



WHITE PAPER

Metrics that Matter

Van Halen, M&Ms and Measurement
in E-Discovery

Gareth Evans, Partner, Gibson, Dunn & Crutcher

David Grant, Senior Managing Director, FTI Consulting

GIBSON DUNN



F T I
CONSULTING

TECHNOLOGY

What do Van Halen, M&Ms and metrics in e-discovery have to do with one another? More than you might think.

During Van Halen's heyday, the band was frequently on tour and had a rather elaborate stage show. It included heavy sound and light equipment, as well as harnesses so that David Lee Roth could fly over the audience.

In each city, the band was dependent upon different local companies to unload and set up the equipment. High-quality work was paramount. If the local production company overlooked something, it could mean serious injury for band members and, potentially, the audience.

Other than following each roadie as they did their tasks, how could the band ensure their instructions were followed? Cleverly, the band devised a strategy that would enable them to quickly assess whether the local company was likely following instructions.

Buried deep in the band's long list of instructions in its 53-page tour rider, Van Halen required a bowl of M&Ms in the dressing room, but with one important caveat: All of the brown M&Ms had to be removed. Failure to comply would result in the concert promoter's forfeiture of the entire show at full pay.

You might think of this as an example of 1970s rock band outlandishness. But it actually served an important practical function.

If Van Halen arrived on the day of the concert and found brown M&Ms in their dressing room, they knew that their instructions weren't followed closely. So they would conduct a thorough quality review of every step

of the stage setup. If the brown M&Ms were removed, as instructed, the band had greater comfort that their instructions had been followed, and they could focus on testing a few key pieces of equipment.

A great video of Diamond Dave—"at his most impish and charming," according to NPR—providing an official history of the ingenious safety measure is available [here](#).¹ And a copy of the tour rider is available [here](#).² Roth provided an example of a concert where the brown M&Ms had not been removed and, not coincidentally, an improperly installed stage caused nearly \$500,000 of property damage.

Just as the colored M&Ms provided a quick visual metric of quality, similar calculations can be made to test the effectiveness of search and review. With a few quick calculations, it's possible both to have effective statistical arguments for minimizing the burden of the discovery, and also to understand whether you are on the right track or need to revisit earlier steps for quality analysis.

Doing this is increasingly important as large document volumes are changing the approaches litigants must take if they wish to avoid being overwhelmed with e-discovery burden and expense. In short, using the right metrics can help ensure defensibility and cost-effectiveness.

Key Metrics for Search and Review

Key metrics for use in search and review include:

- Prevalence (also called "Richness")
- Recall
- Precision
- Depth for Recall
- Confidence Level
- Confidence Interval (*i.e.*, margin of error)

¹ <http://www.npr.org/blogs/therecord/2012/02/14/146880432/the-truth-about-van-halen-and-those-brown-m-ms>.

² <http://www.thesmokinggun.com/documents/crime/van-halens-legendary-mms-rider>.

These metrics all involve statistical sampling, which allows one to use a representative random sample to affordably draw reliable conclusions about the overall document population without the time and expense of reviewing the entire document set.

As a bonus, the size of the sample needed is not directly proportional to the population size, and so

generally requires review of a relatively small number of documents.³

We refer to “responsive” or “relevant” documents in the examples below, but the same metrics can be used with respect to other aspects of a document population, such as attorney-client privileged documents, key documents, foreign language documents and so on.

