# What Should You Do with an Abundance of Validation Data?

**Lucinda Smart, Benjamin Wright**

Kiefner and Associates, Inc

**PPIM 2025**
PIPELINE PIGGING & INTEGRITY MANAGEMENT
JANUARY 27-31 • HOUSTON

## Abstract

The concept of in-line inspection (ILI) tool validation could be summarized as follows: identify dig locations based on ILI-reported features, measure these features in the ditch, compare in-ditch results with ILI predictions, log the results for future ILI assessments at specified intervals, and then proceed to the next pipeline segment. However, the complexities involved in ILI tool validation make it interesting. Each pipeline and ILI run is unique, presenting its own set of challenges that must be addressed.

This paper will provide insight into the analysis decisions made in a case in which a pipeline underwent multiple axial magnetic flux leakage (MFL-A) ILI runs, along with a circumferential MFL (MFL-C) ILI, which reported thousands of defects of varying geometries and depths. The focus will be on the efforts that may be performed to validate this tool run, how to resolve the differences in reporting between the two tool technologies, and how to ensure the pipeline's safety for continued operation.

Understanding the levels of validation as outlined in API 1163 is essential in whether the tool run can be deemed validated. When validation digs are performed to assess tool performance, particularly when a wide array of defects are reported, it is important to account for the types of defects that may be found and to have appropriate technology available to capture them. More data can explore many facets, and handling that appropriately can be challenging. The results can give a wealth of insight into the tool's performance and the line's condition. This paper will continue to provide recommendations on what factors an operator should consider when obtaining a large data set of correlated features from validation digs. Understanding ILI tool performance and the true condition of the pipeline is critical to maintaining the integrity of pipeline assets.

## Introduction

The state of aging oil and gas infrastructure necessitates regular inspections to evaluate the condition and safety of our pipelines. Over time, pipeline evaluation technologies have improved, leading to more measurement methods to gain copious amounts of data. Modern ILI tools often combine multiple technologies, such as Axial MFL and Deformation, Circumferential MFL, UT, MFL for hard spot identification, or EMAT, to maximize data collection and provide comprehensive assessments. Combination tools allow operators to detect and assess multiple types of anomalies in a single inspection run, reducing operational downtime and inspection costs.

Data utilization in modern pipeline integrity management should take advantage of the large amounts of data collected by ILI tools, both past and present, assuming the data is verified and validated. They should be analyzed to gain insights into the condition of pipelines. Advanced software and algorithms can be used to identify and classify anomalies by type, size, and location, prioritize maintenance and repair activities based on the severity of detected defects, monitor changes over time to assess the effectiveness of mitigation measures, and improve predictive models for pipeline integrity and risk assessment. It can be difficult to compare two tool technologies; however,

such is most commonly the case with MFL-A versus MFL-C, where both tools may detect and size metal losses, but each tool has a different specification for the probability of detection/identification and the sizing accuracy specification.

## Evolution of API 1163

The API Standard 1163, "In-line Inspection Systems Qualification," has evolved significantly since its first release to improve ILI systems' qualification and performance validation. Each edition has introduced enhancements to ensure greater consistency and reliability in assessing pipeline integrity.

The 1st Edition, published in August 2005, established foundational guidelines for qualifying ILI systems. It emphasized validating the performance of both the tools and the personnel interpreting the data, laying the groundwork for subsequent, more detailed standard iterations.[1] The validation process provided in this edition was fairly prescriptive, using tables to support the justification of a rejectable tool run based on whether or not the data was within the tool vendor's tool specification. However, with the concept of ever-changing and improving technologies, the original document provided enough generalities to be applied in a myriad of conditions.

The 2nd Edition, released in April 2013, expanded on this framework by introducing a structured, tiered approach to validating ILI performance, moving away from the tables for determining an acceptable tool performance. This included three distinct levels of validation: a basic comparison with records from previous inspections, statistical analysis of ILI data against known anomalies, and a comprehensive assessment involving field verification and detailed analysis.[2] These enhancements provided operators with a more systematic methodology for evaluating inspection data.

