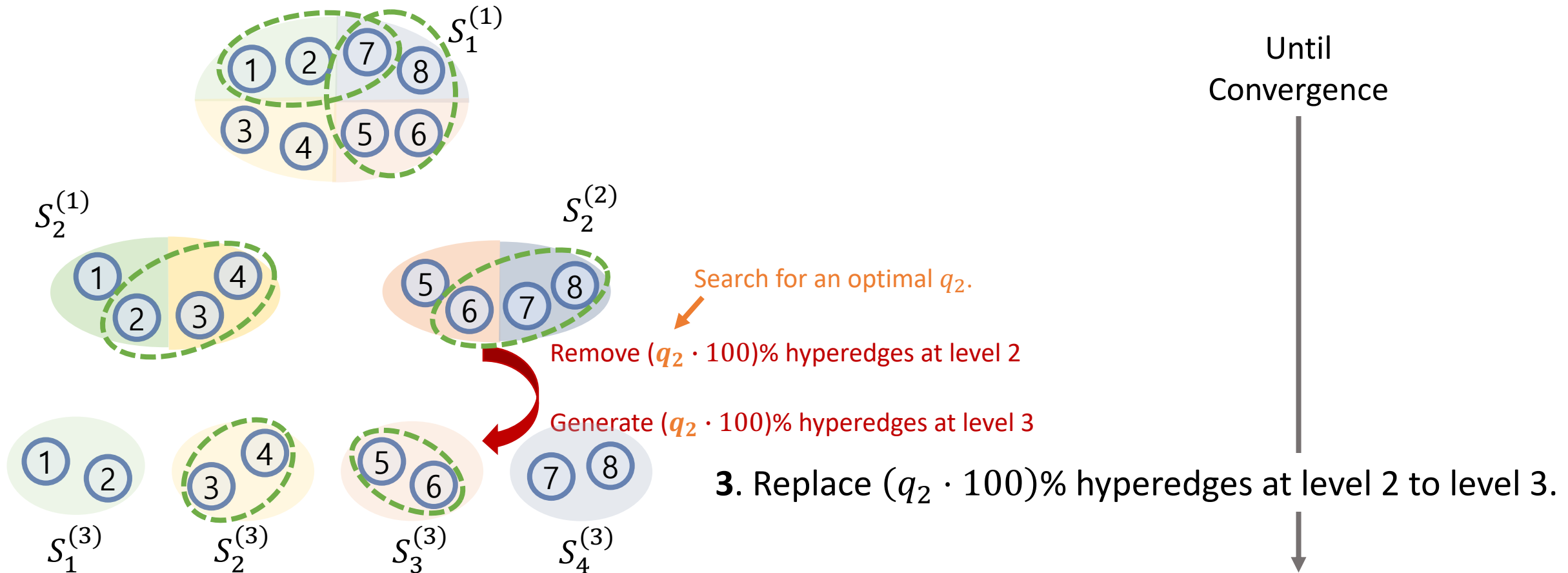
# Our Model: HyperLap<sup>+</sup> (cont.)

**HyperLap<sup>+</sup>** automatically tunes the parameters of HyperLap.



# Our Model: HyperLap<sup>+</sup> (cont.)

HyperLap<sup>+</sup> automatically tunes the parameters of HyperLap.



# Evaluation of Our Model

**HyperLap<sup>+</sup>** reproduces most accurately the distributions of  
(1) egonet density (2) egonet overlapness (3) hyperedge homogeneity

Measure: the similarity between the distributions derived from the real-world and the generated hypergraph by Kolmogorov-Smirnov D-statistics.

