# INVESTIGATIVE SINGLE CELL GENETICS: DETERMINING THE PROBABILITY OF A GENOTYPE AFTER OBSERVING SINGLE-CELL DATA ENABLES ROBUST DATABASE SEARCHES ACROSS MIXTURE COMPLEXITIES

LFTDI

Qhawe A. Bhembe[a] M.S.; Nidhi C. Sheth[a] M.S.; Leah O'Donnell[c] M.S.; Ken R. Duffy[b] Ph.D.; Desmond S. Lun[a] Ph.D.; Catherine M. Grgicak[a] Ph.D.

[a]Rutgers University, Camden, US; [b]Northeastern University, Boston, US; [c]Maynooth University, Ireland
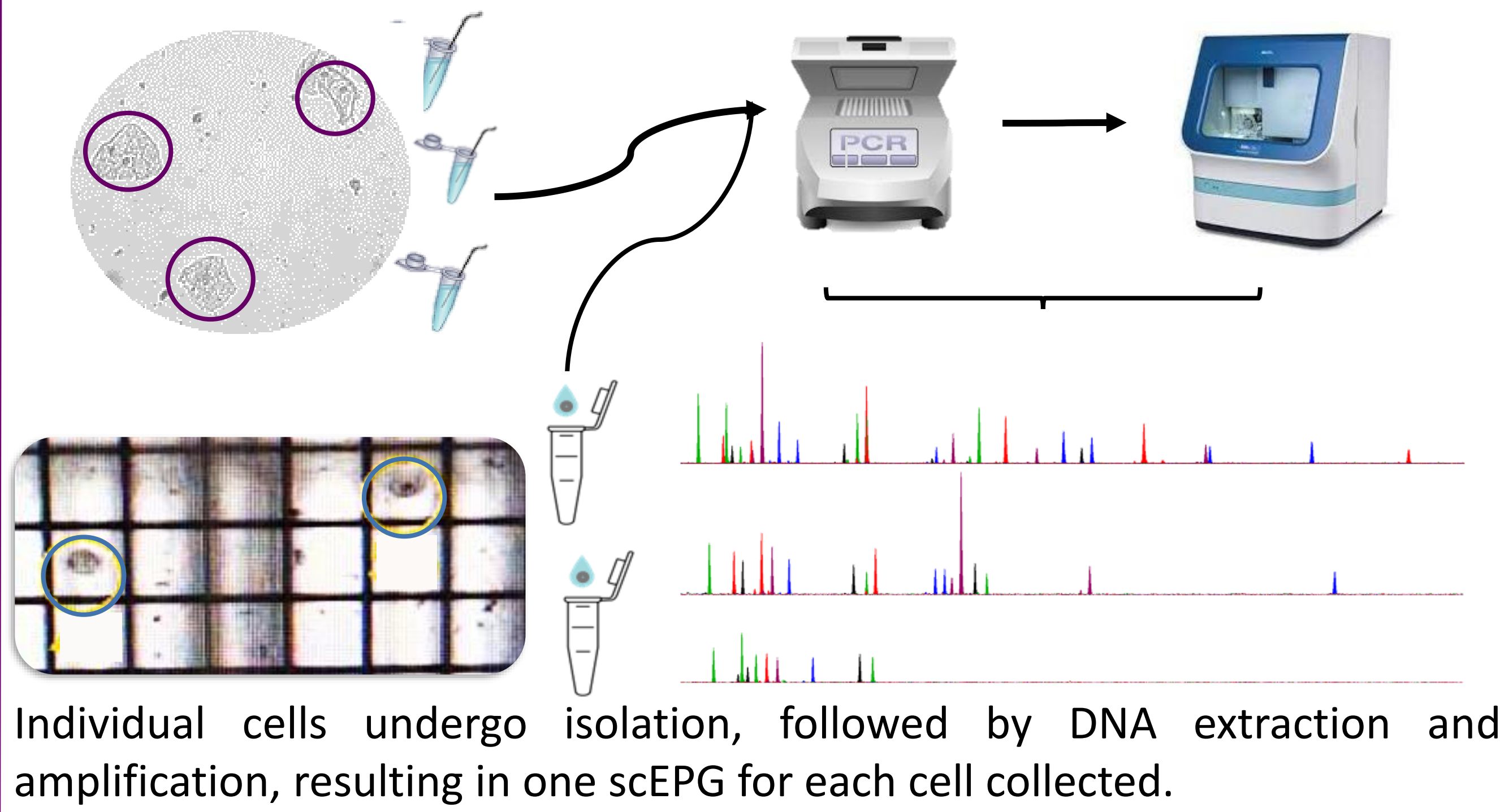
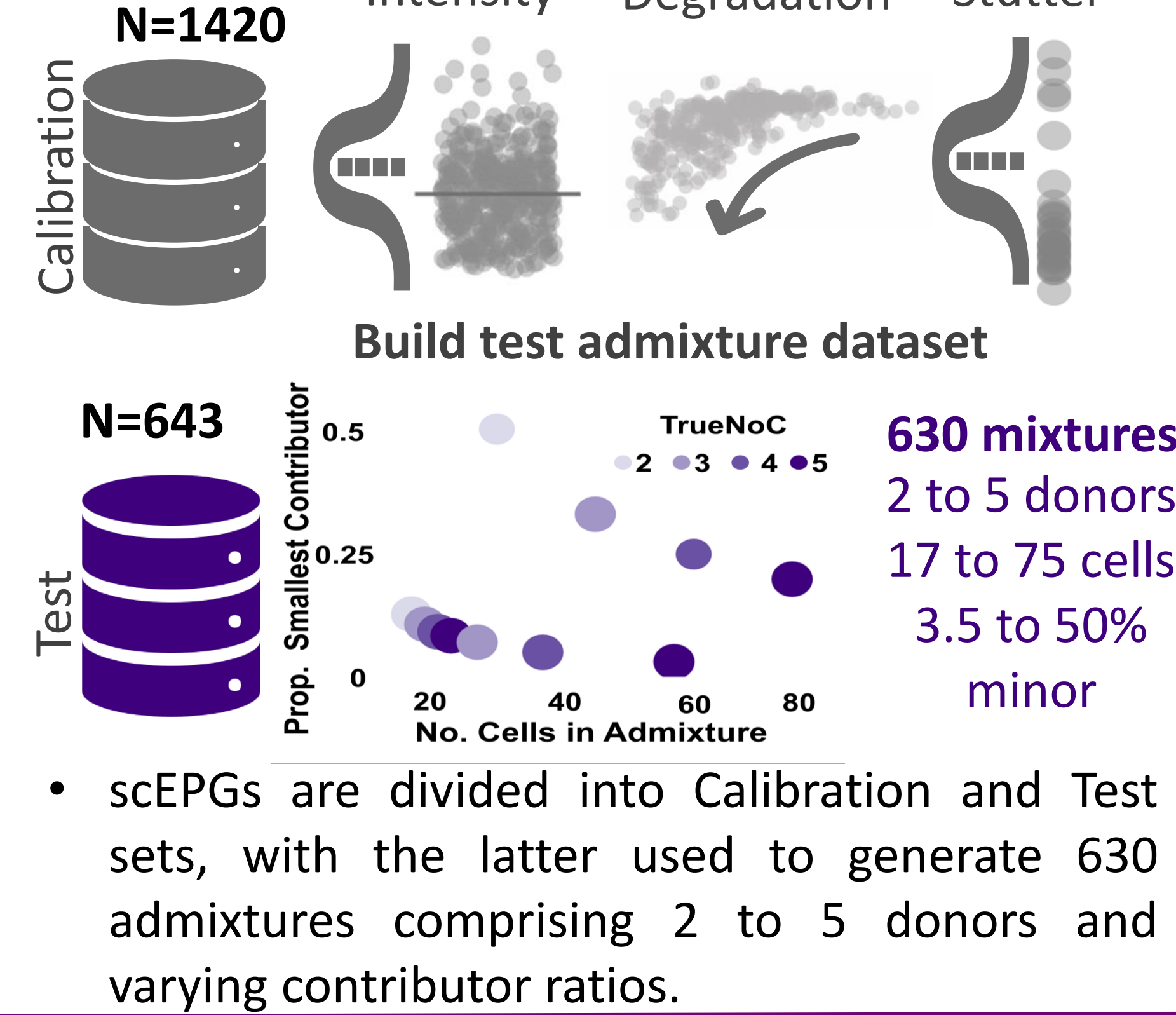RUTGERS UNIVERSITY | CAMDEN

