

Quantifying Network Similarity using Graph Cumulants

Gecia Bravo-Hermsdorff¹ Lee M. Gunderson¹ Pierre-André Maugis² Carey E. Priebe³

¹Department of Statistical Science, University College London ²Google Research Zürich ³Department of Applied Mathematics and Statistics, Johns Hopkins University

TL;DR:

Graph cumulants perform better and are more intuitive than the typical subgraph statistics

Using only **moments** is awkward

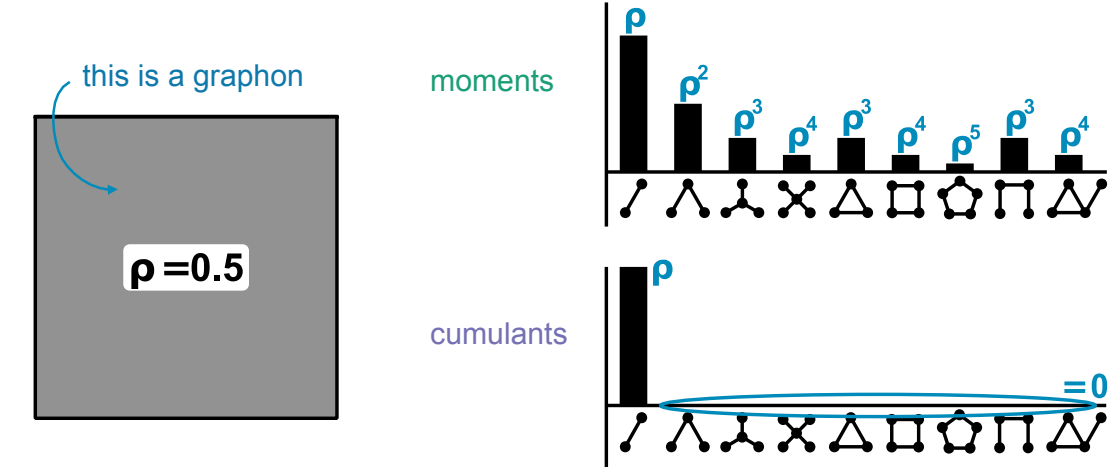
“The length of a human averages 1.7 meters, and their **average squared length** is 2.9 square meters”

Using **cumulants** is **easier** to understand

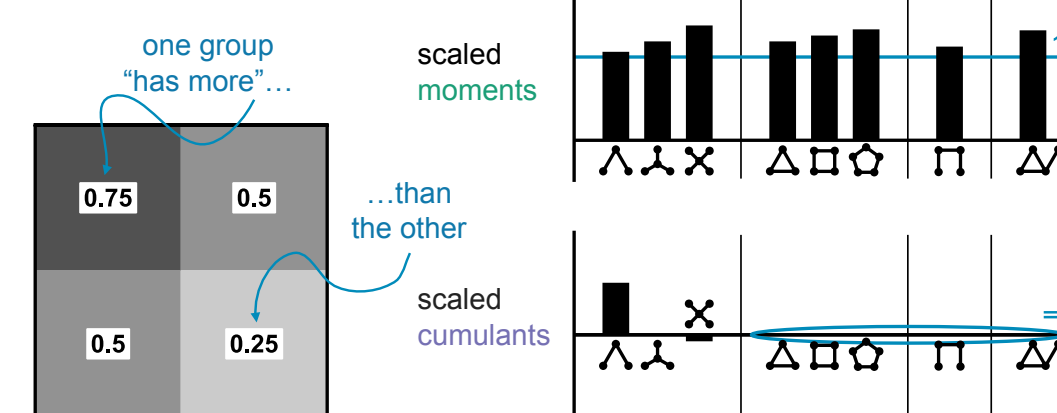
“The length of a human has: a **variance** of 0.01 meters, a **standard deviation** of 0.1 meters, a **relative fluctuation** of 6%”