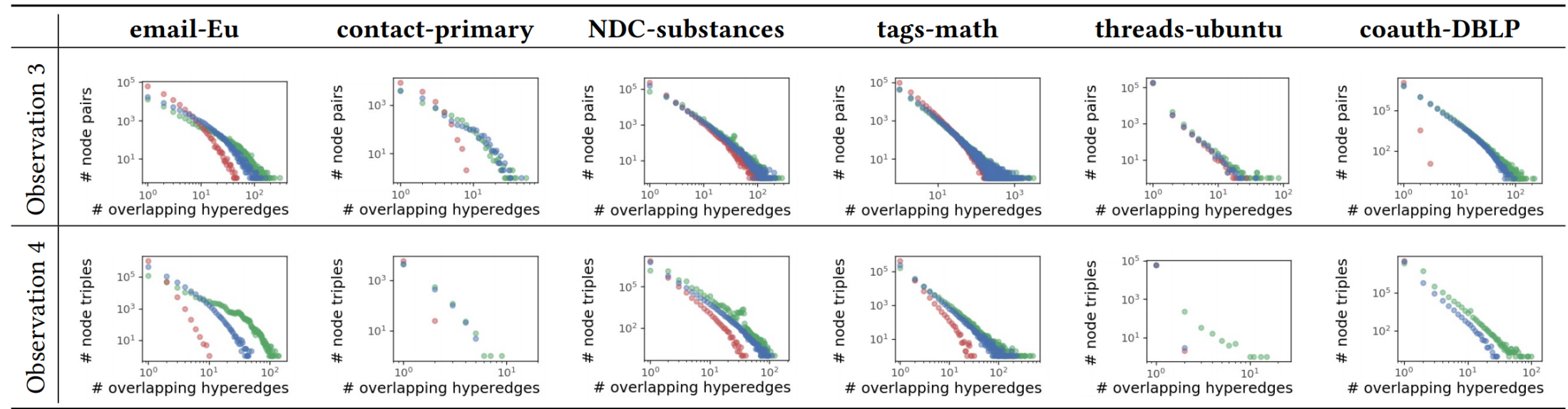
Dataset	Density of Egonets (Obs. 1)					Overlapness of Egonets (Obs. 2)					Homogeneity of Hyperedges (Obs. 5)				
	H-CL	H-PA	H-FF	H-LAP	H-LAP <sup>+</sup>	H-CL	H-PA	H-FF	H-LAP	H-LAP <sup>+</sup>	H-CL	H-PA	H-FF	H-LAP	H-LAP <sup>+</sup>
email-Enron	0.545	0.202	0.391	0.405	<b>0.125</b>	0.517	0.398	0.398	0.391	<b>0.111</b>	0.498	0.241	0.656	0.191	<b>0.136</b>
email-Eu	0.724	-	0.402	0.577	<b>0.310</b>	0.534	-	0.639	0.432	<b>0.197</b>	0.505	-	0.688	0.247	<b>0.168</b>
contact-primary	0.896	0.537	0.975	0.334	<b>0.128</b>	0.867	0.471	0.942	0.285	<b>0.095</b>	0.430	0.236	0.484	<b>0.142</b>	0.188
contact-high	0.948	0.529	0.880	0.522	<b>0.345</b>	0.874	0.431	0.703	0.486	<b>0.296</b>	0.423	0.196	0.336	<b>0.120</b>	0.178
NDC-classes	0.694	0.785	0.731	0.696	<b>0.635</b>	0.302	0.715	0.406	<b>0.231</b>	0.248	0.274	0.410	0.484	0.272	<b>0.225</b>
NDC-substances	0.451	-	0.801	0.426	<b>0.366</b>	0.321	-	0.338	0.243	<b>0.157</b>	0.377	-	0.740	0.262	<b>0.108</b>
tags-ubuntu	0.522	<b>0.162</b>	0.216	0.410	0.300	0.432	<b>0.117</b>	0.398	0.487	0.210	0.245	0.136	0.844	0.105	<b>0.011</b>
tags-math	0.496	0.350	0.561	<b>0.195</b>	0.227	0.460	0.325	0.709	<b>0.151</b>	0.186	0.337	0.217	0.921	0.086	<b>0.015</b>
threads-ubuntu	0.159	0.856	-	0.163	<b>0.159</b>	0.299	0.953	-	0.300	<b>0.297</b>	0.020	0.291	-	0.016	<b>0.011</b>
threads-math	0.137	0.492	-	<b>0.120</b>	0.135	0.232	0.714	-	0.235	<b>0.229</b>	0.060	0.368	-	0.102	<b>0.019</b>
coauth-DBLP	0.228	-	-	0.227	<b>0.132</b>	0.302	-	-	0.267	<b>0.244</b>	0.715	-	-	0.540	<b>0.026</b>
coauth-geology	0.200	-	-	0.202	<b>0.138</b>	<b>0.248</b>	-	-	0.252	0.266	0.624	-	-	0.481	<b>0.044</b>
coauth-history	<b>0.087</b>	-	-	0.090	0.089	<b>0.316</b>	-	-	0.321	0.324	0.154	-	-	0.125	<b>0.020</b>
<b>Average</b>	0.468	0.489	0.619	0.335	<b>0.237</b>	0.439	0.515	0.566	0.313	<b>0.219</b>	0.358	0.261	0.644	0.206	<b>0.088</b>

-: out of time (taking more than 10 hours) or out of memory



# Evaluation of Our Model (cont.)

**HyperLap<sup>+</sup>** reproduces the heavy-tailed distributions of the number of overlapping hyperedges at each *pair* and each *triple* of nodes accurately.



