

# Distributional Statistics of Random Set Models with Fixed Cardinalities

Alice K. Ng<sup>1,2</sup>, Jiahua Xu<sup>1,2</sup>, Paolo Tasca<sup>2</sup>

<sup>1</sup> University College London

<sup>2</sup> Exponential Science

**Abstract.** We develop a unified distributional framework for coverage and overlap in random set models with fixed-cardinality, without-replacement sampling, equivalently represented as random bipartite graphs with a prescribed degree sequence on one vertex set. On the exact side, we introduce a threshold-state dynamic program that tracks degree counts (the numbers of items covered by exactly  $d$  parties) and yields finite- $N$  distributions for exactly- $d$  and at-least- $t$  parties coverage. By aggregation, the same engine recovers the classical recursions for the union, the intersection, and their joint law. On the asymptotic side, we prove a multivariate Central Limit Theorem (CLT) for the full degree-count vector via Stein’s method with an exchangeable-pairs coupling; linear projections (including the union, the intersection, and their difference) are thus Gaussian in the limit, and smooth indices such as the Jaccard similarity admit delta-method approximations with closed-form moments. The exact evaluator supports moderate instances and provides ground truth for verification, while the Gaussian surrogates are accurate across broad regimes. In terms of cost, exact evaluation scales at worst between  $O(mN^2)$  and  $O(mN^3)$  depending on the target, whereas assembling the CLT requires only  $O(m^3)$ . We illustrate how these tools inform the design of probabilistic data structures (Bloom filters, MinHash) and cryptographic mechanisms such as threshold schemes in secret sharing.

**Keywords:** Combinatorics · Distribution Statistics · Probability · Random Sets · Bipartite Graphs · Probabilistic Data Structures

## 1 Introduction

Random set systems are a basic combinatorial model with appearances in probabilistic data structures (e.g., Bloom filters and MinHash), incidence/coverage processes (e.g., bipartite graphs and sensor coverage), and cryptography (e.g., additive/threshold secret sharing over index sets). Consider  $m$  subsets  $P_1, \dots, P_m$  drawn from a finite ground set  $[N] = \{1, \dots, N\}$ , where each  $P_r$  is chosen uniformly at random from  $\binom{[N]}{n_r}$ , independently across  $r$  (fixed-cardinality, without-replacement sampling). In many applications one needs not only expectations but the full distributions of

$$U = \left| \bigcup_{r=1}^m P_r \right| \quad (\text{total union}), \quad I = \left| \bigcap_{r=1}^m P_r \right| \quad (\text{total intersect}), \quad (1)$$

as well as functionals such as the Jaccard similarity  $J = I/U$ .

Alternatively, the system can be viewed as a bipartite graph  $G = (P, [N]; E)$ , with left vertices  $P = \{P_1, P_2, \dots, P_m\}$  and right vertices  $[N] = \{1, 2, \dots, N\}$ , where  $\{P_r, i\} \in E$  iff  $i \in P_r$ . For an item  $i \in [N]$ , its degree  $\deg(i)$  equals the number of sets  $P_r$  containing it. Hence the union  $U$  correspond to items of degree at least 1, and intersection  $I$  correspond to items of degree exactly  $m$ .

A convenient way to organise these statistics is via the *item degree counts*

$$\begin{aligned} Z_d &:= |\{i \in [N] : \deg(i) = d\}| \\ &= |\{i \in [N] : i \text{ belongs to exactly } d \text{ of } P_1, \dots, P_m\}|, \quad d = 0, 1, \dots, m. \end{aligned}$$

Then many quantities of interest are linear projections of the degree-count vector  $\mathbf{Z} = (Z_0, \dots, Z_m)$ : e.g.,  $U = \sum_{d \geq 1} Z_d$ ,  $I = Z_m$ ,  $U - I = \sum_{d=1}^{m-1} Z_d$ , and threshold counts  $W_t = \sum_{d \geq t} Z_d$ . This  $\mathbf{Z}$ -view unifies  $(U, I)$  and a broad family of threshold or aggregate counts, and it naturally supports Gaussian approximations via multivariate limit theorems.

We study the fixed-size model with arbitrary  $\{n_r\}$  and develop both exact and asymptotic tools around the same state representation.

#### *Contributions.*

- **Exact counting via a threshold state (section 4).** We introduce a Markovian threshold state that tracks bucket sizes by coverage degree up to a threshold  $\tau$ , yielding a forward dynamic program with matrix update  $\mathbf{S}_{k+1}^{(\tau)} = \mathbf{S}_k^{(\tau)} + M_\tau \mathbf{n}_{k+1}^{(\tau)}$ , and multivariate-hypergeometric weights. This produces exact finite- $N$  probability mass functions (PMFs) for threshold coverage  $W_t$  (take  $\tau = t$ ) and degree counts  $Z_d$  (take  $\tau = d+1$ ). We give tight one-step feasibility bounds for dynamic programming (DP) and show that, after aggregation, this engine *reduces* to the classical global recursions  $U, I$ , in Appendix A. These identities let us verify asymptotics directly against exact laws in section 4.
- **General multivariate CLT for degree counts (section 5).** We prove a multivariate CLT for  $(Z_0, \dots, Z_m)$  with  $O(N^{-1/2})$  accuracy. Any fixed linear projection  $(u^\top Z, v^\top Z, \dots)$ , including  $(U, I)$ ,  $U - I$ ,  $W_t$ , or multi-slice vectors, is thus asymptotically Gaussian with explicit means and covariances obtained by block-summing the covariance of  $Z$ .
- **Smooth indices via the delta method (section 6).** We give a delta-method Gaussian approximation (closed-form mean/variance) for  $J = I/U$  and, more generally, for any smooth functional  $g$  of a linear projection of  $Z$ . We also provide the exact finite- $N$  distribution of  $J$ , for any  $m \geq 2$ , via recursion function  $F_m$ .
- **Algorithms, validations, and use cases (section 7).** We implement a memoised evaluator for the threshold-state DP to calculate the exact PMFs, and compare it with the CLT/delta approximations across regimes, and illustrate design implications for Bloom filters under merges, MinHash planning, threshold secret-sharing, and incidence-graph coverage targets.

## 2 Related work

*Unions, intersections, and fixed-cardinality models.* Classical results give *univariate* laws and moments for unions (and, less often, intersections) under uniform without-replacement sampling, typically with equal subset sizes; see Barot-de la Peña [1] for unions and Kalinka [7] for intersections. Joint laws for union *and* intersection are not standard, and explicit treatments that allow arbitrary  $\{n_r\}$  are scarce. Our work supplies a finite- $N$  recursion for the joint law that handles general  $\{n_r\}$  and provides sharp feasibility bounds; the univariate laws follow by marginalization.

*Bernoulli presence models and Jaccard.* A substantial line of work analyzes Jaccard under i.i.d. Bernoulli presence/absence, yielding exact nulls and asymptotics [5, 8]. In the fixed-cardinality (hypergeometric) design relevant here, most exact results focus on  $m=2$  (e.g., Real-Vargas [11]). We extend to arbitrary  $m$  via the joint recursion and, in a unified *degree-count* view, provide Gaussian/delta approximations for Jaccard and other smooth indices.

*Stein’s method for multivariate normal approximation.* Exchangeable-pairs approaches for vectors were developed by Chatterjee-Meckes [4] and Reinert-Röllin [12]. We instantiate these techniques in the fixed-cardinality set model: a swap coupling yields a bivariate CLT for (union, intersection) and a multivariate CLT for the full degree-count vector, both with  $O(N^{-1/2})$  rates, from which linear projections inherit Gaussian limits.

*Random bipartite graphs with prescribed degrees.* Degree-constrained bipartite graphs and binary matrices with prescribed margins are widely studied, often focusing on counting/typical structure or on generating uniform samples via degree-preserving randomisation and swap chains [10, 14]. Related work also analyses degree sequences in random bipartite models, including cases where degrees are specified in one colour class [9]. By contrast, our focus is distributional: we give a computable finite- $N$  recursion for the full item-side degree histogram  $\mathbf{Z}$  (and hence union-intersection and threshold functionals), together with multivariate normal approximations that transfer to linear projections and smooth indices.

*Probabilistic data structures* Bloom filters and MinHash drive distributional questions about unions, intersections, and ratios; prior analyses often rely on independence or Poissonisation heuristics. Our exact recursion and degree-count CLTs offer finite- $N$  design rules for merged filters and similarity sketches under without-replacement sampling.

Relative to (i) univariate union/intersection results with equal sizes, (ii) Bernoulli-based Jaccard analyses, and (iii) degree-constrained bipartite-graph work that focuses on sampling or asymptotic graph properties, we provide: a computable finite- $N$  joint recursion for union and intersection with arbitrary

$\{n_r\}$ ; a unifying level-count framework with multivariate CLTs (and delta-method approximations for smooth indices); and a threshold-state dynamic program whose aggregation recovers the classical univariate and bivariate recursions used for verification.

### 3 Model and notation

Let  $[N] = \{1, 2, \dots, N\}$  be a finite universe and let  $m \in \mathbb{N}$  be fixed. For  $r \in [m]$ , party  $r$  samples a subset

$$P_r \subseteq [N], \quad |P_r| = n_r,$$

uniformly without replacement, independently across  $r$ . We allow  $n_r = n_r(N)$  to depend on  $N$ , and write

$$\alpha_r := \frac{n_r}{N} \in (0, 1), \quad r \in [m],$$

with  $\alpha_r$  fixed as  $N \rightarrow \infty$  (the “fixed-proportions” regime). For brevity, write  $S_m = (n_1, \dots, n_m)$  and  $\alpha = (\alpha_1, \dots, \alpha_m)$ .

*Degree-count vector and linear projections.* For each item  $i \in [N]$ , let the per-item coverage count be

$$R(i) := \sum_{r=1}^m \mathbf{1}\{i \in P_r\}.$$

For  $d = 0, 1, \dots, m$  define the degree counts

$$Z_d = |\{i \in [N] : R(i) = d\}|,$$

so  $\mathbf{Z} = (Z_0, \dots, Z_m)$  and  $\sum_{d=0}^m Z_d = N$ .

Many statistics of interest are linear projections of  $\mathbf{Z}$ :

$$U = \sum_{d \geq 1} Z_d, \quad I = Z_m, \quad W_t := \sum_{d \geq t} Z_d,$$

as well as other non-unanimous mass, e.g.,  $U - I = \sum_{d=1}^{m-1} Z_d$ ,  $A_D = \sum_{d \in D} Z_d$  ( $D \subseteq \{0, \dots, m\}$ ).

### 4 Exact counting via a threshold state

We study the exact laws of the degree-count statistics

$$Z_d := |\{i : R_m(i) = d\}|, \quad W_t := |\{i : R_m(i) \geq t\}| = \sum_{d \geq t} Z_d,$$

where for  $k \in [m]$  the running coverage count of item  $i$  is

$$R_k(i) = \sum_{r=1}^k \mathbf{1}\{i \in P_r\}, \quad R_0(i) := 0.$$

We process parties sequentially and keep a *threshold state* that makes the evolution Markov, so the next party's draw decomposes into independent multivariate-hypergeometric choices across buckets.

*Threshold state.* Fix an integer  $\tau \in [1, m]$ . After  $k$  parties, define

$$\mathbf{B}_k^{(\tau)} = (b_k^{(0)}, \dots, b_k^{(\tau-1)})^\top \in \mathbb{Z}_{\geq 0}^{\tau \times 1}, \quad \sum_{s=0}^{\tau-1} b_k^{(s)} \leq N,$$

with

$$b_k^{(s)} = |\{i : R_k(i) = s\}| \quad (0 \leq s \leq \tau-1), \quad b_k^{(\tau)} = |\{i : R_k(i) \geq \tau\}| = N - \sum_{s=0}^{\tau-1} b_k^{(s)}.$$

Thus the bucket  $b_k^{(\tau)}$  aggregates all items at degree  $\tau$  or higher. This single parameter  $\tau$  lets us read off different targets:

$$W_t = b_m^{(\tau)} \text{ when } \tau = t; \quad Z_d = b_m^{(d)} \text{ when } \tau = d + 1.$$

*Counts of ordered selections.* Let  $H_k(\mathbf{B}_k^{(\tau)})$  denote the number of ordered  $k$ -tuples  $(P_1, \dots, P_k)$  that realise the state  $\mathbf{B}_k^{(\tau)}$ . The base case is deterministic:

$$\mathbf{B}_1^{(\tau)} = (N - n_1, n_1, 0, \dots, 0)^\top, \quad H_1(\mathbf{B}_1^{(\tau)}) = \binom{N}{n_1},$$

and  $H_1(\cdot) = 0$  for any other state.