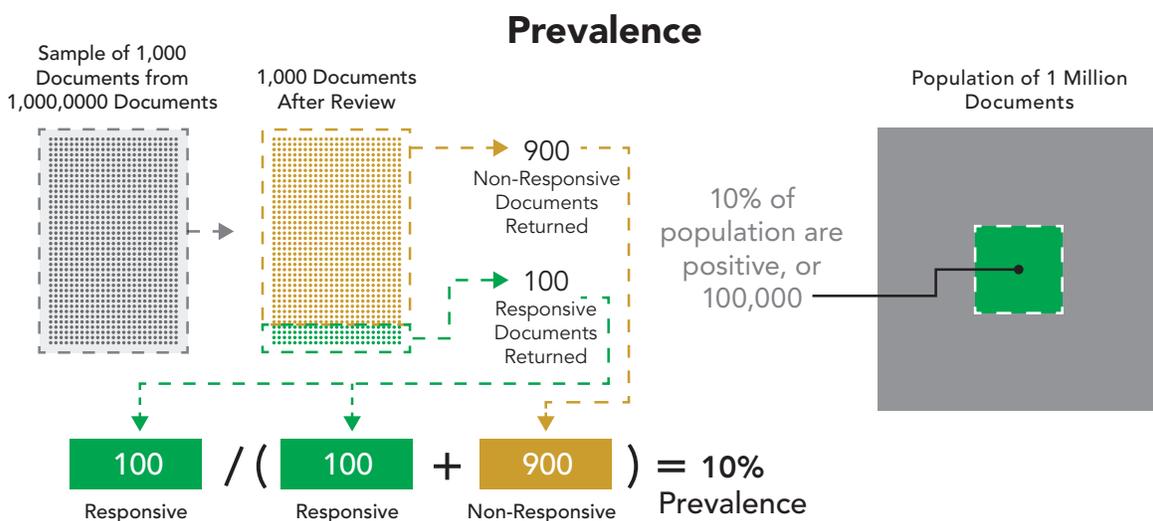
Prevalence / Richness – “How many are there?”

Imagine knowing in advance roughly how many responsive documents there will be—or how many foreign language documents, or privileged documents—so that you can ensure the right resources are in place. In other words, how many brown M&Ms can we expect to be in that big bowl? “Prevalence” (also known as “Richness”) will tell us.

Prevalence measures the proportion of the population that has a particular characteristic based on the proportion found in the sample. For example, if 10% of the sample is responsive, then we can project that 10% of the population from which the sample is drawn will be responsive (within a certain margin of error).

Or, to go back to the M&Ms example, if a statistical sample of M&Ms produces 10% brown M&Ms, we can project that 10% of the total M&Ms are brown.

Knowing Prevalence or Richness not only enables planning for the right resources and employing the best search and review strategy, but it also can help you establish your “goal” for retrieval processes (for example, how many responsive documents exist that your keywords, technology assisted review process or reviewers are trying to find).



³ The population from which the sample is drawn should exclude the documents used to develop the search methodology – such as documents used to train a predictive coding model. Otherwise the sample may not properly represent how the process would work against documents not seen before by the process.

Recall – “How many of (the Brown M&Ms) am I finding?”

Imagine knowing with a high level of confidence how well your search and review process is working to find responsive (or key or privileged) documents. In other words, what proportion of the total number of brown M&Ms are we finding? “Recall” will tell us.

Recall is the percentage of the total responsive documents in a document population that a search or review process actually finds. It is probably the most important search and review metric. To get the Recall percentage, divide the number of responsive documents in a sample that a process finds by the actual total number of responsive documents in the sample.

Recall permits one to reliably evaluate the effectiveness of proposed search terms—*i.e.*, how well they perform in actually finding relevant and responsive documents. This information can be critical in negotiating search terms or ensuring those that you are using are defensible.

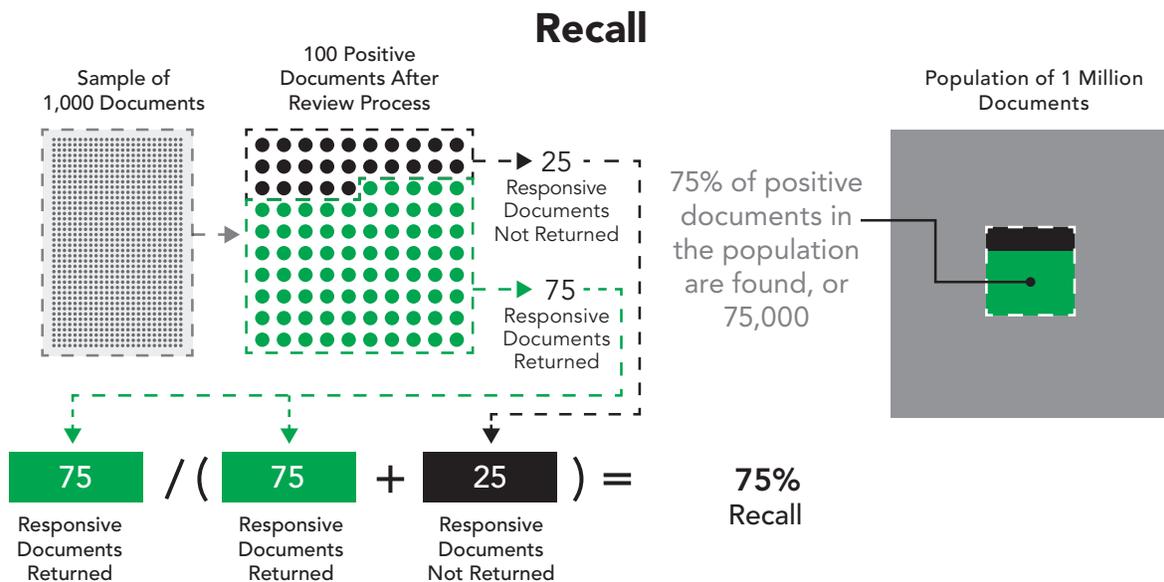
Recall is also the key metric used to target and evaluate performance of predictive coding protocols.