## Highlights

- We introduce a Bayesian approach that interprets forensic single-cell data in a suspect-agnostic manner, by
  - grouping single-cell electropherograms (scEPGs) via model-based clustering (MBC), based on scEPG similarity
  - using EESCIt™ – *Evidentiary Evaluation of Single Cells* – which determines the probability of observing the cluster of scEPGs given all possible genotypes
  - then determining the probability of a genotype, g, given the cluster, C, of scEPGs:

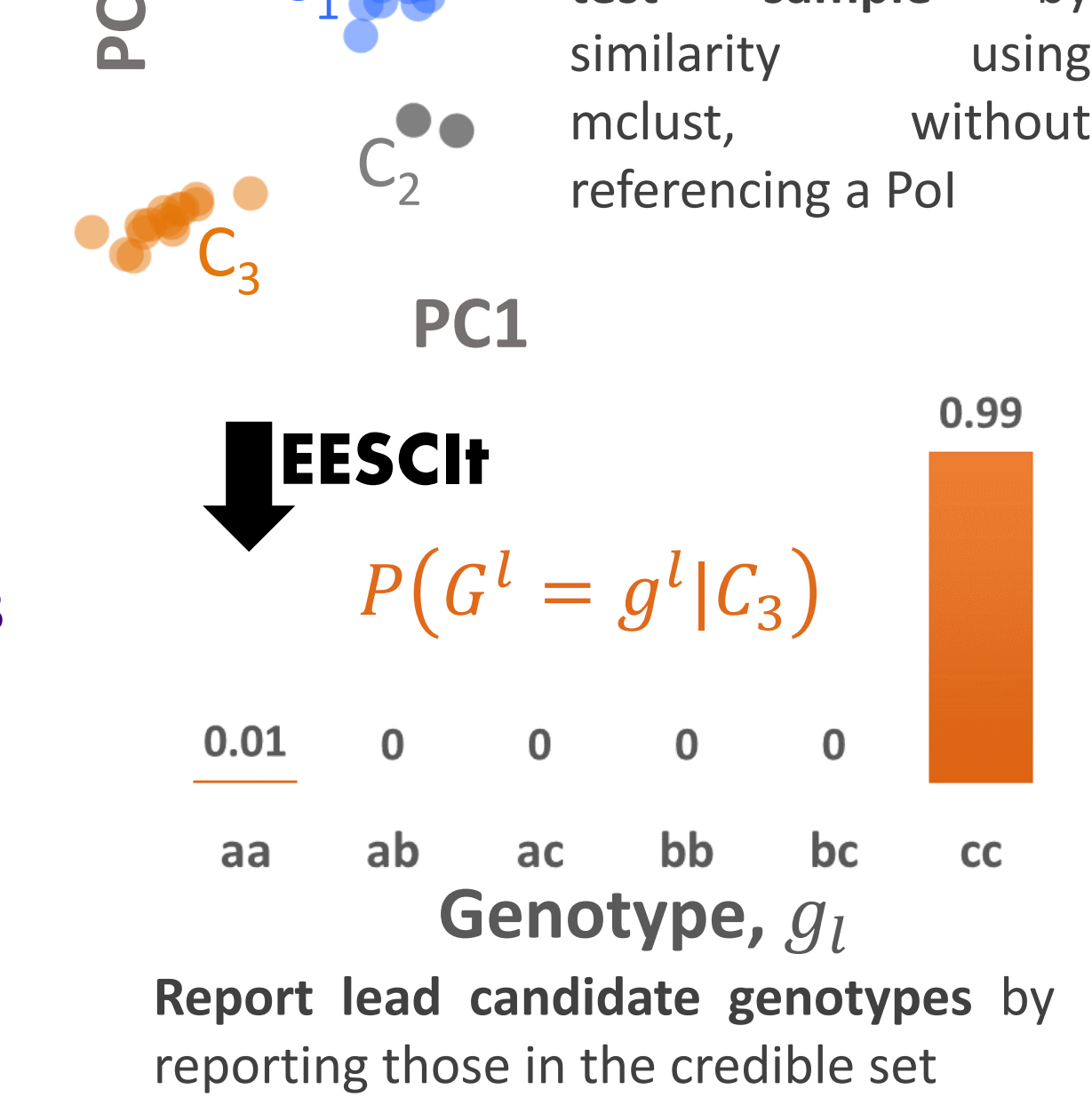$$P(G = g|C) = \frac{P(C|G = g)P(G = g)}{\sum P(C|G = g_i)P(G = g_i)}$$
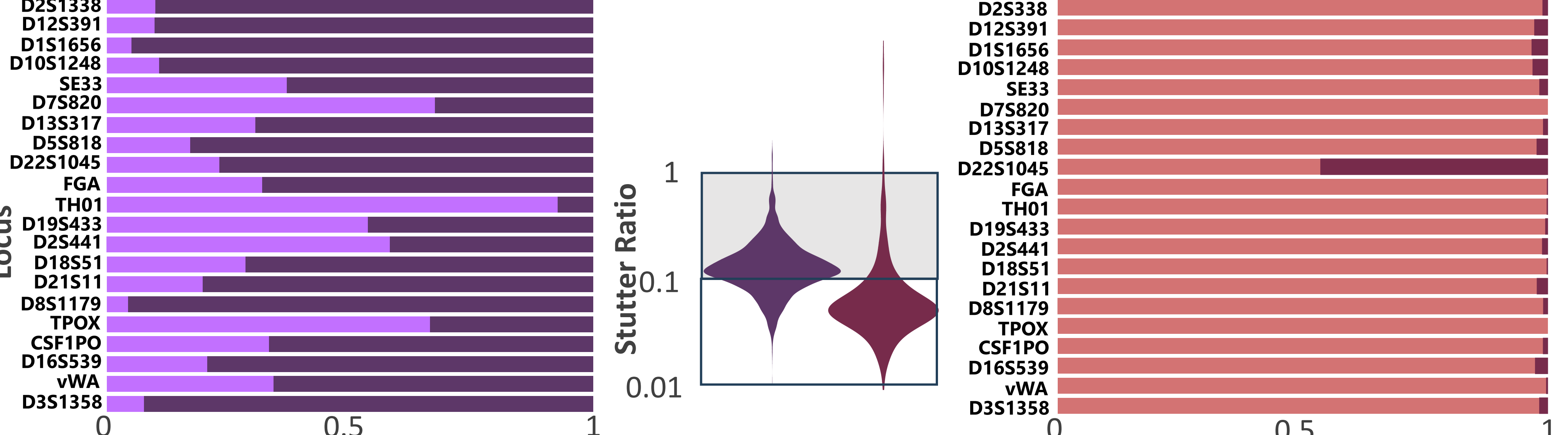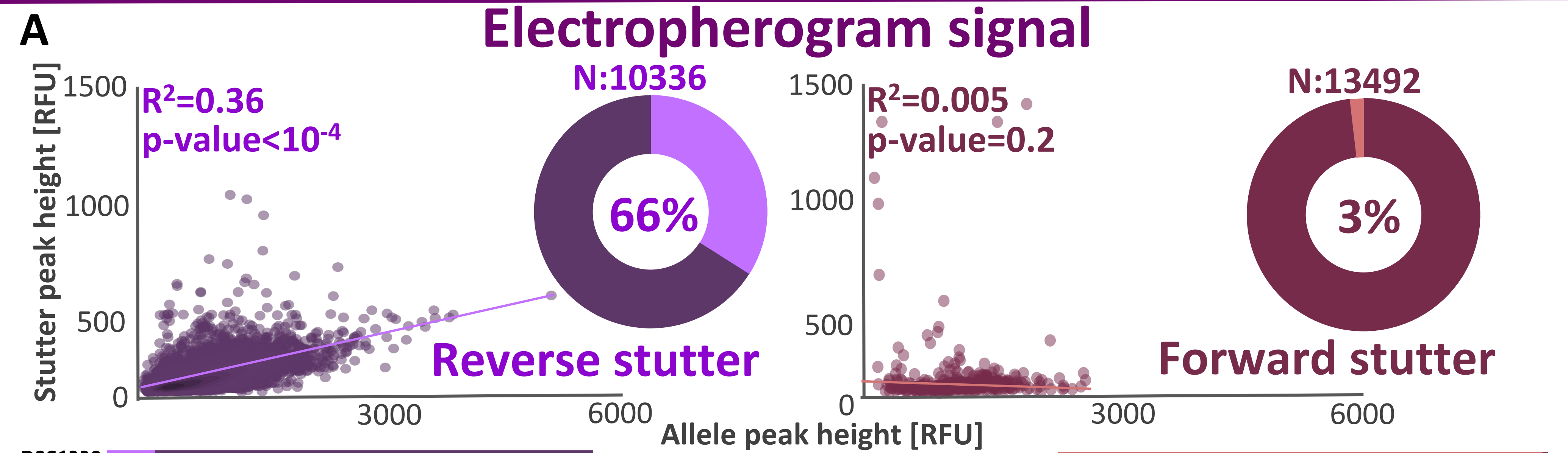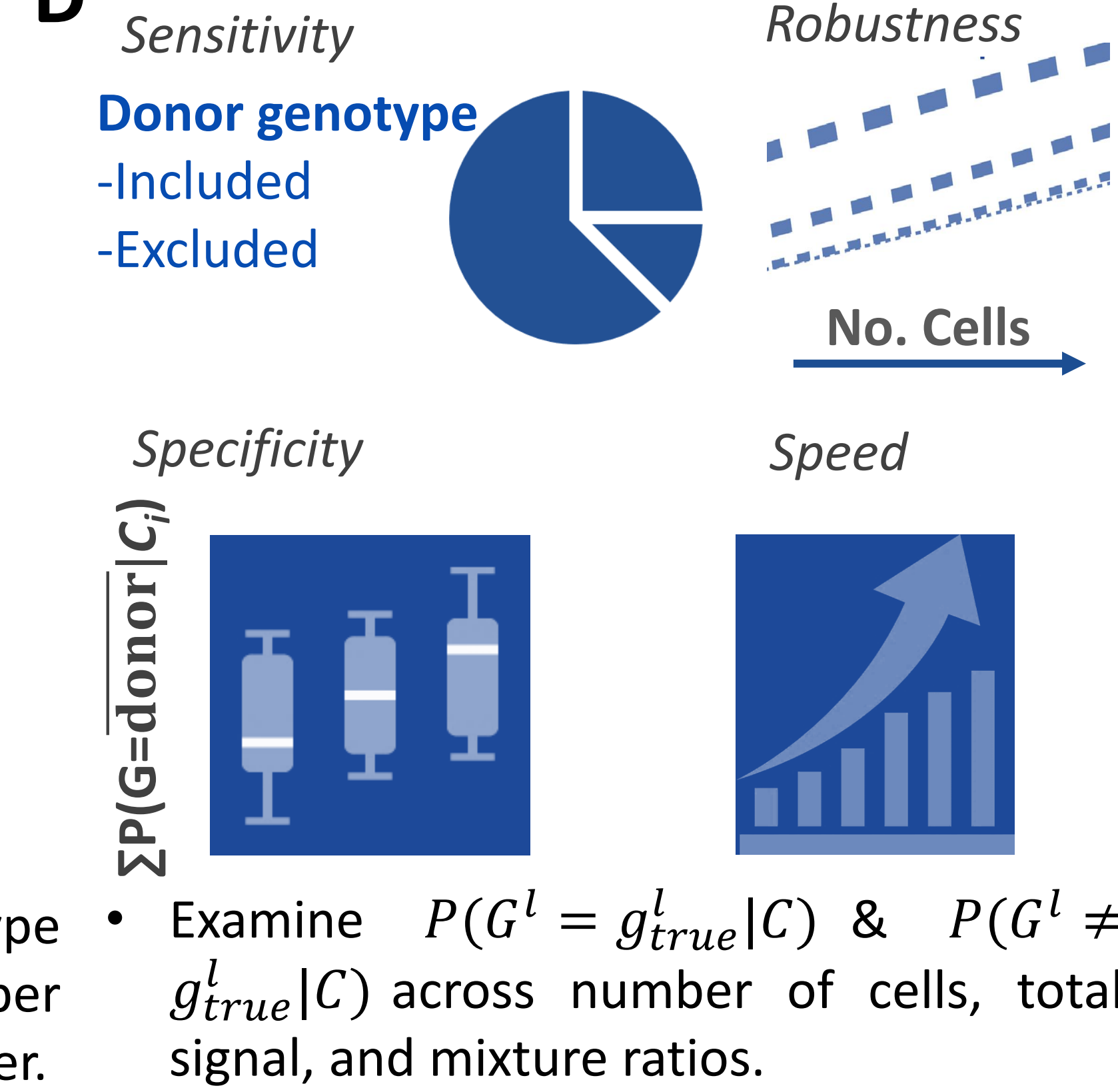
## A — Single-cell forensic workflow