*One-step update.* Given  $\mathbf{B}_k^{(\tau)}$  and adding party  $k+1$  with  $|P_{k+1}| = n_{k+1}$ , split that party's draw by buckets as

$$\mathbf{n}_{k+1}^{(\tau)} = (n_{k+1}^{(0)}, \dots, n_{k+1}^{(\tau-1)})^\top \in \mathbb{Z}_{\geq 0}^{\tau \times 1}, \quad n_{k+1}^{(\tau)} := n_{k+1} - \sum_{s=0}^{\tau-1} n_{k+1}^{(s)}, \quad 0 \leq n_{k+1}^{(s)} \leq b_k^{(s)} \quad \forall s \in [0, \tau].$$

Elements drawn from bucket  $s < \tau$  move up exactly one degree; elements drawn from bucket  $\tau$  stay in bucket  $\tau$ . Hence

$$b_{k+1}^{(0)} = b_k^{(0)} - n_{k+1}^{(0)}, \tag{2a}$$

$$b_{k+1}^{(s)} = b_k^{(s)} - n_{k+1}^{(s)} + n_{k+1}^{(s-1)} \quad (1 \leq s \leq \tau-1), \tag{2b}$$

$$b_{k+1}^{(\tau)} = b_k^{(\tau)} + n_{k+1}^{(\tau-1)}. \tag{2c}$$

We can write

$$\mathbf{B}_{k+1}^{(\tau)} = \mathbf{B}_k^{(\tau)} + M_\tau \mathbf{n}_{k+1}^{(\tau)} \quad (3)$$

where  $M_\tau \in \mathbb{Z}^{\tau \times \tau}$ ,  $M_{j,j} = -1$  for  $j \in [1, \tau]$ , and  $M_{j,j-1} = +1$  for  $j \in [2, \tau]$ , i.e.

$$M_\tau = \begin{pmatrix} -1 & 0 & \cdots & 0 \\ 1 & -1 & \ddots & \vdots \\ 0 & \ddots & \ddots & 0 \\ \vdots & \cdots & 1 & -1 \end{pmatrix} \in \mathbb{Z}^{\tau \times \tau}. \quad (4)$$

Since  $M_\tau$  is lower triangular with diagonal  $-1$ , it is invertible and the draw split is uniquely recovered from the state difference:

$$\mathbf{n}_{k+1}^{(\tau)} = M_\tau^{-1} \left( \mathbf{B}_{k+1}^{(\tau)} - \mathbf{B}_k^{(\tau)} \right). \quad (5)$$

Componentwise, with  $\Delta b^{(s)} := b_{k+1}^{(s)} - b_k^{(s)}$ ,

$$\Delta b^{(0)} = -n_{k+1}^{(0)}, \quad \Delta b^{(s)} = -n_{k+1}^{(s)} + n_{k+1}^{(s-1)} \quad (1 \leq s \leq t-1).$$

Since  $M_\tau^{-1}$  is the lower triangle filled with  $-1$ ,

$$n_{k+1}^{(s)} = -\sum_{r=0}^s \Delta b^{(r)} \quad (0 \leq s \leq t-1), \quad n_{k+1}^{(t)} = n_{k+1} - \sum_{s=0}^{t-1} n_{k+1}^{(s)}.$$

Given a feasible split, the number of ways to realise it is the multivariate-hypergeometric factor

$$\prod_{s=0}^{\tau} \binom{b_k^{(s)}}{n_{k+1}^{(s)}}.$$

Therefore

$$H_{k+1}(\mathbf{B}_{k+1}^{(\tau)}) = \sum_{\mathbf{B}_k^{(\tau)}} \left( \prod_{s=0}^{\tau} \binom{b_k^{(s)}}{n_{k+1}^{(s)}} \right) H_k(\mathbf{B}_k^{(\tau)}), \quad (6)$$

summing over all the feasible predecessors  $\mathbf{B}_k^{(\tau)}$  consistent with boundaries (7a-e).

*Feasible bounds.* A transition  $\mathbf{B}_k^{(\tau)} \rightarrow \mathbf{B}_{k+1}^{(\tau)}$  is feasible iff

$$0 \leq n_{k+1}^{(s)} \leq b_k^{(s)} \quad (0 \leq s \leq \tau), \quad (7a)$$

$$\max\{0, b_k^{(0)} - n_{k+1}\} \leq b_{k+1}^{(0)} \leq b_k^{(0)}, \quad (7b)$$

$$\max\{0, b_k^{(s)} - n_{k+1}\} \leq b_{k+1}^{(s)} \leq b_k^{(s)} + \min\{b_k^{(s-1)}, n_{k+1}\} \quad (1 \leq s \leq \tau-1), \quad (7c)$$

$$b_k^{(\tau)} \leq b_{k+1}^{(\tau)} \leq b_k^{(\tau)} + \min\{b_k^{(\tau-1)}, n_{k+1}\}, \quad (7d)$$

$$\sum_{r=0}^s b_k^{(r)} - \min\{b_k^{(s)}, n_{k+1}\} \leq \sum_{r=0}^s b_{k+1}^{(r)} \leq \sum_{r=0}^s b_k^{(r)} \quad (0 \leq s \leq \tau-1). \quad (7e)$$

**Theorem 1 (PMF for threshold coverage  $W_t$ ).** Fix  $t \geq 1$  and take  $\tau = t$ . Then  $W_t = b_m^{(\tau)}$ , and

$$\mathbb{P}(W_t = w) = \sum_{\mathbf{B}_m^{(t)}: b_m^{(t)} = w} \frac{H_m(\mathbf{B}_m^{(t)})}{\prod_{r=1}^m \binom{N}{n_r}}.$$

**Theorem 2 (PMF for the degree count  $Z_d$ ).** Fix  $d \geq 0$  and take  $\tau = d + 1$ . Then  $Z_d = b_m^{(d)}$ , and

$$\mathbb{P}(Z_d = z) = \sum_{\mathbf{B}_m^{(d+1)}: b_m^{(d)} = z} \frac{H_m(\mathbf{B}_m^{(d+1)})}{\prod_{r=1}^m \binom{N}{n_r}}.$$

*Remark.* For fixed  $\tau$ , complexity is pseudo-polynomial in the size of the feasible grid of  $\mathbf{B}_k^{(\tau)}$  and practical for small  $m, N$ . Aggregating the state yields 1D/2D recursions for specific targets (below), with worst-case time  $T(m, N)$  between  $O(mN^2)$  and  $O(mN^3)$ .

**Lemma 1 (Reduction to the union ( $U = \sum_{l \geq 1}^m Z_l$ )).** Set  $\tau = 1$  so  $\mathbf{B}_k^{(1)} = (b_k^{(0)})$  and  $U_k := N - b_k^{(0)} = |\cup_{r \leq k} P_r|$ . Define aggregated counts from the threshold state

$$C_k(u_k) := \sum_{\mathbf{B}_k^{(1)}: N - b_k^{(0)} = u_k} H_k(\mathbf{B}_k^{(1)}).$$

Then  $C_1(u_1) = \binom{N}{n_1} \mathbf{1}\{u_1 = n_1\}$  and, for  $k \geq 1$ ,

$$C_{k+1}(u_{k+1}) = \sum_{u_k} \binom{u_k}{n_{k+1} + u_k - u_{k+1}} \binom{N - u_k}{u_{k+1} - u_k} C_k(u_k),$$

Hence  $\mathbb{P}(U_m = u_m) = C_m(u_m) / \prod_{r=1}^m \binom{N}{n_r}$ .

*Proof sketch.* With  $\tau = 1$  the split is  $(n^{(0)}, n^{(1)})$  between “new” (from  $N - U_k$ ) and “already covered” (from  $U_k$ ), and  $U_{k+1} = U_k + n^{(0)}$  forces  $n^{(0)} = u_{k+1} - u_k$ . The multiplicity of a step is  $\binom{U_k}{n_{k+1} - n^{(0)}} \binom{N - U_k}{n^{(0)}} = \binom{u_k}{n_{k+1} + u_k - u_{k+1}} \binom{N - u_k}{u_{k+1} - u_k}$ , summing over feasible  $x$  gives the stated recursion. Detailed bounds for  $u_k$  and the full proof are in Appendix A.1. Note that the time complexity for the calculation is  $T(m, N) = O(mN^2)$ .  $\square$

**Lemma 2 (Reduction to the intersection ( $I = Z_m$ )).** Run the threshold scheme with the time-varying choice  $\tau = k$  at step  $k$ , so  $I_k := b_k^{(k)} = |\cap_{r \leq k} P_r|$ . Define aggregated counts

$$D_k(v_k) := \sum_{\mathbf{B}_k^{(k)}: b_k^{(k)} = v_k} H_k(\mathbf{B}_k^{(k)}).$$

Then  $D_1(y; S_1) = \binom{N}{n_1} \mathbf{1}\{y = n_1\}$  and, for  $k \geq 1$ ,

$$D_{k+1}(v_{k+1}) = \sum_{v_k} \binom{v_k}{v_{k+1}} \binom{N - v_k}{n_{k+1} - v_{k+1}} D_k(v_k).$$

Hence  $\mathbb{P}(I_m = v_m) = D_m(v_m) / \prod_{r=1}^m \binom{N}{n_r}$ .

Proof. At step  $k$ , items with coverage  $\geq k$  are exactly the current intersection, of size  $v_k$ . To have  $I_{k+1} = v_{k+1}$ , party  $k+1$  must pick exactly  $v_{k+1}$  from those  $v_k$  and the remaining  $n_{k+1} - v_{k+1}$  from outside, giving  $\binom{v_k}{v_{k+1}} \binom{N - v_k}{n_{k+1} - v_{k+1}}$ , summing over feasible  $v_k$  yields the recursion. Detailed bounds for  $u_k$  and the full proof are in Appendix A.2. Same as union, the time complexity for the calculation is  $T(m, N) = O(mN^2)$ .  $\square$

**Lemma 3 (Reduction to the bivariate recursion  $F_m$ ).** Fix  $\tau = m$  and define aggregated counts from the threshold state by

$$F_k(u_k, v_k) := \sum_{\mathbf{B}_k^{(m)}: N - b_k^{(0)} = u_k, \sum_{s=k}^{m-1} b_k^{(s)} + b_k^{(\geq m)} = v_k} H_k(\mathbf{B}_k^{(m)}).$$

Then  $F_1(u_1, v_1) = \binom{N}{n_1} \mathbf{1}\{u_1 = v_1 = n_1\}$ , and for all  $m \geq 2$ ,

$$F_m(u_m, v_m) = \sum_{v_{m-1}} \sum_{u_{m-1}} \binom{v_{m-1}}{v_m} \binom{u_{m-1} - v_{m-1}}{n_m + u_{m-1} - u_m - v_m} \binom{N - u_{m-1}}{u_m - u_{m-1}} F_{m-1}(u_{m-1}, v_{m-1}),$$

Hence  $p_{X,Y}(u_m, v_m) = F_m(u_m, v_m; S_m) / \prod_{r=1}^m \binom{N}{n_r}$  as in (26).

The feasibility bounds on  $(u_{m-1}, v_{m-1})$  and the full proof are in Appendix A.3. Note that the time complexity for the calculation is  $T(m, N) = O(mN^3)$ .