Graph Cumulants: the Better Subgraph Statistics

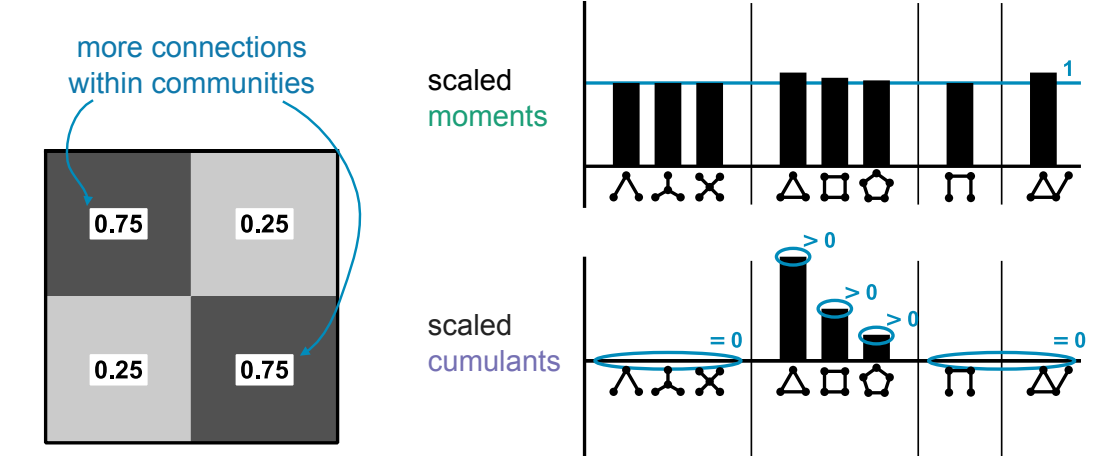
Erdős-Rényi is the new Gaussian...



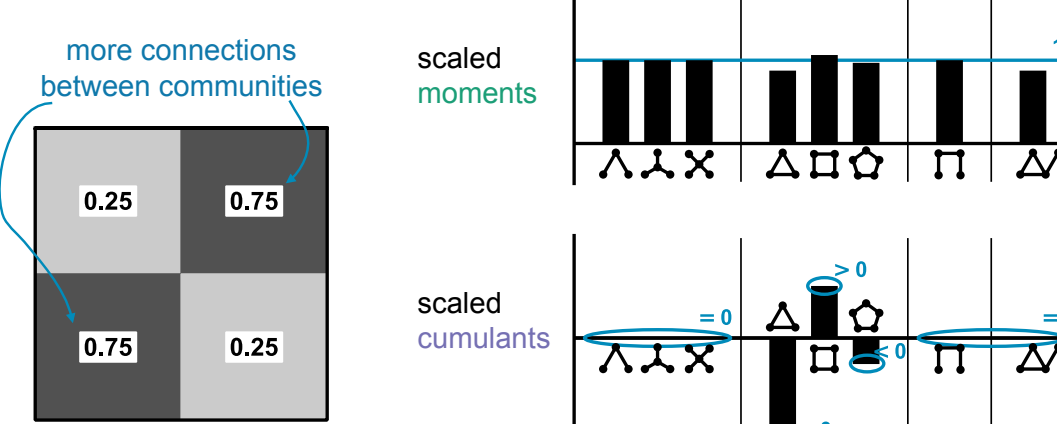
...degrees are encoded in stars...



...clustering is encoded in cycles...

