# ember❂ entropy metrics for biological exploration
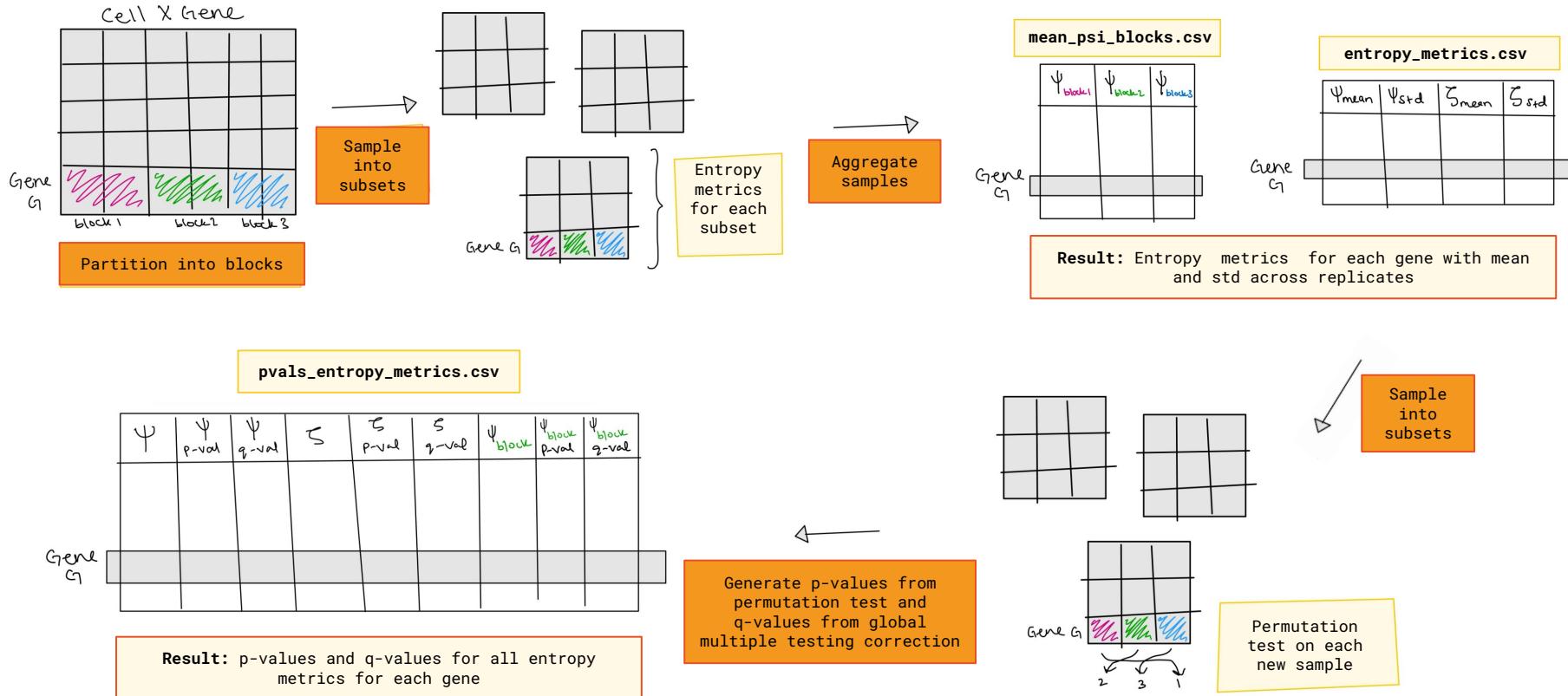
**Nikhila P. Swarna**
**October 9 2025**