Predictive coding protocols often target a particular Recall rate, say 80%. A validation process involving the review of a statistical sample of the entire document set to be predicted on can confirm that the results of the predictive coding process achieved the targeted Recall.

Likewise, Recall can be used to evaluate the effectiveness of a combined search and review process, such as using search terms to initially cull the document set and then applying predictive coding to the search term “hits.”

Recall can also be used to report on reviewer effectiveness – identifying the proportion of documents coded positive by an authoritative reviewer that a first-pass reviewer also codes as positive, for instance.

As an example, if you have a very large bag of M&Ms, the package says it has 100 of each color, and you are searching for all of the brown M&Ms, and you find 80 of the M&Ms, that would mean your Recall on brown M&Ms is 80%.



Precision – “How much effort am I wasting to find (the brown M&Ms)?”

Imagine knowing how efficient your search and review process is—*i.e.*, what percentage of the documents it finds are “false positives.” In other words, when you grab some M&Ms from the bowl, trying to get only the brown ones, what portion are not brown?

Precision measures how well a search or review process yields only responsive documents. It is determined using statistical sampling methods similar to those used in determining Recall. The number of actually responsive documents identified by the search/retrieval process is divided by the total number of documents that the search/retrieval process retrieved. Many will be false positives, *i.e.*, “negatives” that the process suggested would be “positives.” For example, 5% Precision means 19 irrelevant or nonresponsive documents are retrieved for every 1 relevant or responsive document retrieved

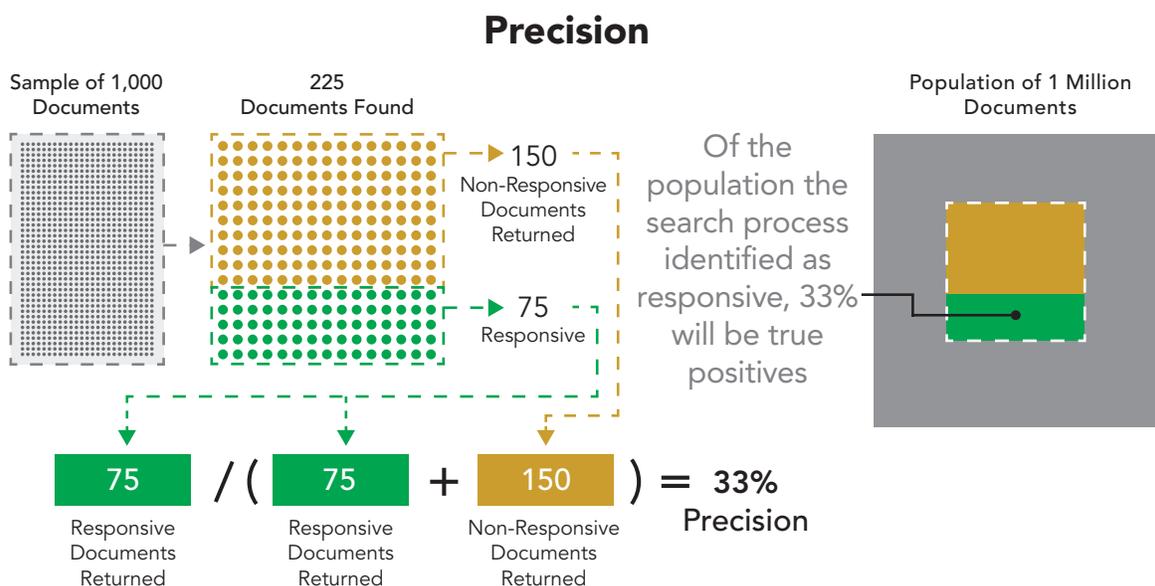
Having search terms or a predictive coding process that yields high levels of Precision is key to avoiding the inefficiencies and high costs of reviewing large numbers of irrelevant and nonresponsive documents.