Individual cells undergo isolation, followed by DNA extraction and amplification, resulting in one scEPG for each cell collected.

## B — Fit EESCIt models to scEPG calibration dataset

Intensity | Degradation | Stutter

Calibration N=1420

**Build test admixture dataset**

Test N=643

**630 mixtures**
2 to 5 donors
17 to 75 cells
3.5 to 50% minor

- scEPGs are divided into Calibration and Test sets, with the latter used to generate 630 admixtures comprising 2 to 5 donors and varying contributor ratios.

## C

Cluster scEPGs in each test sample by similarity using mclust without referencing a PoI

EESCIt

$$P(G^l = g^l|C_3)$$

| 0.01 | 0 | 0 | 0 | 0 | 0.99 |
|---|---|---|---|---|---|
| aa | ab | ac | bb | bc | cc |

Genotype, $g_l$

Report lead candidate genotypes by reporting those in the credible set

- Clustering enables genotype probability determinations per locus, given the scEPGs in a cluster.

## D — Performance assessment

Sensitivity | Robustness

**Donor genotype**
-Included
-Excluded

No. Cells

Specificity | Speed

$$\sum P(G = \overline{donor}|C_i)$$

- Examine $P(G^l = g^l_{true}|C)$ & $P(G^l \neq g^l_{true}|C)$ across number of cells, total signal, and mixture ratios.

---

# Electropherogram signal

## A

$R^2=0.36$, p-value<$10^{-4}$ — **Reverse stutter** — N:10336 — **66%**

$R^2=0.005$, p-value=0.2 — **Forward stutter** — N:13492 — **3%**

Stutter peak height [RFU] vs Allele peak height [RFU]