**Pachter Lab**
**Division of Biology and Bioengineering**
**California Institute of Technology**

# light_ember workflow

**Cell X Gene**

Partition into blocks

block 1   block 2   block 3

Gene G

Sample into subsets

Entropy metrics for each subset

Gene G

Aggregate samples

**mean_psi_blocks.csv**

$\Psi$ block1   $\Psi$ block2   $\Psi$ block3

Gene G

**entropy_metrics.csv**

| $\Psi_{mean}$ | $\Psi_{std}$ | $\zeta_{mean}$ | $\zeta_{std}$ |

Gene G

**Result:** Entropy metrics for each gene with mean and std across replicates

**pvals_entropy_metrics.csv**

| $\Psi$ | $\Psi$ p-val | $\Psi$ q-val | $\zeta$ | $\zeta$ p-val | $\zeta$ q-val | $\Psi$ block p-val | $\Psi$ block p-val | $\Psi$ block q-val |

Gene G

**Result:** p-values and q-values for all entropy metrics for each gene

Generate p-values from permutation test and q-values from global multiple testing correction

Sample into subsets

Gene G

Permutation test on each new sample

2   3   1

# ember ❈ **e**ntropy **m**etrics for **b**iological **e**xplo**r**ation

**light_ember** - one stop shop for generating entropy metrics and p-values

**INPUT**
- **h5ad_dir** (path to adata)
- **partition_label** (col in adata.obs)
- **save_dir** (path to save results)
- sampling= True
  - sample_id_col=None
  - category_col=None
  - condition_col=None
  - num_draws=100
  - save_draws = False
  - seed = 42
- partition_pvals=True
- block_pvals=False
  - block_label=None
  - n_pval_iterations=1000
- n_cpus=1 (for parallel processing of sampling)

**light_ember**

**OUTPUT**
- csv file in save_dir with all entropy metrics
- csv file in Psi_block_df folder with psi block
  - Separate file for pvals
  - Separate files for each partition
  - Alternate file names depending on sampling on or off.

# ember ❖ **e**ntropy **m**etrics for **b**iological **e**xplo**r**ation

**generate_pvals** - manual access to generate pvals after initial investigation using **light_ember**

**INPUT**
- **h5ad_dir** (path to adata)
- **partition_label** (col in adata.obs)
- **entropy_metrics_dir** (path to light_ember output files)
- **save_dir** (path to save results)
- **sample_id_col**
- **category_col**
- **condition_col**
- block_label=None
- seed = 42
- n_iterations=1000
- n_cpus=1

**generate_pvals**

**OUTPUT**
- csv file in save_dir with entropy metrics and corresponding p-values and FDR q-values
  - Separate files for each partition

# ember 🔴 **e**ntropy **m**etrics for **b**iological **e**xplo**r**ation

**highly_specific_to_partition**

*INPUT*
- **partition_label**
- **pvals_dir**
- **save_dir**
- psi_thresh
- zeta_thresh
- q_thresh

*OUTPUT*
CSV file with genes highly specific to a partition
- Ordered from most to least specific

**top_genes**

**highly_specific_to_block**

*INPUT*
- **partition_label**
- **block_label**
- **pvals_dir**
- **save_dir**
- psi_thresh
- zeta_thresh
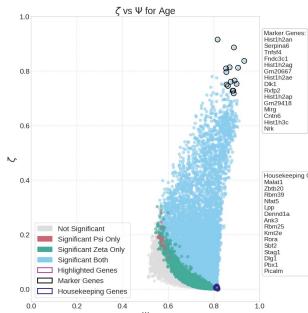- q_thresh

*OUTPUT*
CSV file with **marker genes** for a block
- Ordered from most to least specific

**non_specific_to_partition**

*INPUT*
- **partition_label**
- **pvals_dir**
- **save_dir**
- psi_thresh
- zeta_thresh
- q_thresh

*OUTPUT*
CSV file with non-specific **housekeeping genes** to a partition
- Ordered from most to least specific

# ember ✷ **e**ntropy **m**etrics for **b**iological **e**xplo**r**ation

## plot_partition_specificity

*INPUT*
- **partition_label**
- **pvals_dir**
- **save_dir**
- highlight_genes
- fontsize
- color_palette

## plot_psi_blocks

*INPUT*
- **gene_name**
- **partition_label**
- **psi_block_df_dir**
- **save_dir**
- fontsize

