# Duplicate Cluster Analysis
## (50 CC + 50 ES files sample)
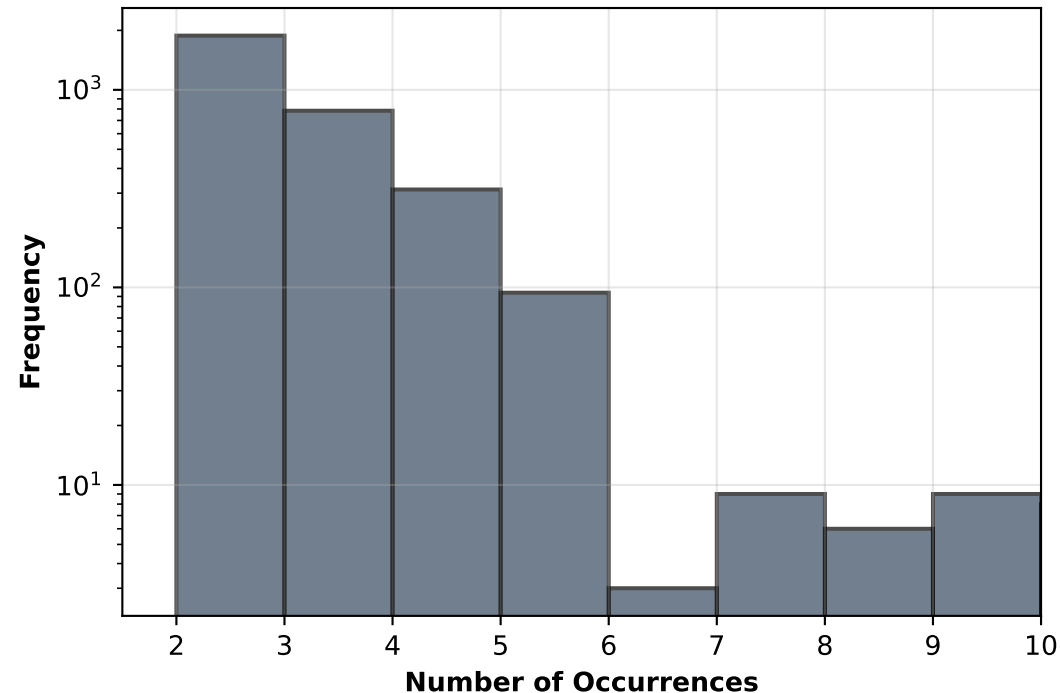
### Unique vs Duplicated Fingerprints



Duplicated
3136
(26.8%)

Unique
8550
(73.2%)

