

We Mean Business:

# Turning Validation Metrics into Defensible Decisions

May 5, 2026



# Speakers

---



**Stephanie Clerkin**

---

Director of Litigation Support  
Korein Tillery



**Matt Mahon**

---

Vice President of Client Solutions  
Level Legal



# Agenda

---

- Foundation
- Key Concepts
- Pre-Process Considerations
- Case Studies
- Q&A

# Hot Takes

---



**Technology won't slow you down. Human sign-off will.**



**We need to stop arguing the mechanics of document review and shift the focus to agreeing on a defensible validation protocol.**

---

# Foundation





---

# Disclaimer

---

- A singular methodology for validating the results of a technology-assisted review workflow (including GenAI workflows) **doesn't exist**.
- You need to consider:
  - Richness
  - Proportionality
  - Edge Cases
  - Risk Tolerance
  - Custodians
  - Data Sources
  - File Types
  - Foreign Language

# The Confusion Matrix

		Actual	
		Positive	Negative
Predicted	Positive	 <b>True Positive</b> Relevant documents tagged as relevant	 <b>False Positive</b> Non-relevant documents tagged as relevant
	Negative	 <b>False Negative</b> Relevant documents tagged as non-relevant	 <b>True Negative</b> Non-relevant documents tagged as non-relevant

---

# Key Concepts

---

# Richness / Prevalence



**Definition:** The percent of a population that has a specific characteristic, such as responsiveness.



- Measures the % of relevant documents within a document population.



- Informs review strategy and workflow decisions.

# Richness / Prevalence

$$\textit{Richness} = \frac{\textit{True Positive} + \textit{False Negative}}{\textit{True Positive} + \textit{False Negative} + \textit{True Negative} + \textit{False Positive}}$$

# Recall



**Definition:** The number of documents retrieved from a search divided by all of the responsive documents in a collection.



- Measure of completeness of your review.



- *Remember:* It's always an estimate unless the full set of documents is reviewed.

# Recall

---

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

# Precision

---



**Definition:** The number of true positives retrieved from a search divided by the total number of results returned.



- Measure of the accuracy of your review.

# Precision

---

$$\textit{Precision} = \frac{\textit{True Positive}}{\textit{True Positive} + \textit{False Negative}}$$

# F-Measure / F Score / F1 Score



**Definition:** A measure of a search's accuracy calculated by using precision and recall.



- A weighted average that balances recall and precision.



- In eDiscovery, recall is often more important.

# F-Measure / F Score / F1 Score

$$\begin{array}{l} F\text{-Measure} \\ F\text{ Score} \\ F1\text{ Score} \end{array} = 2x \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

# Error



**Definition:** The fraction of documents that are incorrectly coded by a search or review effort.

## Keep in mind:



- False positives and false negatives are not equally harmful.



- The error rate itself doesn't tell you the type of error.



- Context matters (e.g., privilege review vs. responsiveness review).

# Error

---

$$\textit{Error} = \frac{\textit{False Positive} + \textit{False Negative}}{\textit{True Positive} + \textit{False Negative} + \textit{True Negative} + \textit{False Positive}}$$

# Random Sampling



**Definition:** The process of selecting data from a population with no bias or input from the person performing the sampling, in which each item has an equal chance of being selected as any other item.

# Stratified Sampling



**Definition:** A method of data sampling where data is initially divided into subgroups (e.g., by age range or a geographic criteria) or strata, and then each group is sampled in order to ensure that each subgroup is properly represented.

## Useful for:



- Different custodians.



- Different document types or sources.



- Uneven document distributions.

# Confidence Level



**Definition:** The percentage of samples for which the results are expected to correctly describe a population parameter within a provided confidence interval.



- Be careful: Percentages can hide large numerical ranges.

# Confidence Interval / Margin of Error



**Definition:** The range of values that is likely to contain the true parameter for a population to the specified confidence level.



- Be careful: The margin of error range grows significantly with large populations.

# Point Estimate



**Definition:** The result of a sample that estimates prevalence in the specific population being sampled.



- Used to estimate collection richness and the total number of relevant documents.



- Be careful: Percentages alone can be misleading.
  - Always convert percentages to document counts.

# Elusion



**Definition:** The percentage of documents of a search's null set that were missed by the search, usually determined with review of a random sample of the null set.



- While widely used, it still has issues:
  - Can be difficult to interpret without total number of relevant documents.
  - The rate can imply different outcomes depending on context.



- Don't rely on elusion alone.

---

# Pre-Process Considerations

---

# Before You Start



## Questions you should ask to guide your methodology:

- What metrics are most important and why?
- What processes address the impact of those metrics?
- Have opposing counsel or the court set expectations around validation methodology or recall targets?

# GenAI Considerations

---

- Will you use GenAI in your workflows?
- How do you ensure that your use of GenAI stands up to mathematical and legal scrutiny?
- How will you validate the GenAI output?
- Does your GenAI workflow require court approval or opposing counsel agreement?