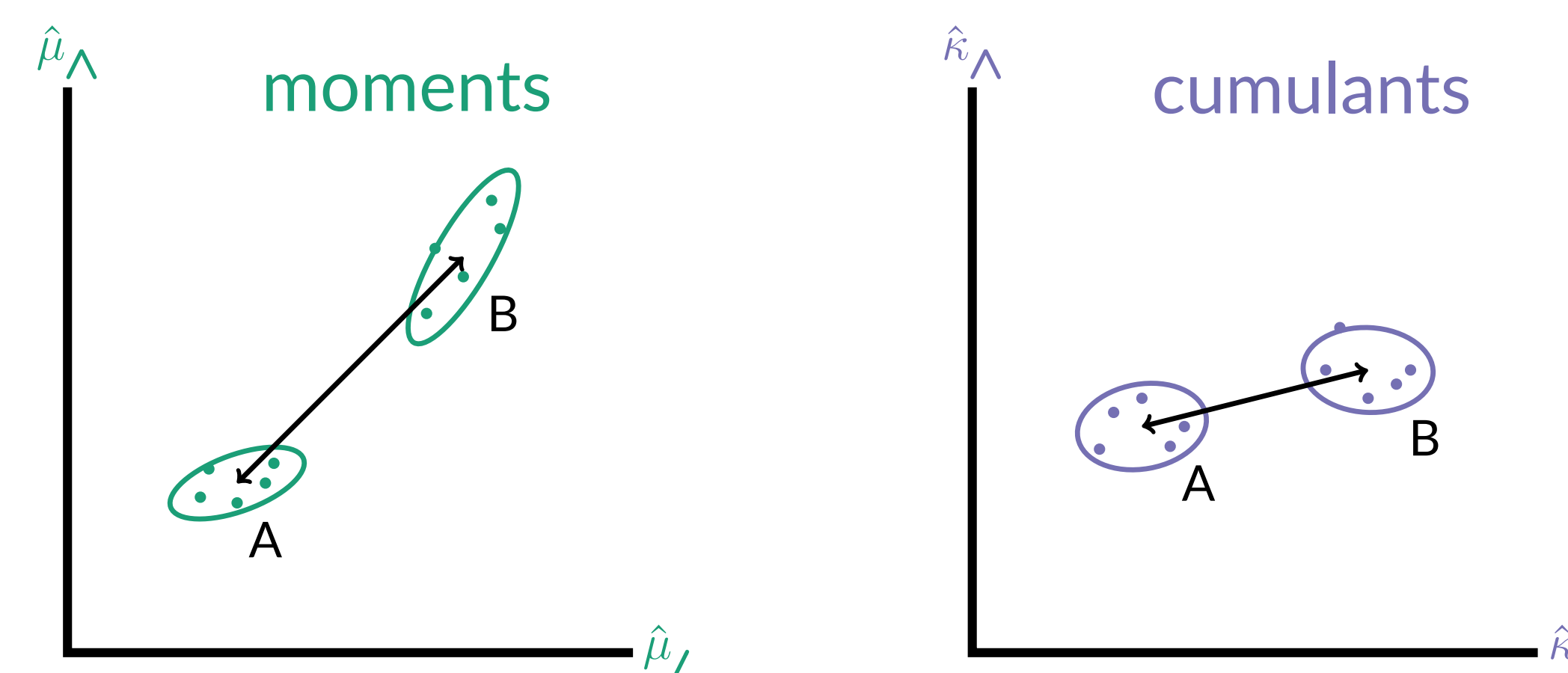


...and bipartite-ness as well!



But are **graph cumulants** better for testing?
(than **subgraph densities**)

Apples-to-Apples: A Two-Sample Test



$$\hat{d}_{\kappa}^2(A, B) = \left(\hat{\kappa}(A) - \hat{\kappa}(B) \right)^{\top} \left(\hat{\Sigma}_{\kappa}(A) + \hat{\Sigma}_{\kappa}(B) \right)^{-1} \left(\hat{\kappa}(A) - \hat{\kappa}(B) \right)$$

$$\hat{d}_{\mu}^2(A, B) = \left(\hat{\mu}(A) - \hat{\mu}(B) \right)^{\top} \left(\hat{\Sigma}_{\mu}(A) + \hat{\Sigma}_{\mu}(B) \right)^{-1} \left(\hat{\mu}(A) - \hat{\mu}(B) \right)$$

“Z²-score” difference metric difference

When estimating the covariance...

$$\text{Cov}(\hat{\mu}_g, \hat{\mu}_{g'}) = \underbrace{\langle \hat{\mu}_g \hat{\mu}_{g'} \rangle}_{\text{“hard”}} - \underbrace{\langle \hat{\mu}_g \rangle \langle \hat{\mu}_{g'} \rangle}_{\text{“easy”}}$$

...the “hard” part uses a combinatorial disjoint union rule

$$c_{\wedge} c_{\vee} = 4c_{\wedge} + 2c_{\Delta} + 2c_{\lambda} + 4c_{\sqcap} + c_{\wedge}$$

Combinatorial Construction of Cumulants

$$\mu_1 = \kappa_1 \quad \langle X \rangle = \text{mean}$$

$$\mu_2 = \kappa_2 + \kappa_1 \kappa_1 \quad \langle X^2 \rangle = \text{variance} + \text{mean}^2$$