# Evaluation of Our Model (cont.)

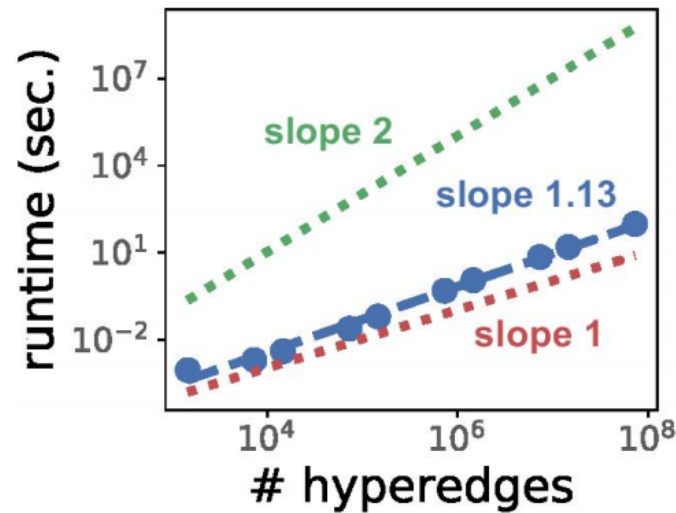
**HyperLap<sup>+</sup>** reproduces the heavy-tailed distributions of the number of overlapping hyperedges at each *pair* and each *triple* of nodes accurately.

Dataset	Pair of Nodes (Obs. 3)								Triple of Nodes (Obs. 4)							
	Distance from Real (D-statistics)					Heavy-tail Test			Distance from Real (D-statistics)					Heavy-tail Test		
	H-CL	H-PA	H-FF	H-LAP	H-LAP <sup>+</sup>	pw	tpw	logn	H-CL	H-PA	H-FF	H-LAP	H-LAP <sup>+</sup>	pw	tpw	logn
email-Enron	0.143	<b>0.056</b>	0.217	0.075	0.139	-2.37	-0.29	-1.53	0.089	0.295	0.136	<b>0.061</b>	0.072	-0.22	<b>0.38</b>	<b>0.24</b>
email-Eu	0.225	-	0.352	0.162	<b>0.066</b>	<b>0.24</b>	<b>2.75</b>	<b>2.53</b>	0.480	-	0.516	0.337	<b>0.206</b>	<b>0.41</b>	<b>2.11</b>	<b>1.96</b>
contact-primary	0.196	0.062	0.223	0.070	<b>0.051</b>	<b>9.53</b>	<b>15.74</b>	<b>13.92</b>	0.137	0.061	0.110	0.053	<b>0.031</b>	-1.86	-1.27	<b>1.23</b>
contact-high	0.277	<b>0.062</b>	0.141	0.127	0.067	-3.09	-0.95	-0.06	0.210	<b>0.131</b>	0.182	0.182	0.193	-3.95	-	<b>0.50</b>
NDC-classes	0.273	0.197	0.196	0.246	<b>0.172</b>	<b>12.15</b>	<b>14.42</b>	<b>14.04</b>	0.376	<b>0.167</b>	0.405	0.349	0.286	<b>3.22</b>	<b>7.92</b>	<b>7.34</b>
NDC-substances	0.272	-	0.244	0.251	<b>0.202</b>	<b>33.69</b>	<b>40.13</b>	<b>39.66</b>	0.521	-	0.591	0.492	<b>0.453</b>	<b>45.30</b>	<b>55.38</b>	<b>54.99</b>
tags-ubuntu	0.091	<b>0.019</b>	0.182	0.034	0.033	<b>42.33</b>	<b>43.70</b>	<b>43.55</b>	0.148	0.067	0.191	<b>0.020</b>	0.074	<b>14.25</b>	<b>15.57</b>	<b>15.43</b>
tags-math	0.095	0.066	0.278	0.073	<b>0.011</b>	<b>42.75</b>	<b>45.60</b>	<b>45.41</b>	0.209	<b>0.053</b>	0.286	0.113	0.079	<b>21.38</b>	<b>23.12</b>	<b>22.99</b>
threads-ubuntu	0.011	0.137	-	<b>0.008</b>	0.009	<b>1.28</b>	<b>1.75</b>	<b>1.75</b>	<b>0.004</b>	0.130	-	<b>0.004</b>	<b>0.004</b>	-1,346	-1.72	-1.72
threads-math	0.041	0.163	-	<b>0.014</b>	0.033	<b>15.79</b>	<b>16.66</b>	<b>16.52</b>	0.006	0.138	-	<b>0.001</b>	0.005	-1.49	-0.98	<b>0.96</b>
coauth-DBLP	0.224	-	-	0.191	<b>0.032</b>	<b>55.86</b>	<b>74.95</b>	<b>73.45</b>	0.215	-	-	0.214	<b>0.192</b>	<b>2.87</b>	<b>6.73</b>	<b>6.46</b>
coauth-geology	0.178	-	-	0.157	<b>0.040</b>	<b>31.13</b>	<b>45.08</b>	<b>44.06</b>	0.086	-	-	0.085	<b>0.069</b>	-0.10	<b>1.10</b>	<b>0.84</b>
coauth-history	0.033	-	-	0.030	<b>0.009</b>	<b>1.74</b>	<b>1.77</b>	<b>1.63</b>	<b>0.001</b>	-	-	<b>0.001</b>	<b>0.001</b>	-0.86	-	<b>0.57</b>
<b>Average</b>	0.158	0.095	0.229	0.110	<b>0.066</b>				0.193	0.130	0.302	0.147	<b>0.128</b>			