**Proposition 1 (Jaccard index using bivariate recursion  $F_m$ ).** Let  $a \geq 1$ ,  $0 \leq b \leq a$  be coprime. Then for any  $m \geq 2$

$$\mathbb{P}\left(J = \frac{b}{a}\right) = \sum_{k \in \mathcal{K}(a,b)} \frac{F_m(ka, kb)}{\prod_{r=1}^m \binom{N}{n_r}},$$

where  $\mathcal{K}(a, b) = \{k \in \mathbb{Z}_{\geq 0} : (ka, kb) \text{ is feasible}\}$ .

## 5 General multivariate CLT

This section provides the Gaussian approximation for the degree-count vector  $\mathbf{Z} = (Z_0, \dots, Z_m)$ . All normal approximations for univariate and multivariate projections (e.g.  $U$ ,  $I$ ,  $(U, I)$ ,  $W_t$ , or any  $(A, B)$ ) then follow as corollaries by block-summing the covariance.



### 5.1 Means and one-item probabilities

For a single item  $i$ , independence across parties gives, for  $d = 0, \dots, m$ ,

$$p_d(\alpha) := \mathbb{P}\{R(i) = d\} = [t^d] \prod_{r=1}^m ((1 - \alpha_r) + \alpha_r t). \quad (8)$$

Here  $[t^d]$  means “coefficient of  $t^d$ ” in a probability generating function (PGF). Identity (8) is the *coefficient-of* form of the elementary symmetric polynomials  $e_d$ :

$$\mathbb{E}(t^{R(i)}) = \prod_{r=1}^m ((1 - \alpha_r) + \alpha_r t) = \sum_{d=0}^m \left( \prod_{r=1}^m (1 - \alpha_r) \right) e_d \left( \frac{\alpha_1}{1 - \alpha_1}, \dots, \frac{\alpha_m}{1 - \alpha_m} \right) t^d.$$

Hence  $p_d$  is explicit and depends on  $\alpha$  only, and computing all  $p_d$  via a rolling 2-term convolution takes  $T(m) = O(m^2)$ .

By exchangeability of items,

$$\mathbb{E}[Z_d] = \sum_{i=1}^N \mathbb{P}\{R(i) = d\} = N p_d(\alpha). \quad (9)$$

### 5.2 Finite- $N$ covariances: exact identities

Let  $I_i^{(d)} := \mathbf{1}\{R(i) = d\}$ , so  $Z_d = \sum_{i=1}^N I_i^{(d)}$ . For distinct items  $u \neq v$ , write

$$q_{a,b} := \mathbb{P}(R(u) = a, R(v) = b), \quad a, b \in \{0, \dots, m\}.$$

Then by expanding variances and covariances of sums of indicators:

$$\text{Var}(Z_a) = \sum_{i=1}^N \text{Var}(I_i^{(a)}) + \sum_{u \neq v} \text{Cov}(I_u^{(a)}, I_v^{(a)}) = N p_a(1 - p_a) + N(N - 1) (q_{a,a} - p_a^2), \quad (10)$$

$$\begin{aligned} \text{Cov}(Z_a, Z_b) &= \sum_{u \neq v} \text{Cov}(I_u^{(a)}, I_v^{(b)}) - N p_a p_b \\ &= N(N - 1) q_{a,b} - N^2 p_a p_b, \quad a \neq b. \end{aligned} \quad (11)$$

For a *fixed* party  $r$ , the pair  $(\mathbf{1}\{u \in P_r\}, \mathbf{1}\{v \in P_r\})$  has the *hypergeometric* probabilities

$$\pi_{00}^{(r)} = \frac{(N - n_r)(N - n_r - 1)}{N(N - 1)}, \quad \pi_{10}^{(r)} = \pi_{01}^{(r)} = \frac{n_r(N - n_r)}{N(N - 1)}, \quad \pi_{11}^{(r)} = \frac{n_r(n_r - 1)}{N(N - 1)}.$$

Define the per-party bivariate PGF

$$\phi_r(z, w) := \pi_{00}^{(r)} + \pi_{10}^{(r)} z + \pi_{01}^{(r)} w + \pi_{11}^{(r)} zw.$$

Independence across parties yields

$$q_{a,b} = [z^a w^b] \prod_{r=1}^m \phi_r(z, w). \quad (12)$$

*Cost.* All  $q_{a,b}$  can be obtained by  $m$  successive  $2 \times 2$  convolutions on an  $(\leq m+1) \times (\leq m+1)$  table, with time complexity  $T(m) = O(m^3)$ . Plugging into (10)–(11) forms  $\Sigma_N = \text{Cov}(\mathbf{Z})$  in the same  $T(m) = O(m^3)$ .

### 5.3 Asymptotic covariance scaling

Expanding each  $\phi_r$  to first order in  $(N-1)^{-1}$  gives

$$\phi_r(z, w) = ((1 - \alpha_r) + \alpha_r z)((1 - \alpha_r) + \alpha_r w) + \frac{\alpha_r(1 - \alpha_r)}{N-1} (-1 + z + w - zw).$$

Multiplying over  $r$  and extracting coefficients as in (12) yields

$$q_{a,b} = p_a p_b + \frac{1}{N-1} S_{a,b}(\alpha) + O((N-1)^{-2}),$$

where  $S_{a,b}(\alpha)$  is an  $O(1)$  quantity depending only on  $\alpha$ . Substituting into (10)–(11) gives

$$\Sigma_N = N \Sigma(\alpha) + O(1), \quad (13)$$

with entries

$$\Sigma_{aa}(\alpha) = p_a(1 - p_a) - \sum_{r=1}^m \alpha_r(1 - \alpha_r) (\Delta_a^{(-r)})^2, \quad (14)$$

$$\Sigma_{ab}(\alpha) = -p_a p_b - \sum_{r=1}^m \alpha_r(1 - \alpha_r) \Delta_a^{(-r)} \Delta_b^{(-r)}, \quad a \neq b. \quad (15)$$

Here  $p_d^{(-r)} := [t^d] \prod_{s \neq r} ((1 - \alpha_s) + \alpha_s t)$ ,  $\Delta_d^{(-r)} := p_d^{(-r)} - p_{d-1}^{(-r)}$ , and  $p_{-1}^{(-r)} := 0$ . All  $p_k^{(-r)}$  for  $r = 1..m$ ,  $d = 0..m$  can be computed in  $O(m^2)$  time; assembling  $\Sigma(\alpha)$  via (14)–(15) requires  $O(m^3)$  time (due to  $O(m^2)$  pairs  $(a, b)$  each with an  $O(m)$  inner sum).

### 5.4 A multivariate CLT for the degree counts

Because  $\sum_{l=0}^m Z_l = N$  is deterministic,  $\Sigma_N$  is singular in the direction of the all-ones vector  $\mathbf{1} = (1, \dots, 1)^\top$ . Let  $\mathcal{T} := \{x \in \mathbb{R}^{m+1} : \mathbf{1}^\top x = 0\}$  be the  $m$ -dimensional *tangent subspace* and let  $\Pi$  be the orthogonal projector onto  $\mathcal{T}$ . Define the centered, scaled vector

$$W_N := \frac{1}{\sqrt{N}} (\Pi \Sigma(\alpha) \Pi^\top)^{-1/2} \Pi (Z - \mathbb{E}(Z)).$$

**Theorem 3 (Multivariate CLT on the simplex tangent).** *Fix  $m$  and  $\alpha_r \in (0, 1)$ . Then  $W_N \rightarrow \mathcal{N}(0, I_m)$  as  $N \rightarrow \infty$ .*

*$\lim_{N \rightarrow \infty} W_N \sim \mathcal{N}(0, I_m)$ .*

*Moreover, for smooth test functions the approximation error is  $O(N^{-1/2})$ .*

*Proof (sketch).* We build an exchangeable pair  $(\mathbf{Z}, \mathbf{Z}')$  by a swap that only modifies  $O(1)$  items while preserving all set sizes  $|P_r| = n_r$ : pick a random item  $K$ ; for each party  $r$ , if  $K \in P_r$  swap it with a uniformly chosen element outside  $P_r$ , and if  $K \notin P_r$  swap it with a uniformly chosen element inside  $P_r$ . This creates small increments  $\Delta := \mathbf{Z}' - \mathbf{Z}$  with bounded third moments and gives the regression

$$\mathbb{E}(\Delta \mid \mathbf{Z}) = -\frac{1}{N} \Pi (\mathbf{Z} - \mathbb{E}(\mathbf{Z})),$$

where  $\Pi$  is the projector onto  $\mathbb{T}$ . Conditional second moments match  $\Sigma(\alpha)$  up to  $O(1/N)$ . The multivariate Stein–exchangeable-pairs theorem then yields the CLT with  $O(N^{-1/2})$  error. Full details in Appendix C.

## 6 Linear projections and smooth indices

This section records the principal projections of the level counts  $\mathbf{Z} = (Z_0, \dots, Z_m)$  and their Gaussian/delta approximations. All limits are immediate from Theorem 3. Exact finite- $N$  pmfs come from the recursions in Appendix A; figures are in Appendix D.

**Corollary 1 (CLTs for linear projections).** *Let  $L$  be any fixed  $\ell \times (m+1)$  matrix and set  $X^{(\ell)} := LZ$ . Then  $X^{(\ell)}$  inherits a joint Gaussian limit. Under Theorem 3,*

$$\frac{X^{(\ell)} - \mathbb{E}(X)^{(\ell)}}{\sqrt{N}} \rightarrow \mathcal{N}_\ell(0, L\Sigma(\alpha)L^\top).$$

*In particular,  $U = \sum_{k \geq 1} Z_k$  and  $I = Z_m$  are asymptotically normal, and  $(U, I)$  is asymptotically bivariate normal with covariance  $L\Sigma(\alpha)L^\top$  for the selector  $L$  that picks  $(\sum_{k \geq 1} Z_k, Z_m)$ .*

Once  $\Sigma(\alpha)$  is precomputed in  $O(m^3)$  time (Section 5.3), forming  $L\Sigma(\alpha)L^\top$  costs  $T(m) = O(\ell m^2)$ .

**Lemma 4 (Delta–method for smooth indices).** *Let  $g : \mathbb{R}^\ell \rightarrow \mathbb{R}$  be  $C^2$  near  $\mathbb{E}(X)^{(\ell)}$  and let  $T = g(X^{(\ell)})$ . Then*

$$\mathbb{E}(T) = g(\mu) + \frac{1}{2} \text{tr}(H_g(\mu) \Sigma_L) + O(N^{-1/2}), \quad \text{Var}(T) = \nabla g(\mu)^\top \Sigma_L \nabla g(\mu) + O(N^{-1/2}),$$

*where  $\mu = \mathbb{E}(X)^{(\ell)}$  and  $\Sigma_L = L\Sigma(\alpha)L^\top$ , and  $(T - \mathbb{E}(T))/\sqrt{\text{Var}(T)} \Rightarrow \mathcal{N}(0, 1)$ .*

*Cost.* Given  $\Sigma_L$ , evaluating the mean/variance expansions above is  $O(\ell^2)$

### 6.1 Univariate $U$ and $I$

*Means.* By (9) and the identities in §5.1,

$$\mu_U := \mathbb{E}(U) = N \left( 1 - \prod_{r=1}^m (1 - \alpha_r) \right), \quad \mu_I := \mathbb{E}(I) = N \prod_{r=1}^m \alpha_r.$$

*Variances (from  $\Sigma(\alpha)$ ).* Let  $L_U \in \mathbb{R}^{1 \times (m+1)}$  be  $L_U = (0, 1, \dots, 1)$  and  $L_I = (0, \dots, 0, 1)$ . Then, from (13)–(15),

$$\text{Var}(U) = N v_U(\alpha) + O(1), \quad v_U(\alpha) = L_U \Sigma(\alpha) L_U^\top, \quad (16)$$

$$\text{Var}(I) = N v_I(\alpha) + O(1), \quad v_I(\alpha) = L_I \Sigma(\alpha) L_I^\top. \quad (17)$$

**Corollary 2 (Univariate CLTs as projections of Theorem 3).** *With  $U = \sum_{\ell \geq 1} Z_\ell$  and  $I = Z_m$ ,*

$$\frac{U - \mu_U}{\sqrt{N}} \rightarrow \mathcal{N}(0, v_U(\alpha)), \quad \frac{I - \mu_I}{\sqrt{N}} \rightarrow \mathcal{N}(0, v_I(\alpha)),$$

*with an  $O(N^{-1/2})$  error for smooth test functions.*

*Cost.* After  $\Sigma(\alpha)$ , obtaining  $v_U, v_I$  is  $O(m^2)$ ; computing  $\mu_U, \mu_I$  is  $O(m)$ .