Loci: D2S1338, D12S391, D1S1656, D10S1248, SE33, D7S820, D13S317, D5S818, D22S1045, FGA, TH01, D19S433, D2S441, D18S51, D21S11, D8S1179, TPOX, CSF1PO, D16S539, vWA, D3S1358
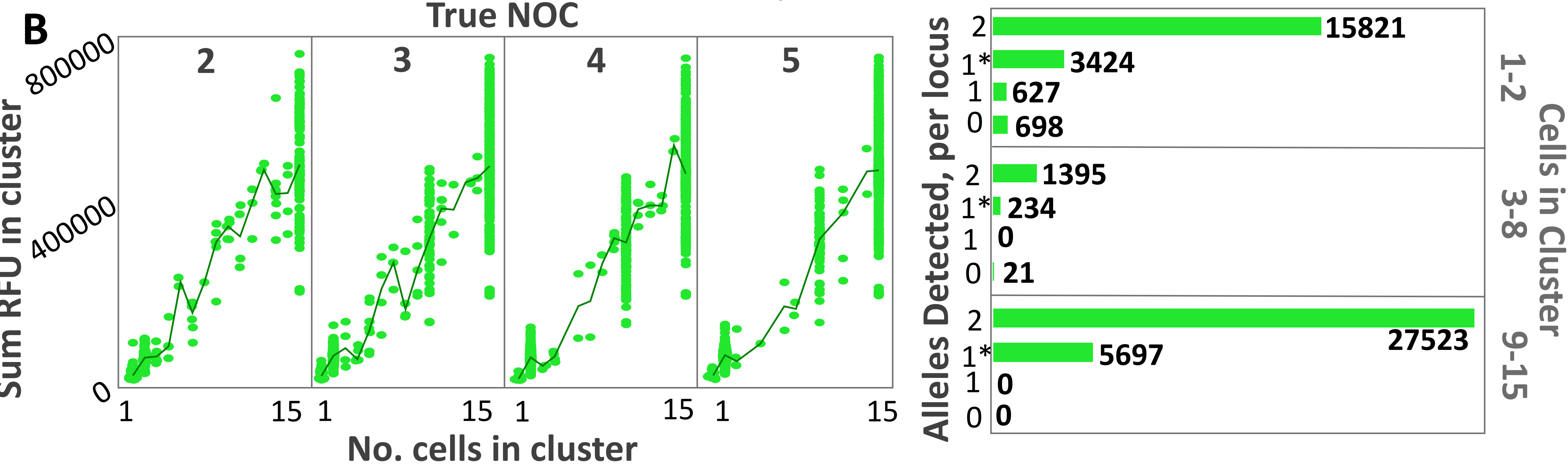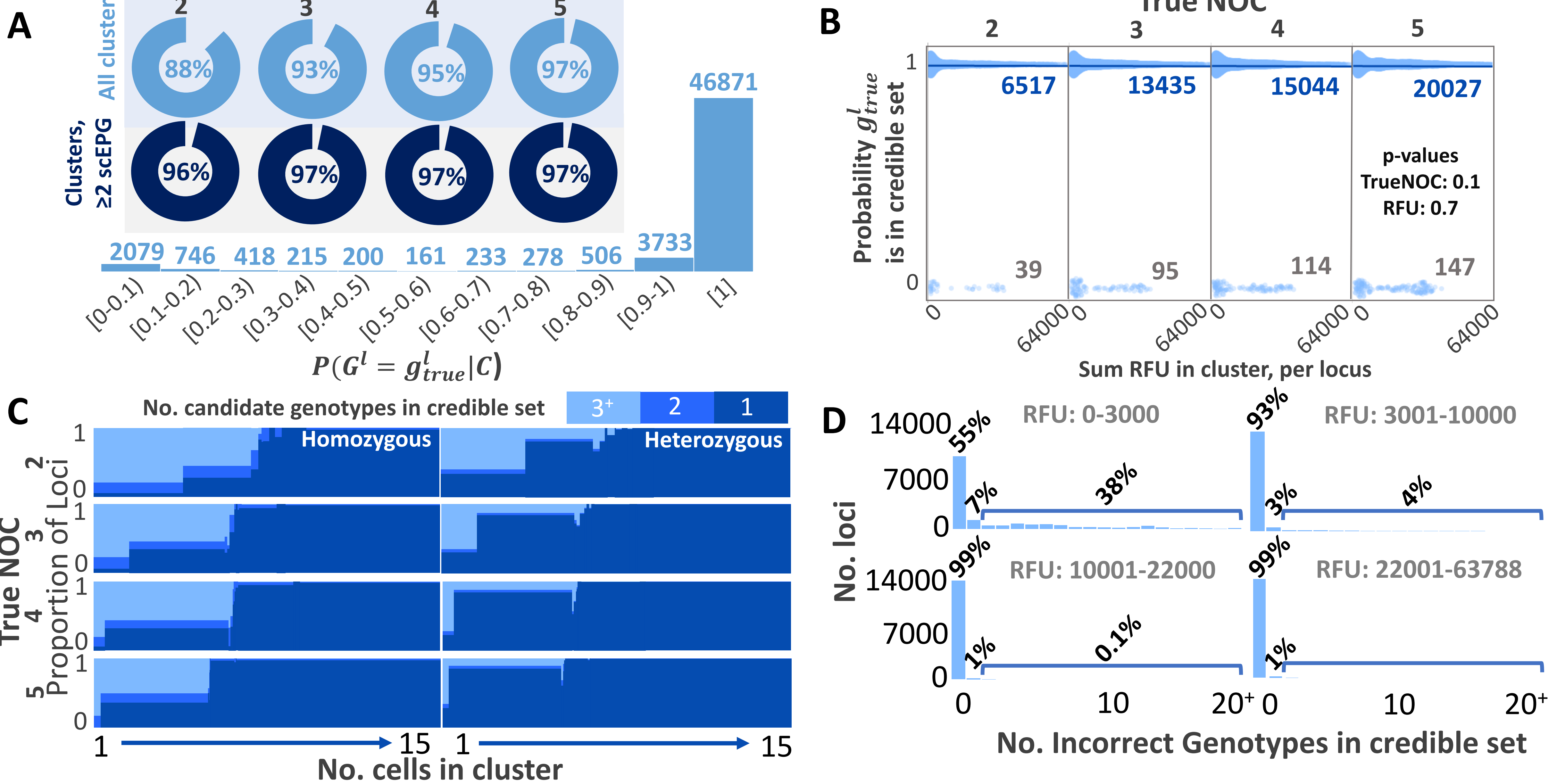
Stutter Ratio