## plots

## plot_sample_counts

*INPUT*
- **h5ad_dir**
- **save_dir**
- **sample_id_col**
- **category_col**
- **condition_col**
- fontsize

## plot_block_specificity

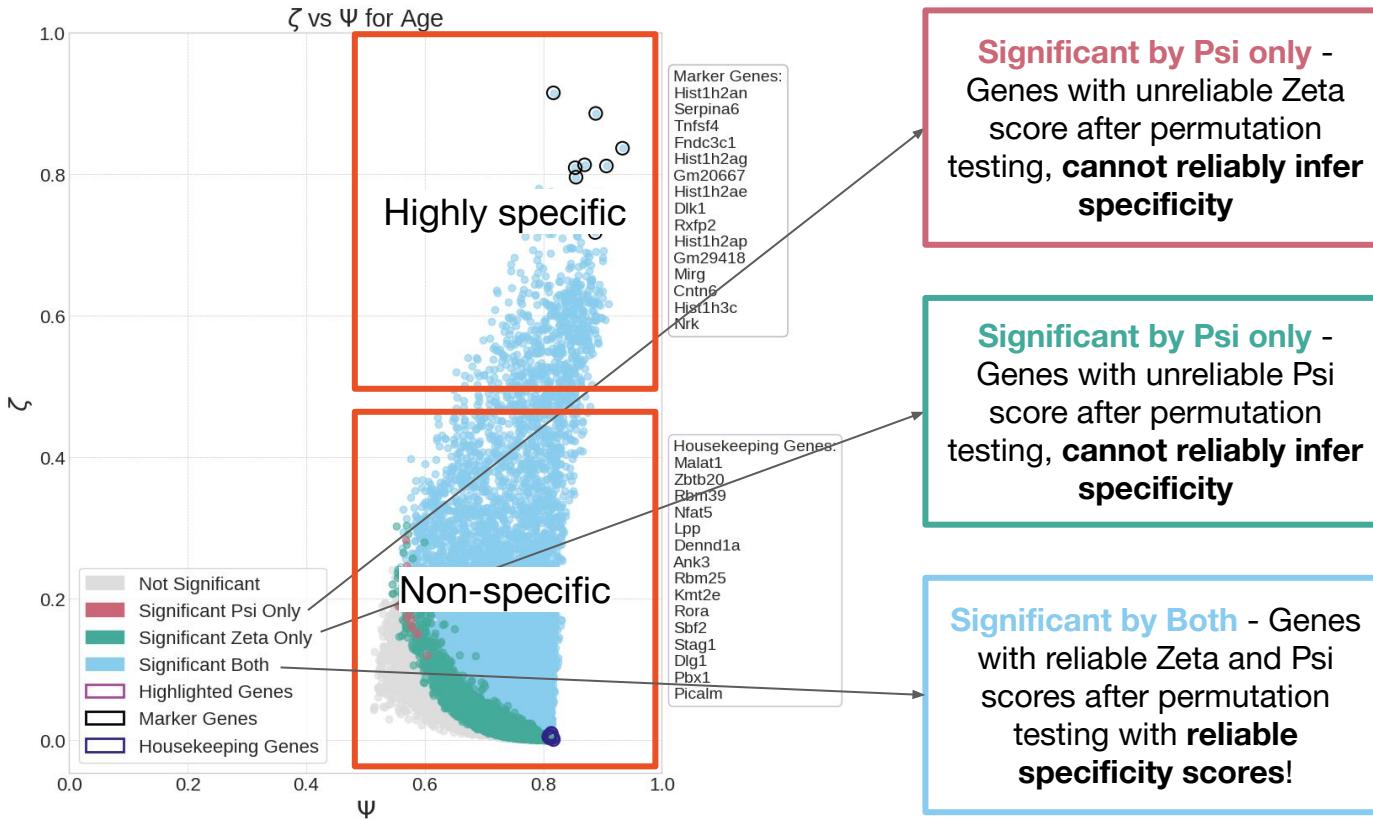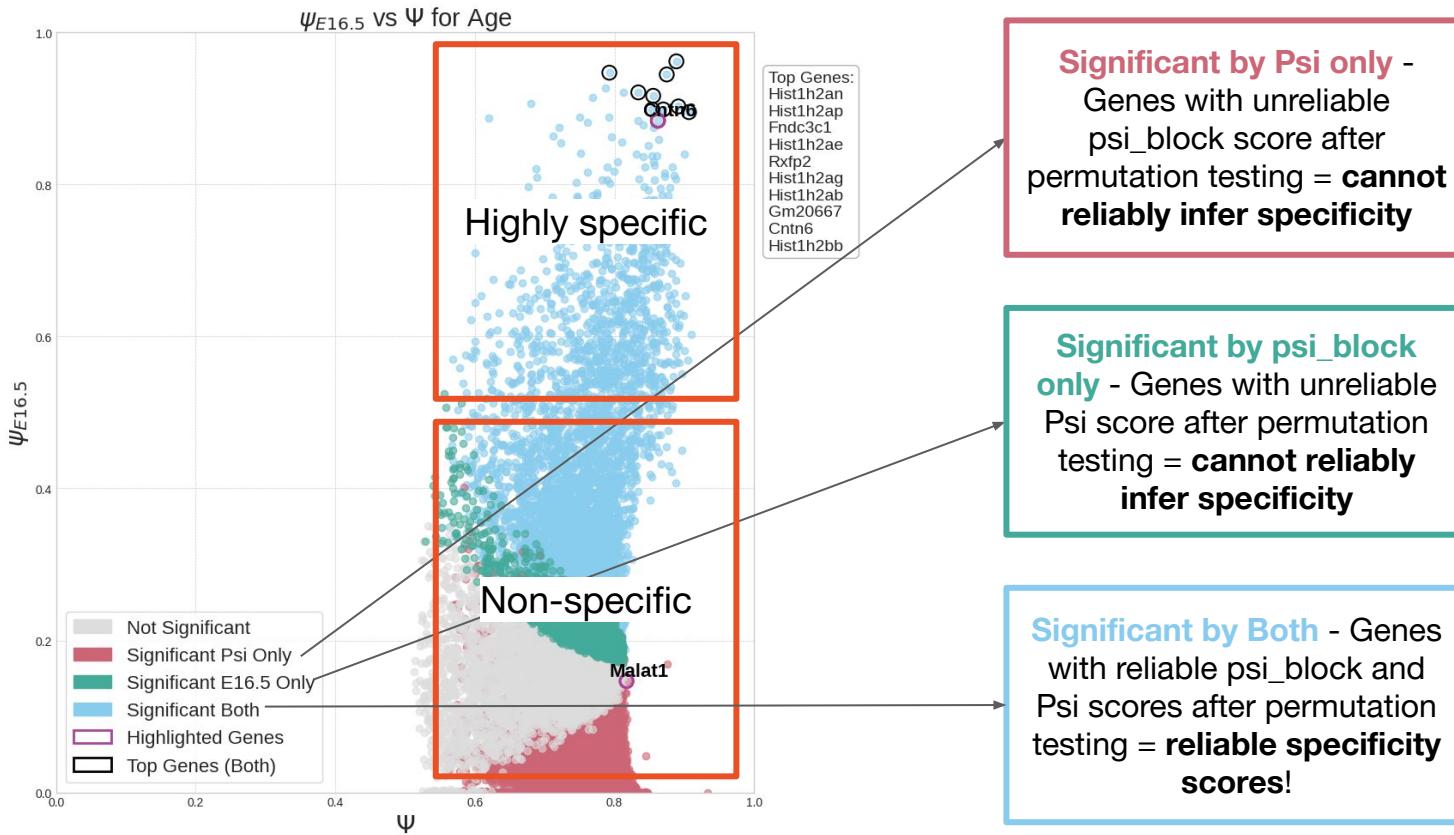*INPUT*
- **partition_label**
- **block_label**
- **pvals_dir**
- **save_dir**
- highlight_genes
- fontsize
- color_palette