### 6.2 General case

Consider any univariate projection of the degree vector defined by an integer subset sum  $A_D := \sum_{d \in D} Z_d$  for  $D \subseteq \{0, \dots, m\}$ . This formulation encompasses key statistics such as the union size ( $D = \{1, \dots, m\}$ ), the intersection count ( $D = \{m\}$ ), or the count of items exceeding a threshold  $t$  ( $D = \{t, \dots, m\}$ ).

While Theorem 3 implies a Gaussian limit for  $A_D$ , the continuous approximation can exhibit discretisation errors in finite samples, especially at edge cases. To obtain a robust discrete refinement, we utilise the property that the fixed-cardinality constraints consistently induce under-dispersion ( $\text{Var}(A_D) < \mathbb{E}[A_D]$ ). This structural feature allows us to approximate  $A_D$  using a **Moment-Matched Binomial** distribution  $B(n_{\text{eff}}, p_{\text{eff}})$ .

First, compute the exact aggregate moments by summing the components of the mean vector and covariance matrix derived in Section 5.2:

$$\mu_D = \sum_{d \in D} \mathbb{E}[Z_d], \quad \sigma_D^2 = \sum_{a \in D} \sum_{b \in D} \text{Cov}(Z_a, Z_b),$$

where the individual terms  $\mathbb{E}[Z_d]$  and  $\text{Cov}(Z_a, Z_b)$  are given by (9) and (10)–(11), respectively.

Then determine the effective binomial parameters by matching these first two moments:

$$p_{\text{eff}} = 1 - \frac{\sigma_D^2}{\mu_D}, \quad n_{\text{eff}} = \frac{\mu_D}{p_{\text{eff}}}. \quad (18)$$

This discrete surrogate preserves the exact moments of the aggregate count  $A_D$  (e.g., the union size  $U$  or intersection count  $Z_m$ ) while respecting the non-negative integer support. By the De Moivre–Laplace theorem, this approximation is asymptotically consistent with the Gaussian limit of  $A_D$  implied by Theorem 3 as  $N \rightarrow \infty$ .

*Note on Multinomial Generalisation.* While this technique accurately models any single univariate projection, generalising to the full vector  $\mathbf{Z}$  via a standard Multinomial distribution is often infeasible.

A standard Multinomial distribution imposes a rigid covariance structure ( $\Sigma_{ab} \propto -\mathbb{E}[Z_a]\mathbb{E}[Z_b]$ ) that cannot generally capture the complex, parameter-dependent covariance matrix  $\Sigma(\alpha)$  derived in (15). Consequently, we recommend the Moment-Matched Binomial for accurate marginal inference of specific counts, while retaining the Multivariate Normal (Theorem 3) for analysing joint dependencies.

### 6.3 Example: Bivariate $(U, I)$

*Mean vector and covariance (as a projection).* Let  $L \in \mathbb{R}^{2 \times (m+1)}$  have first row  $L_U$  and second row  $L_I$ . Then

$$E \begin{bmatrix} U \\ I \end{bmatrix} = \begin{bmatrix} \mu_U \\ \mu_I \end{bmatrix}, \quad \text{Cov} \begin{bmatrix} U \\ I \end{bmatrix} = N \Sigma_{UI} + O(1), \quad \Sigma_{UI} := L \Sigma(\alpha) L^\top.$$

(Explicit formulas follow by substituting (14)–(15).)

**Corollary 3 (Bivariate CLT for  $(U, I)$ ).** *Under Theorem 3,*

$$\frac{1}{\sqrt{N}} \begin{pmatrix} U - \mu_U \\ I - \mu_I \end{pmatrix} \rightarrow \mathcal{N}_2(0, \Sigma_{UI}),$$

*with an  $O(N^{-1/2})$  error for smooth test functions (and hence in convex-set distance).*

*Cost.* Forming  $\Sigma_{UI}$  is  $O(m^2)$  once  $\Sigma(\alpha)$  is available.

### 6.4 Example: Jaccard indices from $(U, I)$

The Jaccard index  $J = U/I \in [0, 1]$  is a smooth function of  $(U, I)$  on the set  $\{u > 0\}$ . In our model  $\mu_U > 0$ , so the delta method applies.

**Corollary 4 (Delta-method approximation for  $J$ ).** *Let  $g(u, v) = v/u$  and write  $\mu = (\mu_U, \mu_I)$ ,  $\Sigma_{UI}$  as above. Then*

$$\frac{J - \mu_J}{\sigma_J} \Rightarrow \mathcal{N}(0, 1), \quad \text{with error } O(N^{-1/2})$$

for smooth test functions, where

$$\mu_J = g(\mu) + \frac{1}{2} \text{tr}(H_g(\mu) \Sigma_{UI}) + O(N^{-1/2}) = \frac{\mu_I}{\mu_U} + \frac{\mu_I \text{Var}(U) - \mu_U \text{Cov}(U, I)}{\mu_U^3} + O(N^{-1/2}), \quad (19)$$

$$\sigma_J^2 = \nabla g(\mu)^\top \Sigma_{UI} \nabla g(\mu) + O(N^{-1/2}) = \frac{\text{Var}(I)}{\mu_U^2} + \frac{\mu_I^2}{\mu_U^4} \text{Var}(U) - \frac{2\mu_I}{\mu_U^3} \text{Cov}(U, I) + O(N^{-1/2}). \quad (20)$$

Here  $\nabla g(\mu) = (-\mu_I/\mu_U^2, 1/\mu_U)$  and  $H_g(\mu) = \begin{bmatrix} 2\mu_I/\mu_U^3 & -1/\mu_U^2 \\ -1/\mu_U^2 & 0 \end{bmatrix}$ .

*Proof (sketch).* By Corollary 3,  $(U, I)$  is asymptotically normal with mean  $\mu$  and covariance  $N \Sigma_{UI} + O(1)$ . Apply the multivariate delta method to  $g(x = u, v) = v/u$ , expanding mean and variance to first order in  $1/N$  to obtain (19)–(20).

*Cost.* After  $\Sigma_{UI}$ , evaluating (19)–(20) is  $O(1)$ .

## 7 Applications

We show how the exact counts and the moment/CLT results for  $(U, I)$  and  $J$  translate into practical procedures. Means/variances/covariances come from section 6. When exact probabilities are feasible we use the dynamic programs of section 4; otherwise we use the (uni/bivariate) normal surrogates with standard continuity corrections.

### 7.1 Probabilistic data structures

Probabilistic data structures (PDS) are compact, mergeable sketches (e.g., Bloom filters, MinHash). Our laws for  $(U, I)$  and  $J$  yield finite- $N$  design rules.

**Bloom filters: sizing under merges** Consider  $m$  parties that insert  $|P_r| = n_r$  distinct elements into a common Bloom filter [3] of length  $M$  with  $h$  hash functions (no deletions). The false-positive rate (FPR), conditional on the *union* size  $U$ , is well-approximated by

$$\text{FPR}(U) \approx \left(1 - e^{-hU/M}\right)^h.$$

Hence, for a target FPR:  $\varepsilon$  and reliability  $1 - \delta$ , choose  $(M, h)$  so that

$$\text{P}(\text{FPR}(X) \leq \varepsilon) = \text{P}\left(X \leq \frac{M}{h} \log \frac{1}{1 - \varepsilon^{1/h}}\right) \geq 1 - \delta. \quad (21)$$

Two evaluation routes:

- *Exact*: compute the RHS by summing the counts for union  $C_m: \sum_{u \leq u^*} \frac{C_m(u)}{\prod_r \binom{N}{n_r}}$ , with  $u^*$  the threshold inside the braces. Complexity  $T(m, N) = O(mN^2)$
- *Gaussian*: use  $U \approx \mathcal{N}(\mu_U, \sigma_U^2)$  and  $\mathbb{P}(U \leq u^*) \approx \Phi((u^* + 0.5 - \mu_U)/\sigma_U)$ . Complexity  $T(m) = O(m^3)$  to precompute  $\Sigma(\alpha)$ .

For merge-heavy workloads, you can invert (21) to solve for the minimal  $M$  given  $h$  (or vice versa) at reliability  $1 - \delta$ .

**MinHash: sample size planning with a  $J$  prior** A MinHash sketch [6] with  $T$  independent hash functions yields  $\hat{J} = \frac{1}{T} \sum_{t=1}^T B_t$  with  $B_t \sim \text{Bernoulli}(J)$  conditionally on  $J$ . Our exact/approximate laws for  $J$  supply a *prior* (or design distribution) for planning  $T$ .

- *Frequentist sizing*: to guarantee a margin  $\eta$  at confidence  $1 - \delta$  for a nominal  $J_0$ , take  $T \geq \left( \frac{z_{1-\delta/2}}{\eta} \right)^2 J_0(1 - J_0)$ . Using  $J_0 = \mu_J$  from (19) is a principled default.
- *Bayesian credible intervals*: treat the law of  $J$  (exact for  $m = 2$ , delta-normal otherwise) as a prior to get posterior bands for  $J$  given  $\hat{J}$ ; this absorbs the finite- $N$  overlap structure among parties.

## 7.2 Secret sharing and access structures

In secret sharing protocols [2, 13], a dealer splits a secret  $D$  into  $N$  labelled *shares* and distributes them to parties. Operationally, multiple shares may be assigned to a single holder (“sharding”). In our model, party  $r$  receives  $n_r$  distinct shares uniformly without replacement, so different parties can receive overlapping shares.

In a threshold  $(k, N)$  scheme, recovery depends on the number of *unique* shares held by a group: even if the group holds more than  $k$  shares in total, overlaps can leave fewer than  $k$  distinct shares, making reconstruction impossible. Thus shares should be spread broadly so that any  $k$  independent parties can recover the secret; excessive overlap undermines both availability and security.

*Threshold  $(k, N)$  schemes.* Recovery occurs iff the union of allocated shares across the group reaches size  $k$ , i.e.  $U = |\cup_r P_r| \geq k$ . Evaluate

$$\mathbb{P}(U \geq k) = \sum_{u=k}^N \sum_v \frac{F_m(u, v)}{\prod_r \binom{N}{n_r}} \quad \text{or} \quad \mathbb{P}(U \geq k) \approx 1 - \Phi\left(\frac{k - 0.5 - \mu_U}{\sigma_U}\right).$$

To also control *over-concentration* (too many identical shares), impose  $I \leq l$  simultaneously and compute  $\mathbb{P}(U \geq k, I \leq l)$  either exactly by summing  $F_m$  on  $\{u \geq k, v \leq l\}$ , or via the bivariate normal CDF for  $(U, I)$ .

*All-or-nothing*  $(N, N)$  *additive sharing (special case)*. Here recovery requires full coverage,  $U = N$ . Compute

$$\mathbb{P}(U = N) = \sum_v \frac{F_m(N, v)}{\prod_r \binom{N}{n_r}} \quad \text{or} \quad \mathbb{P}(U = N) \approx 1 - \Phi\left(\frac{N - 0.5 - \mu_U}{\sigma_U}\right).$$

### 7.3 Blockchain and distributed consensus

*Blockchain data availability and sharding*. With erasure-coded blocks (need any  $k$  of  $N$  chunks), assign chunks to validators as  $P_r$ . Reconstructability is  $\mathbb{P}(U \geq k)$ ; redundancy and sampling fairness are controlled by  $I$  and by tails of  $Z_\ell/W_t$ . This provides explicit, finite- $N$  guarantees for data-availability sampling and shard replication.