Being able to measure Precision can be critical for a responding party in negotiating an acceptable set of search terms or a reasonable Recall level for predictive

coding. It can also be important in persuading the court that the opposing party’s proposed search terms or Recall level are overly broad. (Achieving higher levels of Recall is often done at the expense of lower Precision.)

In M&M terms, let’s say you have a bag of hundreds of assorted M&Ms and want to find the 100 brown M&Ms. If you were to pour out some of the M&Ms in the bag and get a total of 130 M&Ms, including all 100 brown M&Ms, that means that you have an additional 30 M&Ms beyond your goal of 100 brown M&Ms—*i.e.*, 30 “false positive” hits. That would be a Precision rate of 77% (*i.e.*, 100 brown M&Ms divided by the 130 total M&Ms poured out of the bag).

Of course, the “holy grail” of search and review is to have a process that yields both high Recall and high Precision. A carefully tested and calibrated set of search terms and an effective predictive coding process may achieve that goal. The two, however, are often inversely related—*i.e.*, achieving higher Recall may mean lower Precision and vice versa. This means having both statistics is useful for negotiations around proportionality. It allows quantification of the increased burden caused by higher levels of Recall.



Depth for Recall – “How much work is needed to find (the brown M&Ms)?”

At the end of the day, the crucial issue is how many documents you need to review. Depth for Recall measures the proportion of the document population that you must review using a particular search process to achieve a given Recall level. In other words, using our search process, how many M&Ms do I need to look at to confirm that I have found the brown M&Ms at my targeted level of Recall?

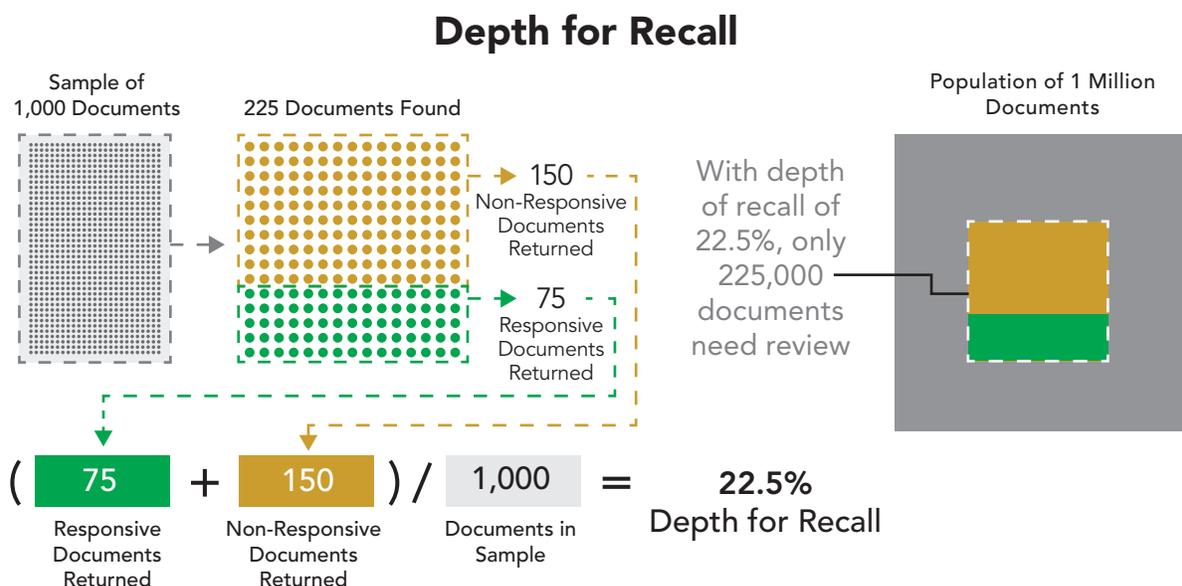
One way to obtain the Depth for Recall figure—*i.e.*, the fraction of the overall document set that you’ll need to review to achieve a targeted Recall level—is to multiply the Prevalence of responsive documents by the targeted Recall, then divide by the Precision of the search process at that Recall level. As discussed above, each of the metrics used in the Depth for Recall calculation can be obtained through sampling.