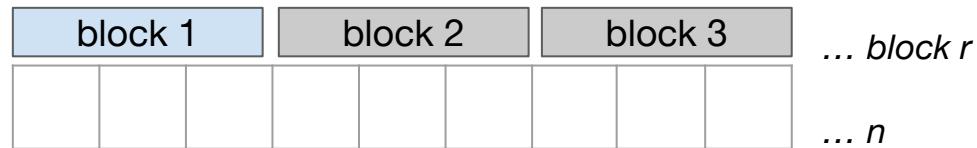# How to interpret `plot_partition_specificity`

# How to interpret `plot_block_specificity`



$\psi_{E16.5}$ vs $\Psi$ for Age

Highly specific

Non-specific

Top Genes:
Hist1h2an
Hist1h2ap
Fndc3c1
Hist1h2ae
Rxfp2
Hist1h2ag
Hist1h2ab
Gm20667
Cntn6
Hist1h2bb

Cntn6

Malat1

Not Significant
Significant Psi Only
Significant E16.5 Only
Significant Both
Highlighted Genes
Top Genes (Both)

$\Psi$

**Significant by Psi only** - Genes with unreliable psi_block score after permutation testing = **cannot reliably infer specificity**

**Significant by psi_block only** - Genes with unreliable Psi score after permutation testing = **cannot reliably infer specificity**

**Significant by Both** - Genes with reliable psi_block and Psi scores after permutation testing = **reliable specificity scores**!

# Defining entropy metrics for biological exploration

For a given gene in a count matrix that can be partitioned into r blocks (based on sex, strain, cell type, tissue, etc), we introduce **3 measures of specificity**:

| block 1 | block 2 | block 3 | *… block r* |
|---------|---------|---------|-------------|
| | | | *… n* |

- Psi ($\Psi$)

- Psi$_{block}$ ($\Psi_{block}$)

- Zeta ($\zeta$)

| $\Psi$ | $\Psi_{block}$ | $\zeta$ |
|--------|----------------|---------|
| Fraction of information explained by partitioning | Specificity to a block | Specificity to a partition |

# Ψ- Information Fraction by Partition

The fraction of information explained by using a
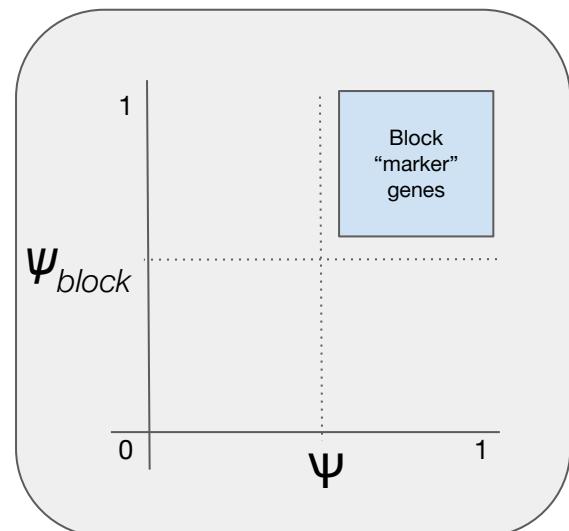particular partition on gene $g$'s counts is given by:

$$\Psi = \frac{E_W}{E_T} = 1 - \frac{E_B}{E_T}$$

The Specificity of Information to Block, denoted by
$\psi_{block}$, is the contribution of each block to $\Psi$:

$$1 = \boxed{\begin{array}{c} \psi_1 \\ \frac{p_1 E_1}{E_W} \end{array}} + \boxed{\begin{array}{c} \psi_2 \\ \frac{p_2 E_2}{E_W} \end{array}} + \cdots + \boxed{\begin{array}{c} \psi_r \\ \frac{p_r E_r}{E_W} \end{array}}$$