The most recent revision, the 3rd Edition, was published in September 2021.[3] This edition refined the validation process further by formalizing the requirements for verifying and validating ILI systems within pipeline integrity management programs. It also expanded the guidance on conducting each validation level, offering specific methodologies and statistical tools. Additionally, this edition introduced recommendations for conducting root-cause analyses in cases of inspection failures, which are critical for identifying and addressing deficiencies in ILI tool performance.

Overall, API 1163's evolution reflects a continuous effort to improve ILI systems' precision, reliability, and safety in pipeline operations.

## Application of API 1163 for Tool Validation

The 2nd Edition is currently incorporated by reference in regulatory requirements and will be the basis of discussion for this paper's tool validation. This edition describes methods in Section 8 that can be applied to validate that reported inspection results are within the performance specification for the inspected pipeline. The Standard distinguishes between results with and without field verification measurements. API 1163 Section 8 (Figure 6) describes a process for validating ILI measurements using three levels of validation. The three levels of validations all consist of the

following steps and differ based on the risk of the pipeline segment and the amount of validation data available:

>A process verification or quality control (§8.2.2 and Annex C.1).
>A comparison with historical data for the pipeline being inspected (§8.2.3).
>A comparison analysis of pipeline component records (§8.2.4).

**Validation Level 1 (Annex C):**
>A comparison with large-scale historic data for pipeline segments similar to the pipeline being inspected (§8.2.3).

Validation Level 1 is designed for pipelines with anomaly populations characterized by a low risk of consequence or probability of failure. This level typically applies when limited or no independent validation measurements are available for the pipeline under inspection. At this validation level, it is assumed that the performance specifications provided by the in-line inspection (ILI) tool are neither accepted nor rejected for the given inspection run. Consequently, the validity of the ILI results cannot be refuted based solely on a Level 1 validation. Instead, higher validation levels, such as Level 2 or Level 3, which involve more rigorous analysis methods, are necessary to formally reject the outcome of an ILI inspection. These advanced levels incorporate direct measurement data and statistical analysis to substantiate or challenge the tool's reported performance metrics.

**Validation Level 2 (Annex C):**
>A comparison with field excavation results is warranted by the reporting of significant indications (§8.2.6).

Validation Level 2 also applies to pipelines with a low risk of consequence or probability of failure due to ILI indications. However, the pipeline has enough validation measurements to confidently state whether the ILI tool performed according to specification and possibly reject the ILI run. A Level 2 validation does not let one confidently state that the ILI tool performed within specification. At least five validation measurements are required. If the ILI tool specification can be rejected, there is the option to progress to a Level 3 validation, which may require additional validation measurements.

**Validation Level 3 (Annex C):**
>A comparison with field excavation results is warranted by the reporting of significant indications (§8.2.6).

Validation Level 3 applies to pipelines with a higher risk of consequence or probability of failure, with significant indications reported by ILI. A statistically significant number of validation measurements is typically required to estimate the as-run tool performance.

Tool performance can be accepted, rejected, or inconclusive depending on the data analysis using the API 1163 decision chart process. If tool performance is deemed inconclusive, it does not mean

the inspection failed. Instead, additional action may be required, such as performing additional validation digs, accepting the determined tool performance, adjusting the depth accuracy applied to the reported ILI features, or requesting the ILI Vendor regrade the data. API 1163 Section 8 (Figure 6) summarizes the system results evaluation process. "Historical data" means data limited to the particular line, while "large-scale historical data" means the data on this line and any similar diameter lines with the same ILI tool type used for inspection.