*Blockchain protocols: PoW, PoS, PoB (selection overlap and propagation)*. When stake-weighted samplings occur over epochs, our intersection machinery ( $D_m$ ) measures expected and tail overlap of validator sets across epochs (sybil/capture risk), and  $Z_\ell$  tails quantify how many validators are repeatedly selected.

*Distributed consensus protocols (committee overlap and quorum safety)*. Across rounds, committees are subsets; the safety of BFT-style protocols hinges on sufficient *intersection* of honest quorums. Model two (or multiple) committees as parties and evaluate the intersection distribution with  $D_m$  (using Lemma 2) or via the CLT for  $(U, I)$  when  $N$  is large. This quantifies the probability of insufficient overlap (safety) or excessive overlap (centralisation risk).

### 7.4 Probabilistic guarantees for coverage and overlap

**Joint SLAs: cover enough, avoid hotspots** Pick targets  $k$  and  $l$  and certify

$$\mathbb{P}(U \geq k, I \leq l) \geq 1 - \delta.$$

Evaluate either by summing  $F_m$  on the rectangle  $\{u \geq k, v \leq l\}$  or with the bivariate normal CDF over  $[k - 0.5, \infty) \times (-\infty, l + 0.5]$ . You can invert this numerically to choose  $\{n_r\}$  (or  $m$ ) at reliability  $1 - \delta$ .

**Conditional redundancy from observed coverage** Given an observed coverage  $\hat{x}$  (e.g., via Bloom filter bit density), the bivariate normal gives

$$\mathbb{E}(I \mid U = \hat{u}) = \mu_I + \rho \frac{\sigma_I}{\sigma_U} (\hat{u} - \mu_U), \quad \text{Var}(I \mid U = \hat{u}) = \sigma_I^2 (1 - \rho^2),$$

with  $\rho = \text{Cov}(U, I)/(\sigma_U \sigma_I)$ . Use  $\mathbb{P}(I > l \mid U = \hat{u})$  as an online alarm for redundancy spikes or collusion risk.



## 8 Conclusion

We developed a unified distributional framework for random set systems under fixed-cardinality, without-replacement sampling. The organising principle is the *degree-count* vector  $\mathbf{Z} = (Z_0, \dots, Z_m)$ , where  $Z_\ell$  counts items covered by exactly  $\ell$  parties. Classic statistics are linear views of  $Z$ , for example  $U = \sum_{\ell \geq 1} Z_\ell$ ,  $I = Z_m$ ,  $W_t = \sum_{\ell \geq t} Z_\ell$ .

We introduced a Markovian threshold state whose one-step update is linear,  $\mathbf{S}_{k+1}^{(\tau)} = \mathbf{S}_k^{(\tau)} + M_\tau \mathbf{n}_{k+1}^{(\tau)}$ , with multivariate-hypergeometric weights and sharp feasibility bounds. This yields exact finite- $N$  PMFs for  $W_t$  (take  $\tau = t$ ) and  $Z_\ell$  (take  $\tau = \ell+1$ ). After aggregation, the same engine collapses to the global recursions e.g.  $\tau = 1 \Rightarrow C_m$  (union),  $\tau = k \Rightarrow D_m$  (intersection at step  $k$ ), and  $\tau = m \Rightarrow F_m$  (bivariate  $(U, I)$ ).

Using an exchangeable-pairs Stein coupling we proved a multivariate CLT for  $Z$  with  $O(N^{-1/2})$  accuracy, which immediately gives asymptotic normality for any fixed linear projection of  $Z$  (including  $(U, I)$ ,  $U-I$ , and  $W_t$ ), and delta-method surrogates for smooth indices such as the Jaccard  $J = U/I$ . Closed-form (finite- $N$ ) expressions for means and covariances come from one- and two-item probabilities.

The exact evaluators are pseudo-polynomial in  $N$  (feasible grid size), while the CLT depends only on  $m$ . Concretely, worst-case update counts are  $T = O(mN^2)$  for univariate cases and  $T = O(mN^3)$  for bivariate cases. By contrast, assembling the CLT needs only the first- and second-order structure of  $Z$ : the one-item pgf yields  $p_\ell$  in  $O(m^2)$ , and the two-item pgf yields all  $p_{\ell, \ell'}^{(2)}$  in  $O(m^3)$ ; hence *full*  $(m+1) \times (m+1)$  covariance construction is  $O(m^3)$ , independent of  $N$ . Thus for large  $N$  the workload drops from  $T = O(mN^2)$  (or worse) to  $T = O(m^3)$  when using the CLT, while retaining explicit, finite- $N$  validation paths whenever exactness is needed.

We demonstrated how the exact/CLT toolkit informs sizing and inference in probabilistic data structures (Bloom merges; MinHash planning), cryptographic mechanisms and blockchain protocols that rely on random subset selection (e.g., threshold recovery via  $W_t$ , committee overlaps via  $Y$ , and redundancy profiles via  $Z_\ell$ ).

## References

1. Barot, M., de la Peña, J.A.: Estimating the size of a union of random subsets of fixed cardinality. *Elemente der Mathematik* **56**(4), 163–169 (12 2001). <https://doi.org/10.1007/PL00000552>, <https://ems.press/journals/em/articles/581>
2. Beimel, A.: Secret-sharing schemes: A survey. In: Third international conference on Coding and cryptology. vol. 6639 LNCS, pp. 11–46 (2011). [https://doi.org/10.1007/978-3-642-20901-7\\_{\\_}2](https://doi.org/10.1007/978-3-642-20901-7_{_}2), [https://www.researchgate.net/publication/220776045\\_Secret-Sharing\\_Schemes\\_A\\_Survey](https://www.researchgate.net/publication/220776045_Secret-Sharing_Schemes_A_Survey)
3. Broder, A., Mitzenmacher, M.: Network applications of bloom filters: A survey. *Internet Mathematics* **1**(4), 485–509 (2004). <https://doi.org/10.1080/15427951.2004.10129096>; CTYPE:STRING:JOURNAL

4. Chatterjee, S., Meckes, E.: Multivariate normal approximation using exchangeable pairs. *Latin American Journal of Probability and Mathematical Statistics* **4**, 237–283 (1 2008), <https://arxiv.org/pdf/math/0701464>
5. Chung, N.C., Miasojedow, B.Z., Startek, M., Gambin, A.: Jaccard/Tanimoto similarity test and estimation methods for biological presence-absence data. *BMC Bioinformatics* **20**(15), 1–11 (12 2019). <https://doi.org/10.1186/S12859-019-3118-5/FIGURES/6>, <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3118-5>
6. Cohen, E.: Min-Hash Sketches. *Encyclopedia of Algorithms*, Second Edition pp. 1282–1287 (1 2016). [https://doi.org/10.1007/978-1-4939-2864-4\\_{\\_}573](https://doi.org/10.1007/978-1-4939-2864-4_{_}573), [https://link.springer.com/rwe/10.1007/978-1-4939-2864-4\\_573](https://link.springer.com/rwe/10.1007/978-1-4939-2864-4_573)
7. Kalinka, A.T.: The probability of drawing intersections: extending the hypergeometric distribution (5 2013), <https://arxiv.org/pdf/1305.0717>
8. McCormick, W.P., Lyons, N.I., Hutcheson, K.: Distributional properties of jaccard’s index of similarity. *Communications in Statistics - Theory and Methods* **21**(1), 51–68 (1 1992). <https://doi.org/10.1080/03610929208830764>;PAGE: STRING:ARTICLE/CHAPTER, <https://www.tandfonline.com/doi/abs/10.1080/03610929208830764>
9. McKay, B.D., Skerman, F.: Degree sequences of random digraphs and bipartite graphs. *Journal of Combinatorics* **7**(1), 21–49 (1 2016). <https://doi.org/10.4310/JOC.2016.V7.N1.A2>, <http://arxiv.org/abs/1302.2446><http://dx.doi.org/10.4310/JOC.2016.v7.n1.a2>
10. Miklós, I., Erdos, P.L., Soukup, L.: Towards Random Uniform Sampling of Bipartite Graphs with given Degree Sequence. *The Electronic Journal of Combinatorics* **20**(1) (1 2013). <https://doi.org/10.37236/3028>, <https://www.combinatorics.org/ojs/index.php/eljc/article/view/v20i1p16>
11. Real, R., Vargas, J.M.: The Probabilistic Basis of Jaccard’s Index of Similarity. *Systematic Biology* **45**(3), 380–385 (9 1996). <https://doi.org/10.1093/SYSBIO/45.3.380>, <https://dx.doi.org/10.1093/sysbio/45.3.380>
12. Reinert, G., Röllin, A.: Multivariate normal approximation with Stein’s method of exchangeable pairs under a general linearity condition. *Annals of Probability* **37**(6), 2150–2173 (12 2009). <https://doi.org/10.1214/09-AOP467>, <http://arxiv.org/abs/0711.1082><http://dx.doi.org/10.1214/09-AOP467>
13. Shamir, A.: How to share a secret. *Communications of the ACM* **22**(11), 612–613 (11 1979). <https://doi.org/10.1145/359168.359176>, <https://dl.acm.org/doi/10.1145/359168.359176>
14. Zachary P Neal, Cadieux, A., Diego Garlaschelli, Gotelli Id, N.J., Saracoid, F., Squartiniid, T., Shuttersid, S.T., Ulrichid, W., Wangid, G., Stronaid, G.: Pattern detection in bipartite networks: A review of terminology, applications, and methods. *PLOS Complex Systems* **1**(2), 1–34 (10 2024). <https://doi.org/10.1371/JOURNAL.PCSY.0000010>, <https://journals.plos.org/complexsystems/article?id=10.1371/journal.pcsy.0000010>

## A Exact recursions for $U, I, (U, I)$

This section give a few exact recursion formula using result from section 4.

### A.1 Univariate recursions for union ( $C_m$ )

*Two parties.* Let  $m = 2$  and set  $v_2 := n_1 + n_2 - u_2$ . Then  $|P_1 \cup P_2| = u_2$  iff  $|P_1 \cap P_2| = v_2$ . Conditioned on  $P_1$ ,

$$\mathbb{P}(U = u_2) = \mathbb{P}(|P_1 \cap P_2| = v_2) = \frac{\binom{n_1}{v_2} \binom{N-n_1}{n_2-v_2}}{\binom{N}{n_2}}, \quad v_2 \in \{\max(0, n_1+n_2-N), \dots, \min(n_1, n_2)\}.$$

Equivalently, the global count of ordered pairs with union  $x$  is

$$C_2(u_2) = \binom{N}{n_1} \binom{n_1}{v_2} \binom{N-n_1}{n_2-v_2}, \quad \text{so} \quad \mathbb{P}(U = u_2) = \frac{C_2(u_2)}{\binom{N}{n_1} \binom{N}{n_2}}.$$

*Three parties.* Conditioned on the union of the first two parties,  $u_2 := |P_1 \cup P_2|$ . To reach final union  $u_3$  after adding  $P_3$ : - pick exactly  $t := n_3 + u_2 - u_3$  elements of  $P_3$  from the previous union (so that  $P_3$  contributes  $u_3 - u_2$  new elements), and - pick  $u_3 - u_2$  new elements from outside the previous union.