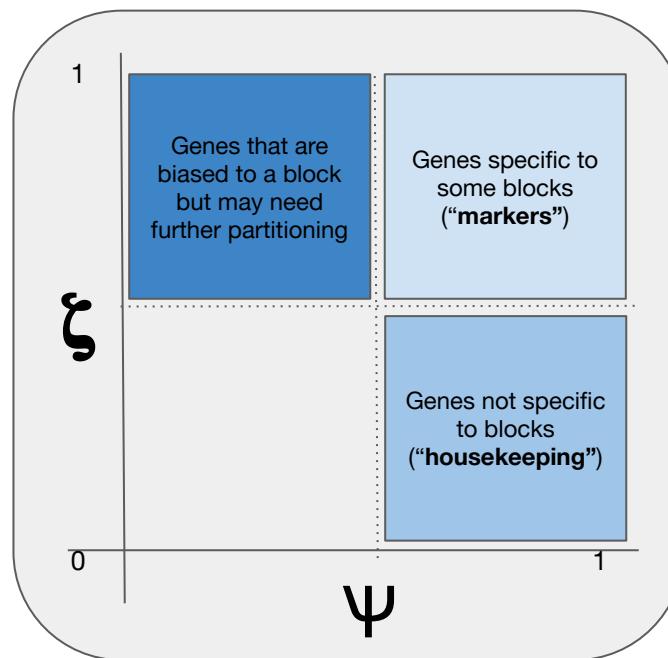| block 1 | block 2 | block 3 |
|---------|---------|---------|

... *block r*

... *n*

# ζ - Specificity of Information to Partition

The specificity of information to a partition ( ζ ) is given by:

$$\zeta = 1 - \frac{H(\psi_{\text{blocks}})}{\log r}$$



Comparison of the SIB distribution to the uniform distribution

# How to select category and condition

Category - Mouse strain
Condition - Sex

| Sample | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|--------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|
| Mouse strain | Strain A | | | | Strain B | | | | Strain C | | | | Strain D | | | |
| Sex | Male | | Female | | Male | | Female | | Male | | Female | | Male | | Female | |

**Number of unique draws**
= (Number of replicates per category-condition group)^(Number of category-condition groups)
= 2^8 = **256**

One example draw:

| Sample | 1 | 3 | 5 | 7 | 9 | 11 | 13 | 15 |
|--------|---|---|---|---|---|----|----|----|
| Mouse strain | Strain A | | Strain B | | Strain C | | Strain D | |
| Sex | Male | Female | Male | Female | Male | Female | Male | Female |

# How to select category and condition

Category - Cell line
Condition - Gene perturbation

| Sample | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cell line | Cell line A | | | | | | Cell line B | | | | | |
| Gene perturbation | Wildtype | | Overexpression | | Knockout | | Wildtype | | Overexpression | | Knockout | |

**Number of unique draws**
= (Number of replicates per category-condition group)^(Number of category-condition groups)
= 2^6 = **64**

One example draw:

| Sample | 1 | 3 | 5 | 7 | 9 | 11 |
|---|---|---|---|---|---|---|
| Cell line | Cell line A | | | Cell line B | | |
| Gene perturbation | WT | OE | KO | WT | OE | KO |

# How to select category and condition

Category - Mouse strain
Condition - Sex

| Sample | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mouse strain | Strain A | | | | Strain B | | | | Strain C | | | | | Strain D | | |
| Sex | Male | Female | | | Male | | Female | | Male | | | Female | | Male | | Female |

**Number of unique draws**
= ∏ (Number of replicates per category-condition group)
= 1*3*2*2*3*2*2*1 = **144**

One example draw:

| Sample | 1 | 2 | 5 | 7 | 9 | 12 | 14 | 16 |
|---|---|---|---|---|---|---|---|---|
| Mouse strain | Strain A | | Strain B | | Strain C | | Strain D | |
| Sex | Male | Female | Male | Female | Male | Female | Male | Female |

Note: #1 and #16 will appear in every draw since there are no replicates for these groups