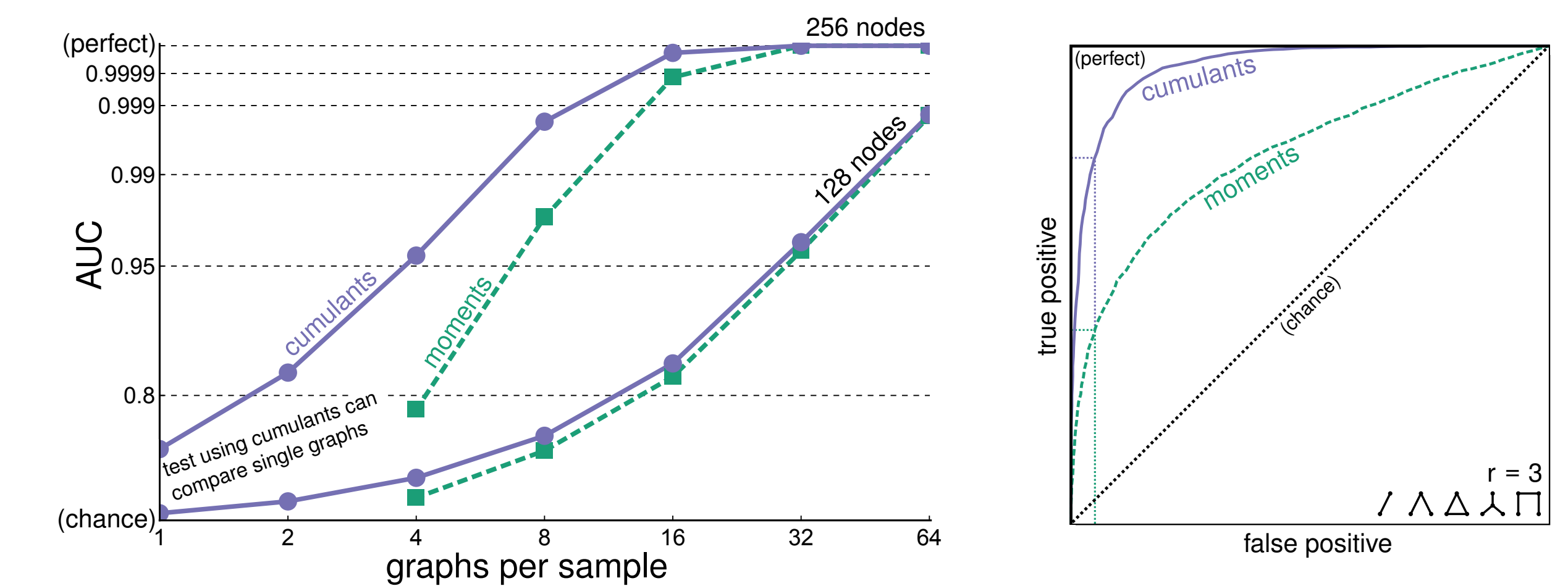
$$\mu_3 = \kappa_3 + \kappa_2 \kappa_1 + \kappa_2 \kappa_1 + \kappa_2 \kappa_1 + \kappa_1 \kappa_1 \kappa_1$$

$$\mu_{\sqcap} = \kappa_{\sqcap} + \kappa_{\wedge} \kappa_{\vee} + \kappa_{\wedge} \kappa_{\vee} + \kappa_{\wedge} \kappa_{\vee} + \kappa_{\vee} \kappa_{\vee} \kappa_{\vee}$$

Graph Cumulants Clearly Conquer

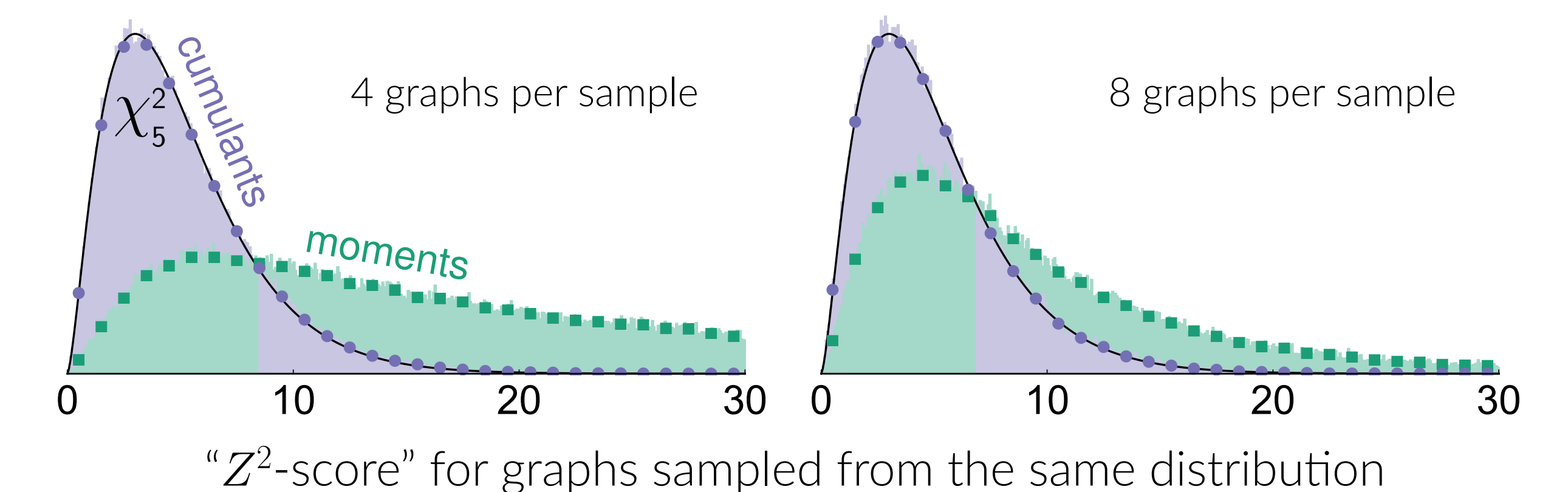
Graph cumulants outperform **subgraph densities** in general...