# Pre-Process Checklist

---

- Agree on review platform and any GenAI / TAR tools to be used.
- Define scope: custodians, date ranges, data sources, and file types.
- Is there an expected richness of the document population?
- Determine confidence level and margin of error for validation sampling.
- Define what constitutes a "responsive" document for this matter.
- Confirm opposing counsel / court agreement on methodology where required.

---

# Case Studies

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

“

*The court ordered the parties to combine three subsamples of 3,000 documents into a single validation sample. The sample was to be reviewed and coded by a subject matter expert, and the producing party was to prepare a table listing each of the 3,000 documents. The requesting party and the special master were to be provided with the table, a copy of each responsive, non-privileged document in the validation sample, and the statistics and recall estimate. The recall estimate was to be computed to inform the decision-making process.*

”

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---



Subcollection  
C(1)

*Documents identified by the review as responsive to at least one Request for Production, including any privileged documents, but not including family members of responsive documents, unless those family members are deemed to be responsive in their own right.*

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---



### Subcollection C(2)

*Documents coded as non-responsive by a human reviewer, regardless of how the documents were selected for review (e.g., by TAR, manual review, or otherwise).*

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---



### Subcollection C(3)

*Documents excluded from manual review as the result of a TAR process. If the review process involved only manual review and no TAR, the Collection will not include Subcollection C(3).*

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Stratified Random Samples

- **Subsample D(1):** 500 documents selected at random from Subcollection C(1).
- **Subsample D(2):** 500 documents selected at random from Subcollection C(2), if TAR was used, otherwise 2,500 documents selected at random from Subcollection C(2), if manual review was used.
- **Subsample D(3):** 2,000 documents selected at random from Subcollection 1(c) if TAR was used.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

Were the Subsample sizes selected at random?

**Sample Size Calculator**

Forward (Required Sample Size) | **Reverse (Achieved Margin of Error)**

Population Size

Sample Size (n)

Confidence Level

Response Distribution (%)

**Calculate** | Reset | Export PDF

**Results (Reverse)**

Achieved Margin of Error: 4.38%

Computed using finite population correction when population is provided.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

Were the Subsample sizes selected at random?

**Sample Size Calculator** 🌙 / ☀️

Forward (Required Sample Size) **Reverse (Achieved Margin of Error)**

Population Size ⓘ 1000000000000

Sample Size (n) ⓘ 500

Confidence Level ⓘ 99%

Response Distribution (%) ⓘ 50

**Calculate** **Reset** **Export PDF**

---

**Results (Reverse)**

Achieved Margin of Error: 5.76%

Computed using finite population correction when population is provided.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

Were the Subsample sizes selected at random?

The screenshot shows a 'Sample Size Calculator' interface. At the top, there are two tabs: 'Forward (Required Sample Size)' and 'Reverse (Achieved Margin of Error)'. The 'Reverse' tab is selected. Below the tabs, there are four input fields: 'Population Size' (1000000000000), 'Sample Size (n)' (2000), 'Confidence Level' (99%), and 'Response Distribution (%)' (50%). Below these fields are three buttons: 'Calculate', 'Reset', and 'Export PDF'. At the bottom, there is a 'Results (Reverse)' section showing 'Achieved Margin of Error: 2.88%'. Orange arrows point to the 'Sample Size (n)' field, the 'Confidence Level' field, and the 'Achieved Margin of Error' result.

**Sample Size Calculator**

Forward (Required Sample Size) | **Reverse (Achieved Margin of Error)**

Population Size (0) | Sample Size (n) (0)

1000000000000 | 2000

Confidence Level (0) | Response Distribution (%) (0)

99% | 50

**Calculate** | Reset | **Export PDF**

**Results (Reverse)**

Achieved Margin of Error: 2.88%

Computed using finite population correction when population is provided.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Validation Sample: Ensure the Review Is Blind

- The sample of 3,000 documents comprised of the documents from Subsamples D(1), D(2), and, if TAR was used, D(3), shall be combined into a single Validation Sample, with no indication of the Subcollection from which the documents were derived or how they were previously coded. The Validation Sample shall be reviewed and coded by a subject matter expert (“SME”) who is knowledgeable about the subject matter of the litigation.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Recall Estimation Method

- The number of responsive documents found  $\approx$  the size of Subcollection C(1)  $\times$  the number of responsive docs in Subsample D(1)  $\div$  500.
- The number of responsive documents coded incorrectly  $\approx$  the size of Subcollection C(2)  $\times$  the number of responsive documents in Subsample D(2)  $\div$  500.
- The number of responsive documents not reviewed  $\approx$  size of Subcollection C(3)  $\times$  the number of responsive documents in Subsample D(3)  $\div$  2,000.
- **Estimated recall**  $\approx$  the number of responsive documents found  $\div$  (the number of responsive documents found + the number of responsive documents coded incorrectly + the number of responsive documents not reviewed).