(A) Pie charts display 66% and 3% of positions in reverse and forward stutter, respectively, with an RFU intensity greater than five, given that the allele position also has an RFU greater than five. The stacked plots separate these stutter proportions by locus, indicating a potential locus effect on reverse and forward stutter. The scatterplots present non-zero stutter peak heights (RFU) against allelic signal (RFU). The violin plot of the stutter ratio on a log scale. Average stutter ratio is 0.09 with 0.23 being the 95th percentile. **(B)** Scatter plots show the total RFU intensity across clusters with the number of cells in a cluster ranging from 1-15 cells, separated by NoC. The bar graph presents the number of detected alleles, grouped by the number of cells in a cluster, with an asterisk (*) denoting homozygous loci. As the number of cells in a cluster increases, the total RFU intensity increases, and hence the information content.

## B

True NOC: 2, 3, 4, 5

Sum RFU in cluster vs No. cells in cluster

Alleles Detected, per locus:

| Cells in Cluster 1-2 | |
|---|---|
| 2 | 15821 |
| 1* | 3424 |
| 1 | 627 |
| 0 | 698 |

| Cells in Cluster 3-8 | |
|---|---|
| 2 | 1395 |
| 1* | 234 |
| 1 | 0 |
| 0 | 21 |

| Cells in Cluster 9-15 | |
|---|---|
| 2 | 27523 |
| 1* | 5697 |
| 1 | 0 |
| 0 | 0 |

---

# Performance assessment

## A

**True NOC**

All clusters: 2 → 88%, 3 → 93%, 4 → 95%, 5 → 97%

Clusters, ≥2 scEPG: 2 → 96%, 3 → 97%, 4 → 97%, 5 → 97%

$P(G^l = g^l_{true}|C)$

[0-0.1]: 2079, 746; [0.1-0.2]: 418; [0.2-0.3]: 215; [0.3-0.4]: 200; [0.4-0.5]: 161; [0.5-0.6]: 233; [0.6-0.7]: 278; [0.7-0.8]: 506; [0.8-0.9]: 3733; [0.9-1]: ; [1]: 46871

## B

**True NOC**: 2, 3, 4, 5

Probability $g^l_{true}$ is in credible set

2 → 6517, 39; 3 → 13435, 95; 4 → 15044, 114; 5 → 20027, 147

p-values
TrueNOC: 0.1
RFU: 0.7

Sum RFU in cluster, per locus (0–64000)

## C

No. candidate genotypes in credible set: 3+, 2, 1

Homozygous | Heterozygous

True NOC: 2, 3, 4, 5 — Proportion of Loci

No. cells in cluster (1 → 15)

## D

RFU: 0-3000 — 55%, 7%, 38%
RFU: 3001-10000 — 93%, 3%, 4%
RFU: 10001-22000 — 99%, 1%, 0.1%
RFU: 22001-63788 — 99%, 1%

No. loci vs No. Incorrect Genotypes in credible set (0, 10, 20+)

(A) Histogram showing the posterior probability associated with the true genotype. Pie charts illustrate the proportion of clusters where the maximum probability aligns with the true genotype across NoC. **(B)** Scatterplots showing the probability that the list of candidate genotypes in the credible set contains the true genotype plotted against the sum of RFU intensity in a cluster per locus separated by NoC. **(C)** Mosaic plots showing the proportion of loci for which there were 1, 2, or 3+ genotypes included in the credible set separated by number of cells in a cluster, which ranged from 1-15 and by whether the true genotype is homozygous or heterozygous. **(D)** Bar chart showing the number of incorrect genotypes included in the credible set separated by NoC and total intensity per locus.

**Discussion:** Our method consistently performs well and includes the true genotype in the set of candidate genotypes that explain the evidence regardless of NoC, number of cells in a cluster and peak heights. Low information content is the only factor that increased the number of candidate genotypes. Low information content is primarily associated with clusters consisting of only one cell.

**Conclusion:** Our findings demonstrate the legitimacy of single-cell data for investigative genetics. We observed consistent inclusion of the true genotype regardless of NoC, number of cells and peak height using a threshold of 0.998 to define the genotype set associated with that posterior probability.