## Working with Validation Data: Statistical Justifications

Most operators pose a question for tool validation: "How much data do I really need, or how many digs do I need to do?" This is an important question because prioritization of time and resources is vital in a quality integrity management program. Depending on the severity of the features reported by ILI, statistical methods used, and the necessity to determine true tool performance, this could be as few as five data points for a simple reject/accept assessment, perhaps 20 or so to gain more confidence in the assessment, or several hundred to support the establishment of true tool performance with a reasonable amount of confidence in the results. The more sample data added to the assessment, the closer we achieve the true tool performance. This fundamentally follows the Law of Large Numbers, which states that as the sample size increases, the empirical distribution of the data becomes a more precise representation of the true population deviation. This precision can highlight subtle irregularities in the data that might not be apparent in small samples.

The Central Limit Theorem (CLT) is a fundamental statistical principle that states that the sampling distribution of the sample mean (or sum) of a sufficiently large number of independent, identically distributed random variables will be approximately normally distributed, regardless of the original distribution of the population. Since we primarily deal with measurement errors, we need to consider this concept in what our data can tell us. As our ability to collect and manage huge amounts of data improves with higher resolution technologies and digitization of the data, we need to be cautious and not hide important information that the data can provide.

Too much validation data can unexpectedly lead to challenges in accurately assuming a normal distribution, as statistical tests are more sensitive to deviations from normality.[4] This issue arises from the interaction between statistical significance and practical significance as sample sizes grow. Statistical tests for normality, such as the Shapiro-Wilk or Kolmogorov-Smirnov tests, are designed to detect departures from normal distribution.[5] When the sample size is large, these tests become extremely sensitive to even minute deviations that may not have practical significance. As a result, the test might reject the null hypothesis of normality, even if the data is effectively normal for most analytical purposes.

A poor result from a normality test might lead analysts to conclude that their data is not normally distributed.[6] This could lead to either using unreasonable, overly complicated applied trend fits to the data or discarding data that might still be useful. Therefore, it is important to understand why

outliers or non-normal data exist and how to account for them, whether by removing them or creating a subset of data that adjusts the analysis accordingly for each data subset.

To reduce the potential for over-reliance on normality tests in large datasets, it can be beneficial to use visual assessments, such as histograms, cumulative distribution plots, or Q-Q plots.[7] A histogram of normal data will produce a bell curve, a normal cumulative distribution plot will produce a smooth s-curve, and a Q-Q plot will produce data that aligns on a 45-degree reference line.

While extensive validation data is beneficial for thorough statistical analysis, ensuring a balanced normality interpretation is important. As analysts, we must consider the practical implications of deviations in the data and any associated unusual events and not rely solely on statistical significance to assess whether the data is sufficiently normal for the intended application.

## Methods to Evaluate Subsets of Data

The volume of data generated in pipeline integrity management, particularly from ILI and field verification processes, can be overwhelming. As noted previously, it may hide important information needed to understand the characteristics of the severity of features reported on the pipeline. Determining subsets of data is a practical and valid approach to derive actionable insights without dealing with overly complex statistical methods. Below, a handful of approaches are defined, along with examples of their application.

### Depth-Based Subsets

Depth is a critical parameter in evaluating metal loss or other pipeline anomalies, as deeper defects typically pose greater risks to structural integrity, and deeper defects have been known to see greater variability in their reported depths[8]. Sorting data by depth involves categorizing anomalies based on their measured depth relative to the pipe wall thickness (e.g., <10%, 10–30%, >30%). This is also useful for evaluating the probability of detection (POD) of the ILI tool, as shallower metal loss depths tend to have a lower POD since the severity of the thickness change is not as great.

### Geometry-Based Subsets

The shape of the metal loss is imperative for analysis, particularly because the ILI Vendor provides tool specifications based on the geometry of the feature, more commonly known as the POF categories[9]. This is due to the physics of how MFL technology functions and the application of that technology will determine the accuracy of the depth call based on the shape the tool can detect.

### Dig Location-Based Subsets

The application of dig location subsets is primarily to address the potential for human error. Depending on the in-ditch measurement techniques applied, issues with the calibration, record-keeping, or correlation of ILI features to the in-ditch found defects may exist. This is particularly relevant for understanding the potential variability between technicians using manual methods, such

as pit gauges. However, it is no longer common practice to use pit gauges for more than initial observations of the deepest features found on site.