...and **graph cumulants** also work for single graph samples!

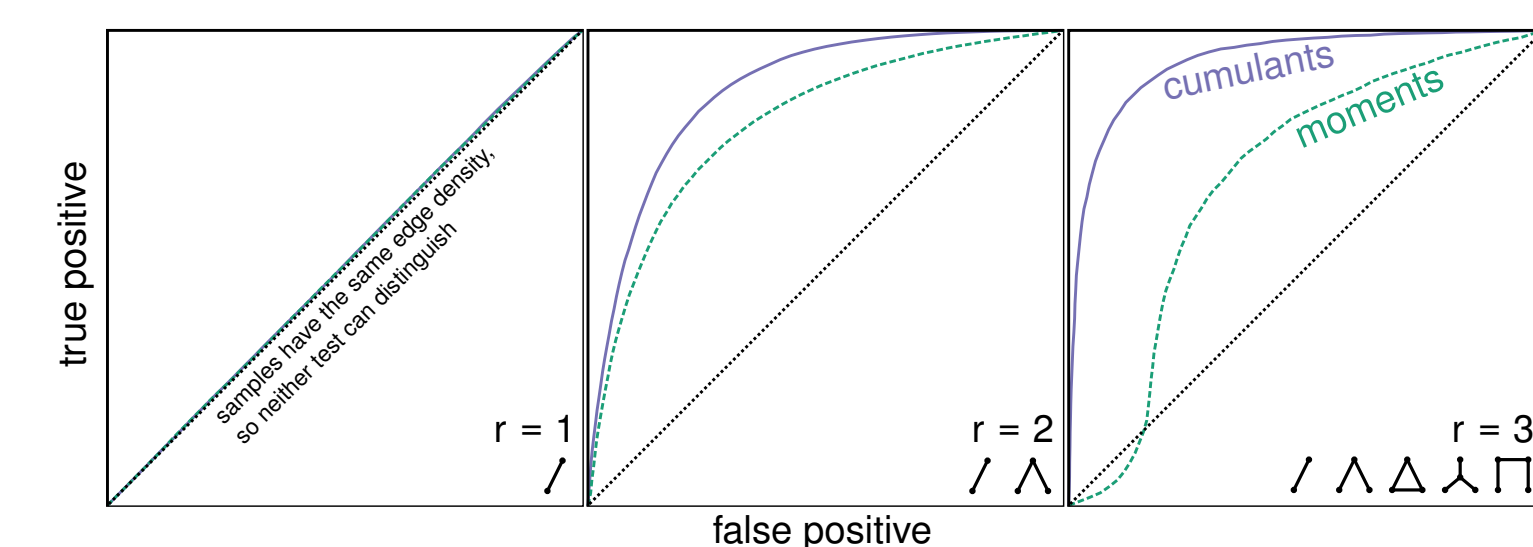


Why do **graph cumulants** perform better?

...because their **fluctuations** look more \mathcal{N} ormal!



Graph Cumulants in the (semi-)Wild!



Varying:
Number of subgraphs
used by both tests

Comparing:
Genetic interaction networks
of Arabidopsis and Mouse

Varying:
Number of graphs per sample
used by both tests

Comparing:
Genetic interaction networks
of Human and Rat