Given a specific previous union of size  $u_2$ , the number of ways to do this is  $\binom{u_2}{t} \binom{N-u_2}{u_3-u_2}$ . Summing over feasible  $u_2$  and weighting by the count of  $(P_1, P_2)$  with union  $u_2$ ,

$$C_3(u_3) = \sum_{u_2} \binom{u_2}{n_3 + u_2 - u_3} \binom{N-u_2}{u_3-u_2} C_2(u_2),$$

hence

$$\mathbb{P}(U = u_3) = \frac{C_3(u_3)}{\prod_{i=1}^3 \binom{N}{n_i}}.$$

With the binomial-zero convention, the sum can be taken over all integers, but one convenient explicit range is

$$u_{2,\min} = \max\{n_1, n_2, u_3 - n_3, 0\}, \quad u_{2,\max} = \min\{n_1 + n_2, N\}.$$

*General  $m$  (global recursion).* Define the global counts  $C_m(u_m)$  of ordered  $m$ -tuples with union size  $u_1$  by the base

$$C_1(u_1) = \begin{cases} \binom{N}{n_1}, & u_1 = n_1, \\ 0, & \text{otherwise,} \end{cases}$$

and the recursion, for  $m \geq 2$ ,

$$C_m(u_m) = \sum_{u_{m-1}} \binom{u_{m-1}}{n_m + u_{m-1} - u_m} \binom{N-u_{m-1}}{u_m-u_{m-1}} C_{m-1}(u_{m-1}). \quad (22)$$

A convenient feasible range is

$$u_{\min} = \max\{n_{m-1}^{\max}, u_m - n_m, s_{m-1} - (m-2)N, 0\}, \quad u_{\max} = \min\{u_m, N, s_{m-1}\},$$

where  $s_{m-1} = \sum_{r=1}^{m-1} n_r$  and  $n_{m-1}^{\max} = \max_{r \leq m-1} n_r$ . Finally,

$$\mathbb{P}(U = u_m) = \frac{C_m(u_m)}{\prod_{i=1}^m \binom{N}{n_i}}. \quad (23)$$

## A.2 Univariate recursions for intersection ( $D_m$ )

*Two parties.* Conditioned on  $P_1$ , the intersection size is hypergeometric:

$$\mathbb{P}(I = v_2) = \frac{\binom{n_1}{v_2} \binom{N-n_1}{n_2-v_2}}{\binom{N}{n_2}}, \quad v_2 \in \{\max(0, n_1 + n_2 - N), \dots, \min(n_1, n_2)\}.$$

Equivalently, the global count is

$$D_2(v_2) = \binom{N}{n_1} \binom{n_1}{v_2} \binom{N-n_1}{n_2-v_2}, \quad \text{so} \quad \mathbb{P}(I = v_2) = \frac{D_2(v_2)}{\binom{N}{n_1} \binom{N}{n_2}}.$$

*Three parties.* Conditioned on the two-way intersection  $v_2 := |P_1 \cap P_2|$ . To have three-way intersection  $v_3$ , pick exactly  $v_3$  of those  $v_2$  common elements for  $P_3$ , and choose the remaining  $n_3 - v_3$  elements of  $P_3$  outside  $P_1 \cap P_2$ . This gives the factor  $\binom{v_2}{v_3} \binom{N-v_2}{n_3-v_3}$ . Summing over feasible  $v_2$ ,

$$D_3(v_3) = \sum_{v_2} \binom{v_2}{v_3} \binom{N-v_2}{n_3-v_3} D_2(v_2), \quad \mathbb{P}(I = v_3) = \frac{D_3(v_3)}{\prod_{i=1}^3 \binom{N}{n_i}}.$$

A convenient explicit range is

$$v_{2,\min} = \max\{v_3, n_1 + n_2 - N, 0\}, \quad v_{2,\max} = \min\{n_1, n_2\}.$$

(Outside this range the binomials vanish.)

*General  $m$  (global recursion).* Define global counts  $D_m(v_m)$  of ordered  $m$ -tuples with intersection size  $v_m$  by the base

$$D_1(v_1) = \begin{cases} \binom{N}{n_1}, & v_1 = n_1, \\ 0, & \text{otherwise,} \end{cases}$$

and the recursion, for  $m \geq 2$ ,

$$D_m(v_m) = \sum_{v_{m-1}} \binom{v_{m-1}}{v_m} \binom{N-v_{m-1}}{n_m-v_m} D_{m-1}(v_{m-1}). \quad (24)$$

A convenient feasible range is

$$v_{m-1,\min} = \max\{v_m, \sum_{r=1}^{m-1} n_r - (m-2)N, 0\}, \quad v_{m-1,\max} = \min\{n_{m-1}^{\min}\},$$

where  $n_{m-1}^{\min} = \min_{r \leq m-1} n_r$ . Finally,

$$\mathbb{P}(I = v_m) = \frac{D_m(v_m)}{\prod_{i=1}^m \binom{N}{n_i}}. \quad (25)$$

From standard feasibility constraints,

$$\max\left\{0, \sum_{r=1}^m n_r - (m-1)N\right\} \leq v_m \leq \min_{r \leq m} n_r, \quad \max_{r \leq m} n_r \leq u_m \leq N.$$

### A.3 Bivariate recursion ( $F_m$ )

For  $k \geq 1$ , write  $S_k = (n_1, \dots, n_k)$  and define the count

$$F_k(u_k, v_k) := \left| \left\{ (P_1, \dots, P_k) : |\cup_{r \leq k} P_r| = u_k, |\cap_{r \leq k} P_r| = v_k \right\} \right|.$$

Then the joint PMF of  $(U, I)$  is

$$\mathbb{P}(U = u_m, I = v_m) = p_{U,I}(u_m, v_m) = \frac{F_m(u_m, v_m)}{\prod_{r=1}^m \binom{N}{n_r}}, \quad (26)$$

for  $(u_m, v_m)$  in the feasible region  $\mathcal{R}$  (below), where  $F_m$  is defined recursively as follows.

*Base case* ( $m = 1$ ).

$$F_1(u_1, v_1) = \mathbf{1}\{u_1 = v_1 = n_1\} \binom{N}{n_1}$$

*Two parties* ( $m = 2$ ).

$$F_2(u_2, v_2) = \mathbf{1}\{u_2 = n_1 + n_2 - v_2\} \binom{N}{n_1} \binom{n_1}{v_2} \binom{N - n_1}{n_2 - v_2}. \quad (27)$$

*General recursion* ( $m \geq 3$ ). Let

$$s_{m-1} := \sum_{r=1}^{m-1} n_r, \quad n_{m-1}^{\min} := \min_{r \leq m-1} n_r, \quad n_{m-1}^{\max} := \max_{r \leq m-1} n_r.$$

Define bounds

$$\begin{aligned} a_{m-1} &:= \max\{v_m, s_{m-1} - (m-2)N\}, \\ b_{m-1} &:= n_{m-1}^{\min}, \\ c_{m-1}(v_{m-1}) &:= \max\left\{n_{m-1}^{\max}, \left\lceil \frac{s_{m-1} - v_{m-1}}{m-2} \right\rceil\right\}, \\ d_{m-1}(v_{m-1}) &:= s_{m-1} - (m-2)v_{m-1}. \end{aligned}$$

Then for all  $m \geq 3$ ,

$$F_m(u_m, v_m) = \sum_{v_{m-1}=a_{m-1}}^{b_{m-1}} \sum_{u_{m-1}=c_{m-1}(v_{m-1})}^{d_{m-1}(v_{m-1})} \binom{v_{m-1}}{v_m} \binom{u_{m-1}-v_{m-1}}{n_m+u_{m-1}-u_m-v_m} \binom{N-u_{m-1}}{u_m-u_{m-1}} \cdot F_{m-1}(u_{m-1}, v_{m-1}). \quad (28)$$

*Feasible region.*

$$\mathcal{R} = \left\{ (u_m, v_m) \in \mathbb{Z}^2 : s_m - (m-2)N \leq v_m \leq n_m^{\min}, \right. \\ \left. \max\{n_m^{\max}, \lceil \frac{s_m - v_m}{m-2} \rceil\} \leq u_m \leq s_m - (m-2)v_m \right\} \quad (29)$$

*Proof.* Condition on  $(u_{m-1}, v_{m-1})$  from the first  $m-1$  parties. To add party  $m$ , (i) keep  $v_m$  elements in the intersection:  $\binom{v_{m-1}}{v_m}$ ; (ii) choose  $s := n_m + u_{m-1} - u_m - v_m$  elements from the “union-only” band  $u_{m-1} - v_{m-1}$ :  $\binom{u_{m-1}-v_{m-1}}{s}$ ; and (iii) add  $t := u_m - u_{m-1}$  new elements outside the previous union:  $\binom{N-u_{m-1}}{t}$ . Note  $n_m = v_m + s + t$ . The bounds  $a_{m-1}, b_{m-1}, c_{m-1}, d_{m-1}$  are exactly those ensuring  $0 \leq v_m \leq v_{m-1}$ ,  $0 \leq s \leq u_{m-1} - v_{m-1}$ ,  $0 \leq t \leq N - u_{m-1}$ , and feasibility for the first  $m-1$  parties. Summing over feasible  $(u_{m-1}, v_{m-1})$  yields (28); normalising by  $\prod_r \binom{N}{n_r}$  gives (26).

## B Moments for the degree counts $Z_k$

We derive  $\mathbb{E}(Z_k)$ ,  $\text{Var}(Z_a)$ , and  $\text{Cov}(Z_a, Z_b)$  starting from single- and two-item probabilities.

*One-item probabilities and means.* For a single item  $i$  and  $d \in \{0, \dots, m\}$ ,

$$p_d(\alpha) = \mathbb{P}(R(i) = d) = [t^d] \prod_{r=1}^m ((1 - \alpha_r) + \alpha_r t).$$

Therefore  $\mathbb{E}(Z_d) = N p_k(\alpha)$ .

*Two-item probabilities.* Fix distinct items  $u \neq v$ . For party  $r$  the pair  $(\mathbf{1}\{u \in P_r\}, \mathbf{1}\{v \in P_r\})$  has probabilities

$$\pi_{00}^{(r)} = \frac{(N - n_r)(N - n_r - 1)}{N(N - 1)}, \quad \pi_{10}^{(r)} = \pi_{01}^{(r)} = \frac{n_r(N - n_r)}{N(N - 1)}, \quad \pi_{11}^{(r)} = \frac{n_r(n_r - 1)}{N(N - 1)}.$$

Define  $\phi_r(z, w) := \pi_{00}^{(r)} z + \pi_{10}^{(r)} w + \pi_{01}^{(r)} w + \pi_{11}^{(r)} zw$ . Across parties,

$$q_{a,b} := \Pr(R(u) = a, R(v) = b) = [z^a w^b] \prod_{r=1}^m \phi_r(z, w).$$

*Finite- $N$  variance-covariance identities.* Recall  $Z_d = \sum_{i=1}^N I_i^{(d)}$ . Then

$$\begin{aligned} \text{Var}(Z_a) &= \sum_{i=1}^N \text{Var}(I_i^{(a)}) + \sum_{u \neq v} \text{Cov}(I_u^{(a)}, I_v^{(a)}) = N p_a(1 - p_a) + N(N - 1)(q_{a,a} - p_a^2), \\ \text{Cov}(Z_a, Z_b) &= \sum_{u \neq v} \text{Cov}(I_u^{(a)}, I_v^{(b)}) - N p_a p_b = N(N - 1) q_{a,b} - N^2 p_a p_b, \quad a \neq b. \end{aligned}$$

These are exact for every  $N$ .