In-Ditch Measurement Techniques
Evaluating data based on the measurement technique used in field verification can also improve the understanding of validation results. The accuracy and resolution of techniques, such as laser profilometry, automated ultrasonic testing (AUT), straight-beam ultrasonic testing (UT), and pit gauges vary significantly.

Laser profilometry provides detailed 3D profiles of defects, enabling precise geometry measurements for depth and volume analysis. AUT provides high-resolution data for wall thickness and defect characterization, is particularly useful for planar defects, and can evaluate internal and external defects. Straight-beam UT is commonly used for measuring wall thickness and defect depth for internal and external defects but lacks the detail of AUT or laser profilometry. Pit gauges are a traditional, manual technique useful for quick measurements but subject to user variability and lower precision.

It is anticipated that most in-ditch measurements will be collected using the same methods throughout a dig program, but when there are differences, the performance of each method will need to be accounted for in the assessment of the ILI tool's performance.

# Validation Case Study

## Background
The case study incorporates two pipeline segments that were inspected within two weeks of each other. Both lines were inspected using the same MFL-A, MFL-C, and Deformation tools. Both segments have the same diameter, vintage of pipe, seam types, wall thickness, and grade. The ILI runs were 31 to 36 miles in length and had reported ILI feature counts of approximately 42,000 and 52,000, respectively. The main threats to these segments were external metal loss corrosion and pre-1970 ERW seams. The ILI runs were verified using API 1163 Section 7 principles and deemed acceptable for analysis. Due to the similarities across the board, it was decided to combine the validation dig data for statistical analysis to obtain results for both the MFL-A and MFL-C tools with high levels of confidence.
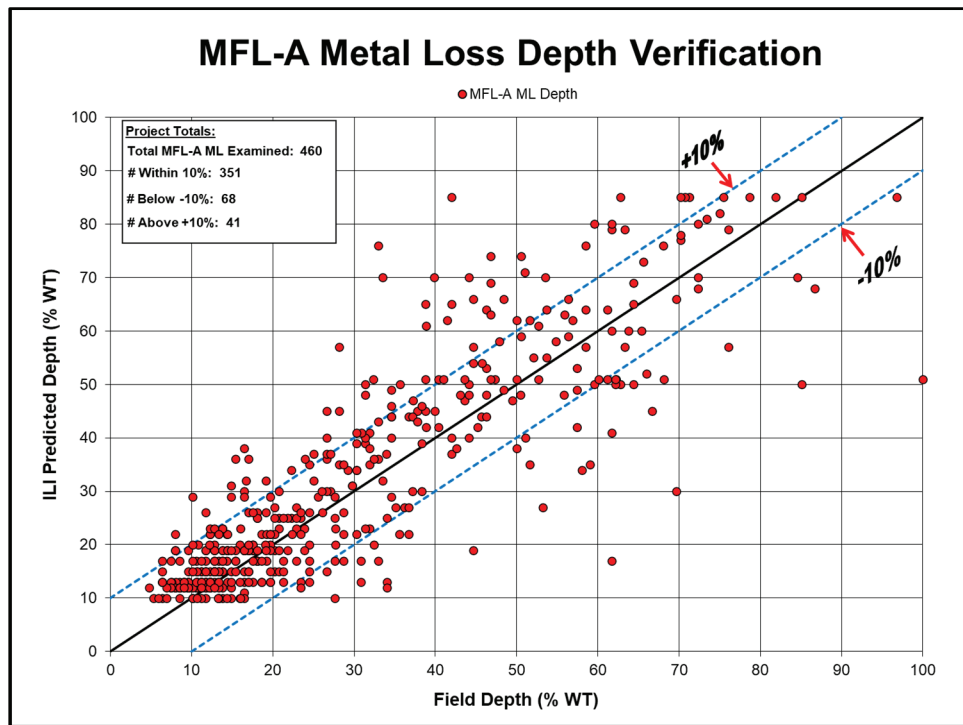
## Obtaining Validation Data
Based on the reported ILI features, over 500 features combined between the two segments needed to be addressed prior to the next ILI. The ILI features that need to be addressed varied between seam weld metal loss on pre-1970 ERW pipe, metal loss with depth over 80% WT with the addition of tool tolerance, dents with metal loss, and corrosion-related metal loss with response times occurring prior to the next re-inspection interval. The features were addressed through either in-ditch investigations or pipe replacements. In-ditch depths of features were initially checked using a pit