An alternative is to simply divide the number of documents in the sample marked responsive by the search or retrieval process by the total number of documents in the sample. If 15% of the documents in the sample would be marked responsive by the process, so would 15% of the documents in the population – so Depth for Recall is 15%.

To illustrate the former calculation, let’s say we have a population of 100,000 documents. If only 1% (or 1,000 documents) is projected to be responsive (*i.e.*, 1% Prevalence), then a target of 75% Recall using our search process means that we have to find 750 responsive documents to achieve the target. If the Precision of our search process is only 5%, how many documents will we need to review to find 750 responsive documents? Even though our search process yields only 5% Precision, the Depth for Recall metric (Prevalence x Recall divided by Precision) tells us that, using our search process, we’ll only need to review 15% (*i.e.*, 1% x 75% divided by 5%) of the population (*i.e.*, 15,000 documents) to find 750 responsive documents (*i.e.*, 75% Recall).

Depth for Recall is important in two different ways. First, Precision alone doesn’t actually indicate the effort required, because in cases of very low Prevalence even a low Precision can lead to a small proportion of the population requiring review.

Thus, in the example above, the 15% Depth for Recall number tells us that even though we have a very low Precision figure, our process still does a good job of



culling out 85% of the irrelevant or nonresponsive documents, leaving only 15% remaining for review.

Second, Depth for Recall enables comparison between various achievable Recall levels using predictive coding (or other search methods). It allows the additional effort required to achieve higher Recall levels to be

quantified, along with any benefits gained in terms of the number of additional responsive documents found. In other words, how many more documents will we need to review to achieve higher Recall levels? Depth for Recall thus enables analysis of the proportionality of the different Recall levels.

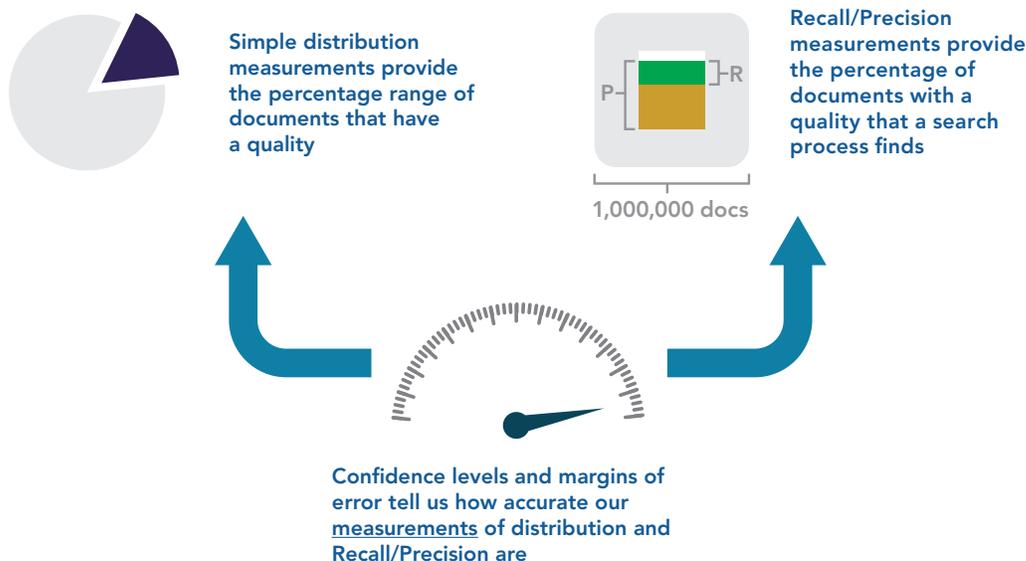
Confidence Level and Confidence Interval (Margin of Error) – “How accurate are my Prevalence / Recall / Precision measurements?” or “What’s the range of results I could achieve”?