*Asymptotic expansion of  $q_{a,b}$ .* Write  $\alpha_r = n_r/N$ . Expand each  $\phi_r$  to first order in  $(N - 1)^{-1}$ :

$$\phi_r(z, w) = ((1 - \alpha_r) + \alpha_r z)((1 - \alpha_r) + \alpha_r w) + \frac{\alpha_r(1 - \alpha_r)}{N - 1}(-1 + z + w - zw).$$

Multiplying over  $r$  and extracting coefficients,

$$q_{a,b} = p_a p_b + \frac{1}{N - 1} S_{a,b}(\alpha) + O((N - 1)^{-2}),$$

where

$$S_{a,b}(\alpha) := \sum_{r=1}^m \alpha_r(1 - \alpha_r) \{ p_{a-1}^{(-r)} p_b^{(-r)} + p_a^{(-r)} p_{b-1}^{(-r)} - p_a^{(-r)} p_b^{(-r)} - p_{a-1}^{(-r)} p_{b-1}^{(-r)} \}. \quad (30)$$

Here  $p_d^{(-r)} := [t^d] \prod_{s \neq r} ((1 - \alpha_s) + \alpha_s t)$  and  $p_{-1}^{(-r)} := 0$ . Equivalently, in forward-difference form,

$$S_{a,b}(\alpha) = - \sum_{r=1}^m \alpha_r(1 - \alpha_r) \Delta_a^{(-r)} \Delta_b^{(-r)}, \quad \Delta_d^{(-r)} := p_d^{(-r)} - p_{d-1}^{(-r)}.$$

*Leading-order covariance.* Substitute  $q_{a,b}$  into the finite- $N$  identities to obtain

$$\Sigma_N = N \Sigma(\alpha) + O(1),$$

with entries

$$\Sigma_{aa}(\alpha) = p_a(1 - p_a) - \sum_{r=1}^m \alpha_r(1 - \alpha_r) (\Delta_a^{(-r)})^2, \quad \Sigma_{ab}(\alpha) = -p_a p_b - \sum_{r=1}^m \alpha_r(1 - \alpha_r) \Delta_a^{(-r)} \Delta_b^{(-r)} \quad (a \neq b).$$

The correction term is negative semidefinite (a sum of outer products with a minus sign); this reflects the weak negative dependence induced by sampling without replacement.

*Sanity checks.* (i) If parties included each item independently with probabilities  $\alpha_r$  (*with replacement*), then  $\phi_r(z, w) = ((1 - \alpha_r) + \alpha_r z)((1 - \alpha_r) + \alpha_r w)$ ,  $S_{a,b} \equiv 0$ , and  $\Sigma(\alpha) = \text{diag}(p) - pp^\top$  (multinomial covariance).

(ii) In the ultra-sparse regime where  $p_d \rightarrow 0$  and  $N p_d \rightarrow \lambda \in (0, \infty)$ , the count  $Z_d$  converges in law to  $\text{Poisson}(\lambda)$ .

## C Stein proof of the level-count CLT

This appendix proves Theorem 3 in full detail. We first explain the swap coupling; then we verify the three standard Stein ingredients: *regression*, *bounded increments*, and *matching second moments*. Finally, we invoke the multivariate Stein bound.

*The swap coupling (preserves  $|P_r| = n_r$ )* Given a configuration  $(P_1, \dots, P_m)$ , construct  $(P'_1, \dots, P'_m)$  as follows.

1. Pick an index  $K \sim \text{Unif}([N])$ .
2. For each party  $r = 1, \dots, m$  *independently of other parties*:
  - If  $K \in P_r$ , choose  $U_r \in [N] \setminus P_r$  uniformly and swap  $K \leftrightarrow U_r$  (remove  $K$  from  $P_r$ , add  $U_r$ ).
  - If  $K \notin P_r$ , choose  $V_r \in P_r$  uniformly and swap  $K \leftrightarrow V_r$  (add  $K$  to  $P_r$ , remove  $V_r$ ).

Each swap preserves  $|P_r| = n_r$ , so the joint law of  $(P'_1, \dots, P'_m)$  is again uniform over the same state space. The mapping is an involution (doing it twice returns to the starting configuration), so the pair

$$\mathbf{Z} := (Z_0, \dots, Z_m), \quad \mathbf{Z}' := (Z'_0, \dots, Z'_m)$$

is *exchangeable*. Define the increment  $\Delta := \mathbf{Z}' - \mathbf{Z}$ .

*How many coordinates can change?* Only the tagged item  $K$  and, for each party, at most one partner ( $U_r$  or  $V_r$ ) can change their coverage  $R(\cdot)$  by  $\pm 1$ . Hence at most  $2m$  items change level. When a single item changes its coverage by  $\pm 1$ , the level counts move one unit between two *adjacent* bins. Therefore, for each  $k$ ,  $\Delta_k \in \{-2m, \dots, 2m\}$  and, more sharply,

$$\|\Delta\|_1 \leq 2m, \quad \|\Delta\|_2 \leq 2\sqrt{m}.$$

These deterministic bounds yield moment bounds after scaling by  $N^{-1/2}$ .

*The regression property* Write  $\bar{\mathbf{Z}} := \mathbf{Z} - \mathbb{E}(\mathbf{Z})$  and let  $\mathbf{T} = \{x : \mathbf{1}^\top x = 0\}$  be the tangent subspace, with projector  $\Pi$  onto  $\mathbf{T}$ . We prove

$$\mathbb{E}[\Delta \mid \mathbf{Z}] = -\frac{1}{N} \Pi \bar{\mathbf{Z}}. \quad (31)$$

*Step 1: contribution from the tagged item  $K$ .* Condition on the entire configuration. The level of  $K$  *before* the swap, say  $R(K) = a$ , contributes  $-e_a$  to  $\Delta$  if  $K$  leaves level  $a$ , where  $e_a$  is the  $a$ -th standard basis vector in  $\mathbb{R}^{m+1}$ . *After* the full set of party swaps, the coverage of  $K$  is re-drawn from its *stationary one-item law* given the rest, whose mean indicator vector is the population average  $\mathbb{E}(\mathbf{Z})/N$ . Therefore, the conditional expectation of the *new* level vector for  $K$  equals  $(\mathbb{E}(\mathbf{Z}))/N$ . Averaging over  $K \sim \text{Unif}([N])$  and using  $\mathbb{P}(R(K) = a \mid \mathbf{Z}) = Z_a/N$  (exchangeability of items), we get

$$\mathbb{E}[(\text{change caused by } K) \mid \mathbf{Z}] = \frac{\mathbb{E}(\mathbf{Z})}{N} - \frac{\mathbf{Z}}{N} = -\frac{1}{N} \bar{\mathbf{Z}}.$$



*Step 2: contributions from partner items.* For a fixed party  $r$ , the partner ( $U_r$  or  $V_r$ ) is chosen uniformly from the complement or from  $P_r$  respectively. Conditioned on  $K$  and the sets  $P_r$ , the two cases are symmetric and have opposite signs in expectation. Therefore, their *net* conditional drift is 0. Summing over parties yields zero expected contribution from all partners. Combining with Step 1 gives (31). Since  $\mathbf{1}^\top \Delta = 0$  deterministically, we may insert  $\Pi$  without changing the equality.

*Bounded increments and third moments* Define the scaled, projected statistic

$$W := \frac{1}{\sqrt{N}} (\Pi \Sigma(\alpha) \Pi^\top)^{-1/2} \Pi \bar{\mathbf{Z}}, \quad W' := W + \frac{1}{\sqrt{N}} (\Pi \Sigma(\alpha) \Pi^\top)^{-1/2} \Pi \Delta.$$

Using  $\|\Pi\| \leq 1$ ,  $\|(\Pi \Sigma(\alpha) \Pi^\top)^{-1/2}\| = O(1)$  (fixed  $m$ , fixed  $\alpha$ ), and  $\|\Delta\|_2 \leq 2\sqrt{m}$ ,

$$\|W' - W\| \leq C \frac{\|\Delta\|_2}{\sqrt{N}} \leq \frac{C'}{\sqrt{N}}, \quad \mathbb{E}\|W' - W\|^3 = O(N^{-3/2}).$$

*Conditional second moments match  $\Sigma(\alpha)$*  Let  $V := \mathbb{E}[\Delta \Delta^\top \mid Z]$  and  $V_0 := \mathbb{E}V$ . We claim

$$\frac{1}{N} V_0 = \Pi \Sigma(\alpha) \Pi^\top, \quad \mathbb{E}\|V - V_0\|_{\text{HS}} = O(N^{-1}). \quad (32)$$

Intuition: each swap touches only  $O(1)$  items, so the change in quadratic forms depends only on  $O(1)$  coordinates of  $\mathbf{Z}$ , leading to  $O(1/N)$  fluctuations.

*Derivation of  $V_0$ .* Pick  $K$  uniformly. The increment  $\Delta$  is a sum of at most  $2m$  elementary moves that transfer one unit between adjacent levels. Averaging over the random positions of these moves and over the random  $K$ , the expected outer product  $\mathbb{E}[\Delta \Delta^\top]$  reproduces the *per-item* contributions to the covariance of  $Z$  described in §5.3, scaled by  $1/N$ . Carrying out the coefficient extraction (Appendix B) gives the first relation in (32). The second relation follows because  $\Delta$  depends on  $Z$  only through the levels of at most  $2m$  items and the choice of  $K$ , giving  $O(1/N)$  variability in conditional second moments.

*Conclusion via multivariate Stein* We now apply the exchangeable-pairs version of multivariate Stein's method: if  $(W, W')$  is exchangeable,  $\mathbb{E}[W' - W \mid W] = -\lambda W + R$  with  $\lambda = 1/N$  and  $\|R\|$  negligible, and if the conditional second moments match the target covariance up to  $O(1/N)$  while  $\mathbb{E}\|W' - W\|^3 = O(N^{-3/2})$ , then for smooth test functions  $h$ ,

$$|\mathbb{E}(h(W)) - \mathbb{E}(h(Z))| \leq \frac{C}{\sqrt{N}},$$

where  $Z \sim \mathcal{N}(0, I_m)$  and  $C$  depends only on  $m$  and  $\alpha$ . A standard smoothing step transfers this bound to the distance induced by indicators of convex sets. This yields Theorem 3.  $\square$

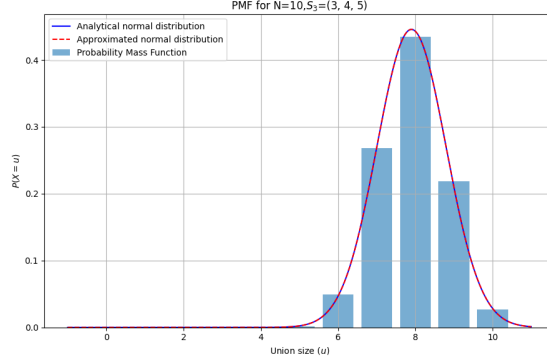
## D Numerical validation

All figures in this section use the *exact* finite- $N$  counts from the paper: the bivariate recursion  $F_m$  in A.3, Equation 28 for  $(U, I)$  and its univariate specialisations  $C_m$  (unions) and  $D_m$  (intersections) from Appendix A.1,A.2. Gaussian overlays are parameterised by the closed-form moments from §6.3, and for Jacard by the delta-method  $(\mu_J, \sigma_J^2)$  in §6.4. Where a continuous normal density is drawn over a discrete PMF, we apply a half-cell continuity correction.

### D.1 Univariate distribution

The univariate plots use the global recursions  $C_m$  (union) and  $D_m$  (intersection) from Appendix A, which are the marginals of  $F_m$  (26). We show (i) exact PMFs, (ii) tail probabilities relevant for design, and (iii) large- $N$  Gaussian overlays.

*Union.* Figure 1 shows the exact PMF  $\mathbb{P}(U = u)$  for  $N = 10$  and  $S_3 = (3, 4, 5)$ . Mass at  $u = N$  corresponds to full coverage (cf. the recovery event in §7); the mode lies close to  $\mu_U$  from §6.3.

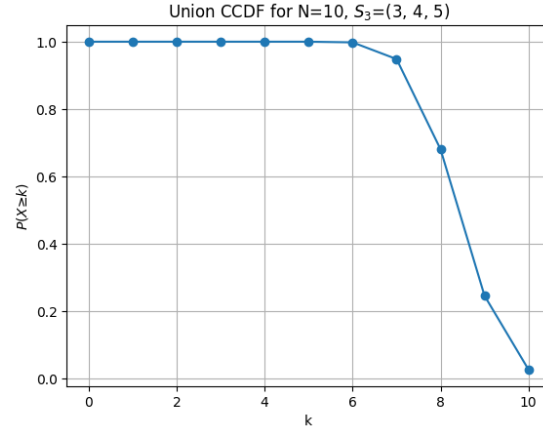


**Fig. 1.** Exact PMF of  $U$  for  $N = 10$  and  $S_3 = (3, 4, 5)$  (via recursion  $C_3$ ).

For tail requirements in §7, we evaluate

$$\mathbb{P}(U \geq k) = \sum_{u=k}^N \frac{C_m(u)}{\prod_{r=1}^m \binom{N}{n_r}}, \quad (33)$$

displayed in Figure 2 for the same parameters.