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Hypothetical Calculation

- Assume the number of documents in the collection (defined as including all documents identified for review for responsiveness and/or privilege following the application of keywords or other culling criteria) = 100,000 documents

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Hypothetical Calculation

- Collection = 100,000 documents
  - Subcollection C(1) = 30,000 documents
  - Subcollection C(2) = 20,000 documents
  - Subcollection C(3) = 50,000 documents

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Hypothetical Calculation

- True Positive = 30,000 x 450
  - Subcollection C(1) = 30,000 documents
    - Subsample D(1) = 500 documents
      - **# of responsive documents in D(1) = 450 = True Positive**
  - Subcollection C(2) = 20,000 documents
    - Subsample D(2) = 500 documents
      - **# of responsive documents in D(2) = 25 = False Negative<sub>D(2)</sub>**
  - Subcollection C(3) = 50,000 documents
    - Subsample D(3) = 2,000 documents
      - **# of responsive documents in D(3) = 200 = False Negative<sub>D(3)</sub>**

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Hypothetical Calculation

- The # of responsive documents found  $\approx$   
# of documents in C(1) x # of responsive documents in D(1)  $\div$  500  
 $\Rightarrow 30,000 \times 450 \div 500 = 27,000$  documents = True Positive
- The # of responsive documents coded incorrectly  $\approx$   
# of documents in C(2) x # of responsive documents in D(2)  $\div$  500  
 $\Rightarrow 20,000 \times 25 \div 500 = 1,000$  documents = False Negative<sub>D(2)</sub>
- The # of responsive documents not reviewed  $\approx$   
# of documents in C(3) x # of responsive documents in D(3)  $\div$  500  
 $\Rightarrow 50,000 \times 200 \div 2,000 = 5,000$  documents = False Negative<sub>D(3)</sub>

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Hypothetical Calculation: Precision and Recall

- Recall = True Positive ÷ (True Positive + False Negative)  
=> Recall = 27,000 ÷ (27,000 + 1,000 + 5,000) = 81.82%
- Precision = True Positive ÷ (True Positive + False Positive)  
=> Precision = 27,000 ÷ (27,000 + 3,000) = 90%

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

### Hypothetical Calculation: Confidence Interval / Margin of Error

**Sample Size Calculator**

Forward (Required Sample Size)    **Reverse (Achieved Margin of Error)**

Population Size ⓘ    Sample Size (n) ⓘ

100000    3000

Confidence Level ⓘ    Response Distribution (%) ⓘ

99%    50

**Calculate**    Reset    **Export PDF**

---

**Results (Reverse)**

Achieved Margin of Error: 2.32%

Computed using finite population correction when population is provided.

# Case Study #1:

## In re Broiler Chicken Antitrust Litigation

---

### Did We Get The Validation We Need?

- Evaluate the Relevance of the False Negatives
- Calculate the Recall Range
  - Point Estimate is 33,000 documents
  - Margin of Error (Confidence Interval) = +/- 2.32%
  - Estimated # of Relevant Documents =  $33,000 \pm 2.32\% \times 100,000$   
=>  $33,000 \pm 2,320$   
=> Estimate # of Relevant Documents = [30,680, 35,320]
  - Recall Range =  $[27,000 \div 35,320, 27,000 \div 30,680]$   
=> Recall Range = [76.44%, 88.01%]

# Case Study #2:

## Applying Broiler Chicken To GenAI Reviews

---

### Validation Framework: Subcollections

- **Subcollection C(1):** Documents coded as responsive by a human reviewer.
- **Subcollection C(2):** Documents coded as responsive by GenAI.
- **Subcollection C(3):** Documents coded as non-responsive by a human reviewer.
- **Subcollection C(4):** Documents coded as non-responsive by GenAI.
- **Subcollection C(5):** Documents in the review population excluded from manual and GenAI review.

# Case Study #2:

## Applying Broiler Chicken To GenAI Reviews

---

### Validation Framework: Stratified Random Samples

- **Subsample D(1):** 500 documents selected at random from Subcollection C(1).
- **Subsample D(2):** 500 documents selected at random from Subcollection C(2).
- **Subsample D(3):** 500 documents selected at random from Subcollection C(3).
- **Subsample D(4):** 500 documents selected at random from Subcollection C(4).
- **Subsample D(5):** 2,000 documents selected at random from Subcollection C(5).

# Case Study #2:

## Applying Broiler Chicken To GenAI Reviews

---

### Validation Framework: Other Subsample Considerations

- Key Custodian(s)
- Data Source(s)
- File Type(s)
- Date Range(s)
- Foreign Language
- Search Term Hits

# Case Study #3:

## Inbound Production

---

### Sampling and Validation For Documents Produced by Third Parties

- **800,000 documents received from a defendant's production:** TAR was applied instead of linear review for deposition prep and expert work.
- **Random + Stratified Sampling:** random for richness baseline; stratified by custodian.
- **3% Elusion:** within threshold; model accepted.
- **Worthless form Templates:** pulled into the model, caught by precision sampling, excluded without hurting recall.
- **Validation:** makes TAR an easier sell to case team and more defensible.

---

# Q&A

---



Level Legal is an eDiscovery, managed review, and consulting company that delights law firms, corporations, and government agencies through industry-best customer service. This concierge approach to outsourced legal services delivers peace of mind.

© Level Legal. All Rights Reserved.