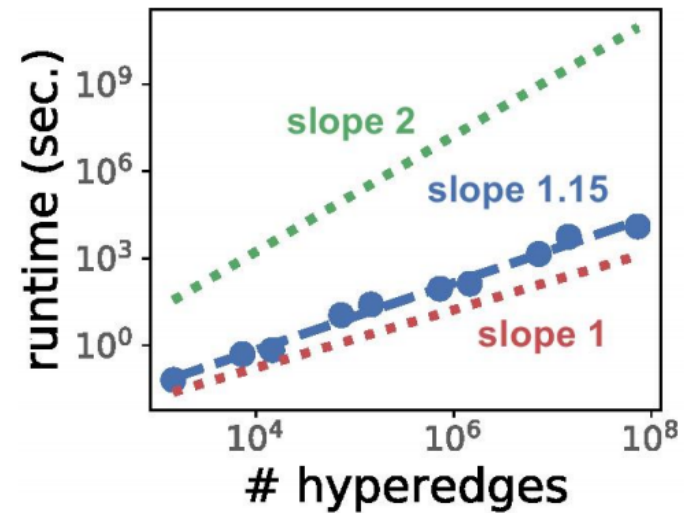
--: out of time (taking more than 10 hours) or out of memory

# Scalability of Our Model

**HyperLap** and **HyperLap<sup>+</sup>** scale near linearly with the size of the considered hypergraph.



(a) HYPERLAP (generation)

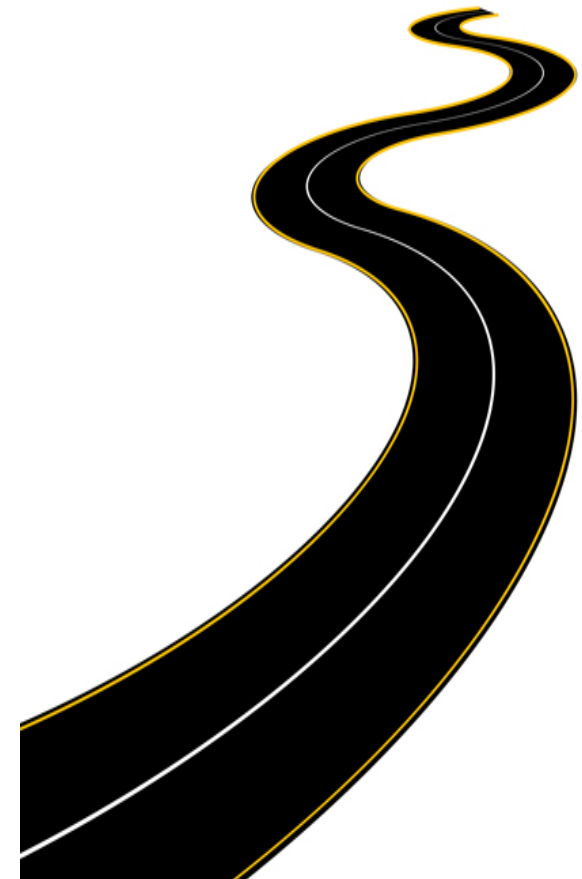
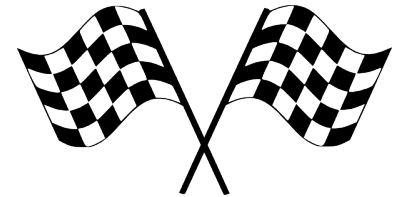


(b) HYPERLAP<sup>+</sup> (generation & fitting)

# Roadmap

---

1. Observation: Egonet Level
2. Observation: Pair/Triple of Nodes Level
3. Observation: Hyperedge Level
4. Generators
5. **Conclusions**



# Conclusions

---

Our contributions in this work:

- ✓ Observations in Real-world Hypergraphs
  1. Egonet Level Observation
  2. Pair/Triple of Nodes Level Observation
  3. Hyperedge Level Observation
- ✓ Novel Measures
- ✓ Realistic Generative Models (HyperLap & HyperLap<sup>+</sup>)

The code and datasets used in the paper are available at  
<https://github.com/young917/www21-hyperlap>



# How Do Hyperedges Overlap in Real-World Hypergraphs? Patterns, Measures, and Generators

---



Geon Lee\*



Minyoung Choe\*



Kijung Shin