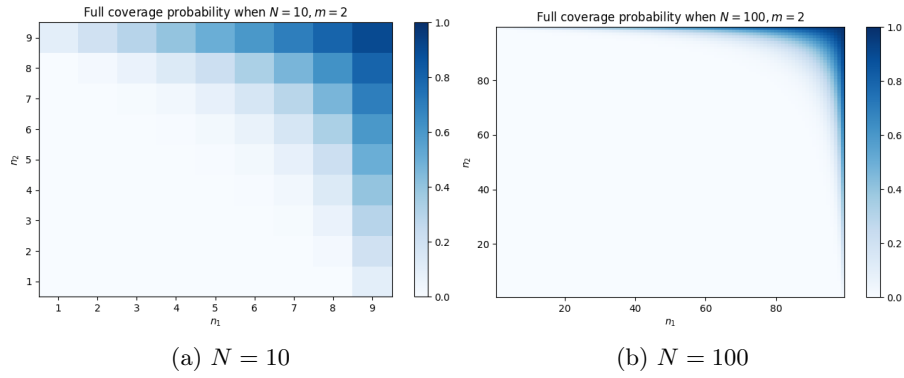


**Fig. 2.** Tail probability  $\mathbb{P}(U \geq k)$  for  $N = 10$  and  $S_3 = (3, 4, 5)$  (via  $C_3$ ).

The full-coverage probability is

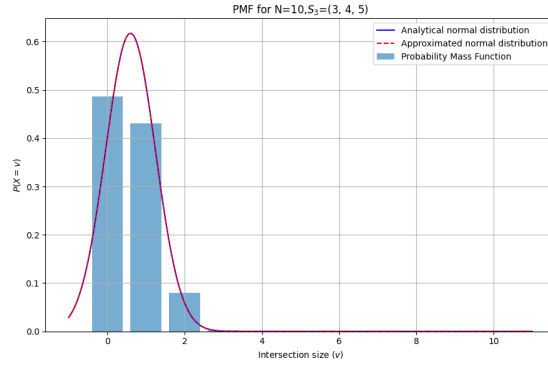
$$\mathbb{P}(U = N) = \frac{C_m(N)}{\prod_{r=1}^m \binom{N}{n_r}},$$

and its dependence on  $(n_1, n_2)$  for  $m = 2$  is summarised in Figure 3 (useful for merged-filter sizing in §7).



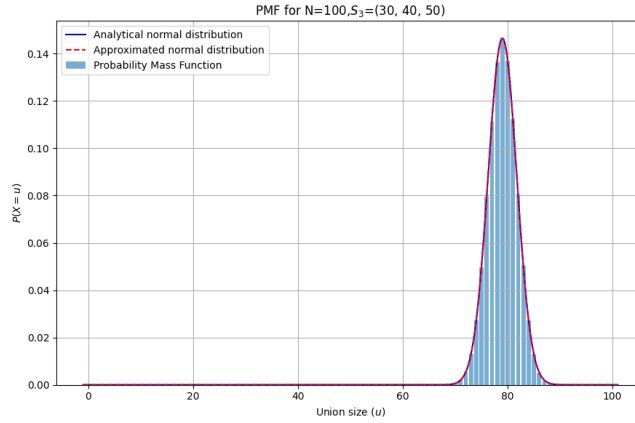
**Fig. 3.** Full-coverage probability  $\mathbb{P}(U = N)$  for  $m = 2$  across  $(n_1, n_2)$  (via  $C_2$ ).

*Intersection.* Similarly, Figure 4 shows the exact PMF  $\mathbb{P}(I = v)$  for  $N = 10$  and  $S_3 = (3, 4, 5)$ .

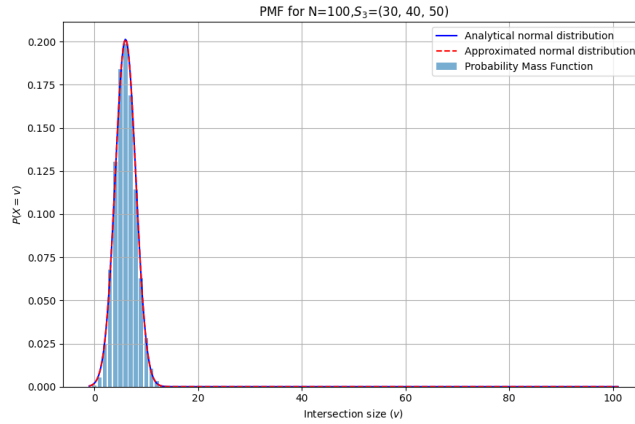


**Fig. 4.** Exact PMF of  $I$  for  $N = 10$  and  $S_3 = (3, 4, 5)$  (via recursion  $D_3$ ).

*Large- $N$  behaviour.* Theorem 3 gives an  $O(N^{-1/2})$  quantitative CLT for  $(U, I)$  after standardisation, univariate normal approximations follow. Figures 5 and 6 overlay  $\mathcal{N}(\mu_U, \sigma_U^2)$  and  $\mathcal{N}(\mu_I, \sigma_I^2)$  (moments from §6.3) on the exact PMFs. Agreement is excellent away from the boundaries  $u \approx N$  and  $v \approx 0$  or  $\min_r n_r$ ; in extremely sparse-overlap regimes (very small  $p_I$ ), using a moment-matched binomial distribution can be sharper (§6.2).



**Fig. 5.** Exact PMF of  $U$  for  $N = 100$  and  $S_3 = (30, 40, 50)$  with normal overlay  $\mathcal{N}(\mu_U, \sigma_U^2)$ ; continuity correction applied.



**Fig. 6.** Exact PMF of  $I$  for  $N = 100$  and  $S_3 = (30, 40, 50)$  with normal overlay  $\mathcal{N}(\mu_I, \sigma_I^2)$ ; continuity correction applied.

## D.2 General univariate distribution

## D.3 Bivariate distribution

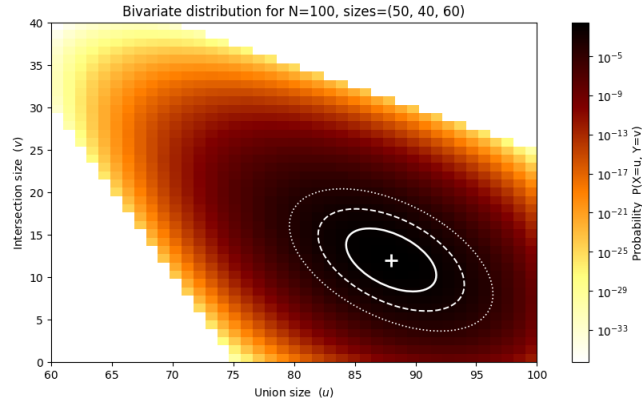
Figure 7 shows the *exact* joint PMF  $p_{U,I}$  from (26) (via  $F_m$ ) as a heat map, with isocontours of the Gaussian surrogate  $\mathcal{N}_2(\mu, \Sigma)$  overlaid. The orientation and eccentricity of the ellipses are determined by the covariance  $\text{Cov}(U, I)$  in §6.3; the close alignment of isocontours with PMF degree sets empirically supports the approximation guaranteed by the bivariate CLT, which provides an  $O(N^{-1/2})$  finite- $N$  error for the standardised vector. As expected, departures are most visible (not shown) when  $(u, v)$  sits near the boundary of the feasible region  $\mathcal{R}$  (defined in Appendix A.3), where lattice effects and truncation become non-negligible.

## D.4 Jaccard distribution

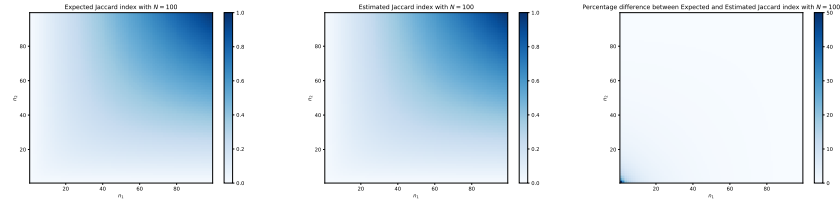
We validate both the *exact* Jaccard constructions from §6.4 and the delta-method surrogate.

*Expected value over  $(n_1, n_2)$ .* Panel (a) of Figure 8 plots the exact  $\mathbb{E}[J_N]$  for  $m = 2$ , obtained by summing the hypergeometric law of  $I = |P_1 \cap P_2|$  (Proposition 1). Panel (b) shows the delta-method mean  $\mu_J$  derived from  $(\mu_N, \Sigma_N)$ ; panel (c) reports the relative error.

*Empirical histogram vs two normal fits.* Figure 9 compares the *empirical* distribution of  $J_N$  (histogram from Monte Carlo draws of  $(P_1, \dots, P_m)$  under fixed-cardinality sampling) with two normal curves: (i) a moment-fit normal using the empirical mean/SD of  $J_N$ ; and (ii) the *delta-method* normal  $\mathcal{N}(\mu_J, \sigma_J^2)$  computed from



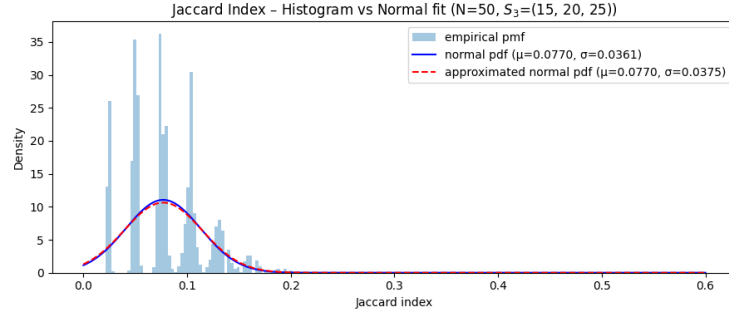
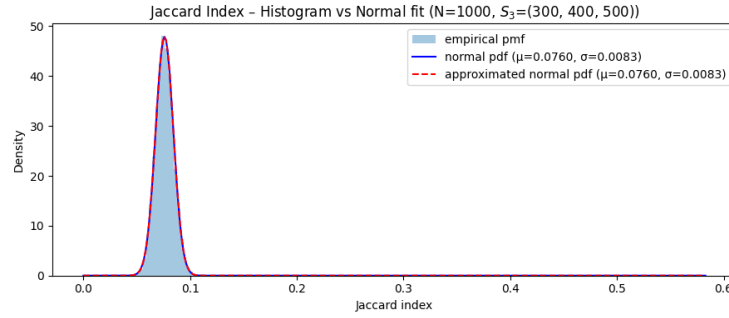
**Fig. 7.** Exact joint PMF of  $(U, I)$  (heat map via  $F_m$ ) with  $\mathcal{N}_2(\mu_N, \Sigma_N)$  isocontours (moments from §6.3) overlaid; continuity correction applied.



(a) Exact  $\mathbb{E}[J_N]$  (Prop. 1) (b) Delta-method  $\mu_J$  (§6.4) (c)  $100 \times |\mathbb{E}[J_N] - \mu_J| / \mathbb{E}[J_N]$

**Fig. 8.** Expected Jaccard index across  $(n_1, n_2)$  with  $N = 100$ . Exact (a) vs delta-method (b); relative error (c).

$(\mu_N, \Sigma_N)$  in §6.3. As  $N$  increases, the normal curves track the histogram increasingly closely and the distribution narrows, consistent with the  $N^{-1/2}$  scaling predicted by the CLT and delta method.

(a)  $N = 50$ .(b)  $N = 1000$ .

**Fig. 9.** Exact PMF of  $J_N$  (via Prop. 1) for a representative parameter choice versus the delta-method  $\mathcal{N}(\mu_J, \sigma_J^2)$ . Continuity correction and truncation to  $[0, 1]$  applied to the overlay.