### Duplicate Types



### Duplication Multiplicity Distribution



```
SUMMARY STATISTICS

Sample size: 50 CC + 50 ES files

Total clusters: 22,516
Unique fingerprints: 11,686

Clusters appearing once: 8,550
Clusters appearing multiple times: 7,613

DUPLICATE RATE: 33.81%

Duplicate breakdown:
  • CC-CC only: 1,200 groups
  • ES-ES only: 1,400 groups
  • CC-ES mixed: 536 groups

Implication:
~1 in 3 clusters is a duplicate
```
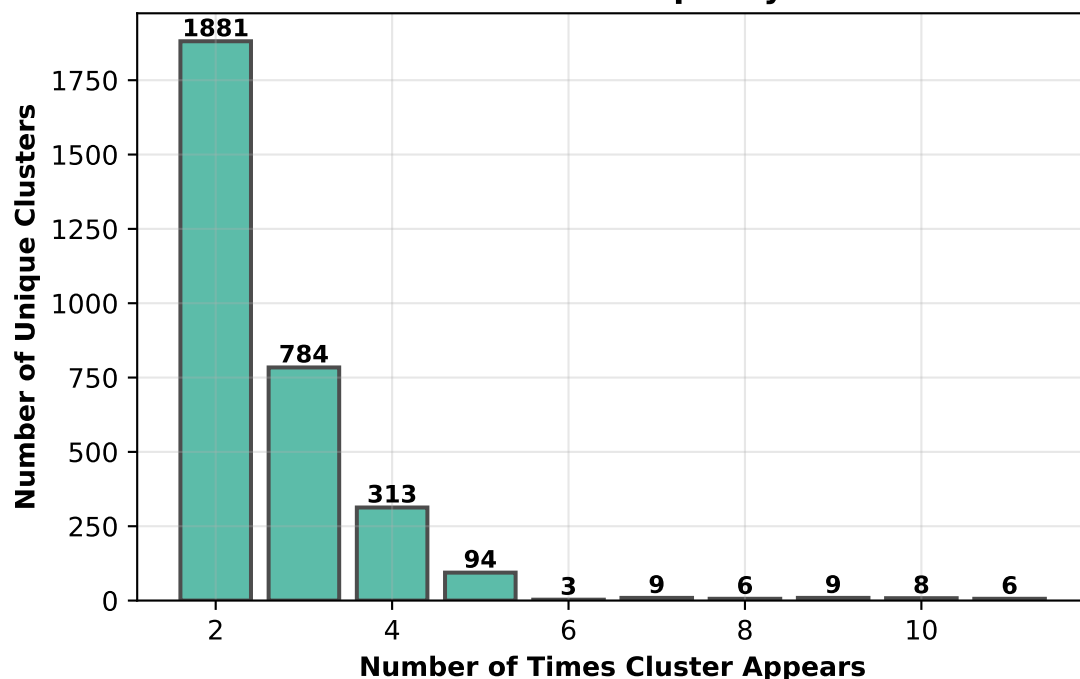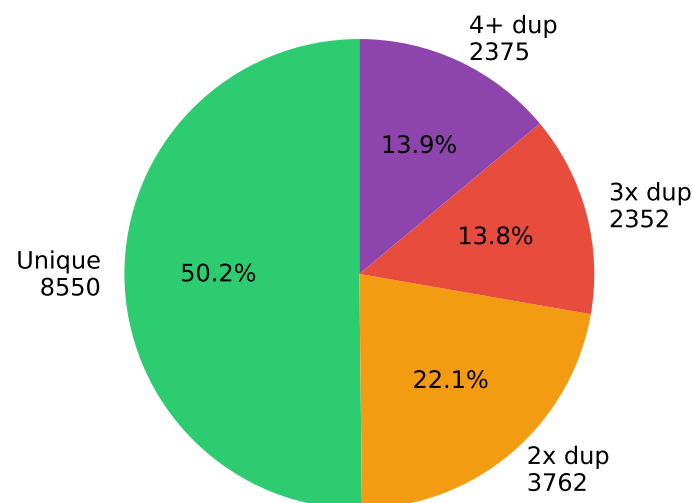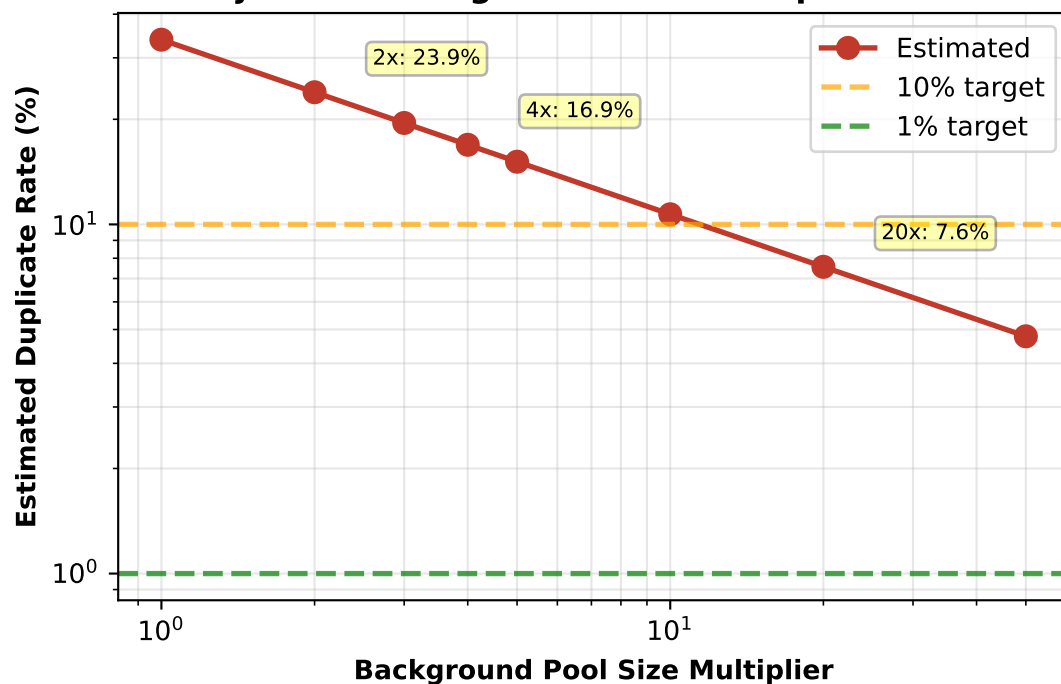
# Duplication Patterns and Projections

## Cluster Multiplicity



## Total Clusters by Duplication



Unique 8550 — 50.2%
2x dup 3762 — 22.1%
3x dup 2352 — 13.8%
4+ dup 2375 — 13.9%

## Projection: Background Size vs Duplicate Rate



2x: 23.9%
4x: 16.9%
20x: 7.6%

Legend:
- Estimated
- 10% target
- 1% target

```
RECOMMENDATIONS

Current status:
 • Duplicate rate: 33.8%
 • ~1 in 3 clusters is duplicated

To reduce to <10%:
 • Need ~11x larger background pool
 • From ~100 files to ~1,100 files

To reduce to <1%:
 • Need ~1,100x larger background pool
 • Practically challenging

Impact assessment:
 ✓ Less critical than file-level leakage
 ✓ Signal portions are different
 ✓ Model sees clusters in varied contexts

 ⚠ May cause overfitting to specific
   background patterns

Best practices:
 1. Monitor validation performance
 2. Use strong regularization
 3. File-level train/val/test split (✓ done)
 4. Consider data augmentation
```