Finally, imagine knowing how accurate your key metrics are. When reviewing only a sample from the population, the measurements are not 100% accurate (that could only be achieved, in theory, by review of every single document in the document population with perfect reviewers). With respect to Prevalence, Recall or Precision, for example, a 2,000 document sample may provide 99% confidence (the Confidence Level) with approximately a +/- 3% margin of error (the Confidence Interval), or better.

Generally, the accuracy of Prevalence, Recall and Precision measurements are related to the size of the sample used to generate the metric. So, if a 2,000 document sample shows that 50% of the documents are responsive, we could say we are 99% confident that between 47% and 53% of the documents in the population are

responsive. By contrast, a 42-document sample would have a +/- 20% margin of error at 99% Confidence—i.e., a much larger Confidence Interval (we could say we are 99% confident that between 30% and 70% of the documents in the population are responsive).

A few things to keep in mind about Confidence Levels and Confidence Intervals: First, a common mistake is to assume that the sample size for a given Confidence Level and Confidence Interval is proportional to the population size. In reality, however, the sample size required for a given Confidence Level and Confidence Interval is not significantly affected by the population size. In other words, the sample needed for a very large document population will not be much larger than the sample size needed for a much smaller document population.



This means that sampling can be relatively inexpensive even for very large document populations. Thus, while many might be reluctant to introduce sampling and metrics into their search and review process because of concerns about additional burden, in most cases the burden may be relatively small (and many of the documents in a sample would have needed to be reviewed anyway).

Second, the Confidence Interval (or margin of error) is a crucial part of your measurement. For example, if the Confidence Interval for the 15% Depth for Recall example above is +/- 5%, that means you will actually need to review anywhere between 10% to 20% of the document population to achieve the targeted Recall level (in the example of a population of 100,000 documents, that means a range of between 10,000 to 20,000 documents requiring review to achieve the targeted Recall).

Third, the relevant sample size required to achieve a given Confidence Level and Confidence Interval is tied to what you are measuring. For example, when measuring Recall, and seeking a particular

Confidence Level and Interval, the necessary sample size will be impacted by the Prevalence of responsive documents in the sample. This can sometimes mean the sample must be larger to find enough responsive documents in the sample to match your goals for the Confidence Level and Confidence Interval of your Recall measurements.

For example, to confirm a Recall level of 50% with a Confidence Level of 95% and a Confidence Interval/margin of error of +/- 5%, you would need 385 responsive documents in your sample. If the responsive documents are 10% of the population, you would need a sample of 3,850 documents to get 385 responsive documents, but if the responsive documents are 33% of the population you would need a sample of 1,155 documents to get 385 responsive documents. In other words, a higher Prevalence rate will result in needing a smaller sample size than if you had a lower Prevalence rate.

Fourth, don't forget that Prevalence, Recall and Precision measure actual performance, while the Confidence Level and Confidence Interval are simply measuring the accuracy of those metrics.

What do brown M&Ms have to do with it again?

In sum, a few metrics—generated from review of what can be surprisingly small samples of the overall document population—can make a world of difference in helping you to develop an efficient and cost-effective search and review process, and to defend the ultimate results upon making your document production.

Having counsel and service providers involved in your document search and review process who are familiar with and can accurately generate such metrics can yield significant cost savings and verify the efficacy of the process.

After all, Van Halen used a convenient metric to help avoid a disaster befalling band members and the audience in their elaborate stage show. You can do so, too, in your document search and review process.

Gareth Evans is a Partner at Gibson, Dunn & Crutcher LLP.

David Grant is a Senior Managing Director at FTI Consulting.

Gibson Dunn and FTI Technology help clients manage the risk and complexity of e-discovery. For more information, please visit: www.gibsondunn.com and www.ftitechnology.com.

GIBSON DUNN

 FTI CONSULTING | TECHNOLOGY