gauge, and the final data for analysis was collected using HandyScan laser profilometry. There were 168 individual anomaly investigation digs performed in total, with several digs documenting more than one target feature. It is important to gain as much information as is available during the investigation process, so scanning the entirety of the region exposed beyond the single target feature allows for a more complete evaluation of the pipe and the data.

In-ditch data from previous ILI assessments was also available for 46 anomaly investigations and were included as validation data points where appropriate based on repair type.

### Initial Statistical Analysis

The field measurements' depths, lengths, and failure pressures (calculated using modified B31G) were correlated to the current ILI assessment using spatial mapping software rather than manual correlation for improved correlation accuracy. This method correlated 458 MFL-A metal loss features, 218 MFL-C metal loss features, and 106 MFL-C seam weld metal loss features. **Figure 1** and **Figure 2** express the in-ditch and ILI data pairs as unity plots. Note that the MFL-A plot indicates 460 correlated metal loss features. However, two of these were removed because they were outliers.
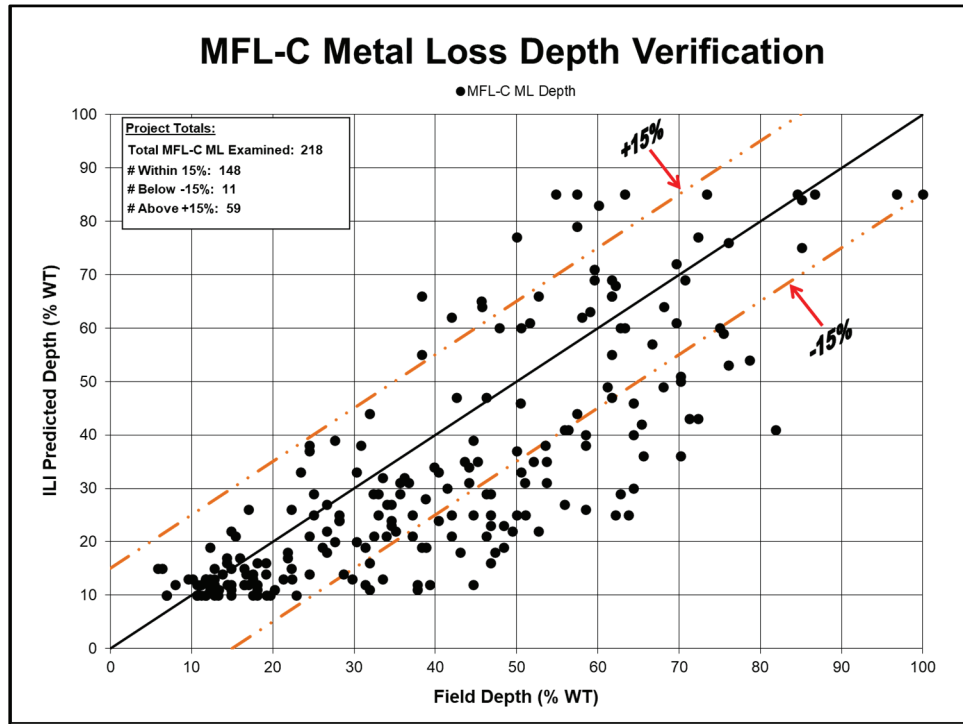
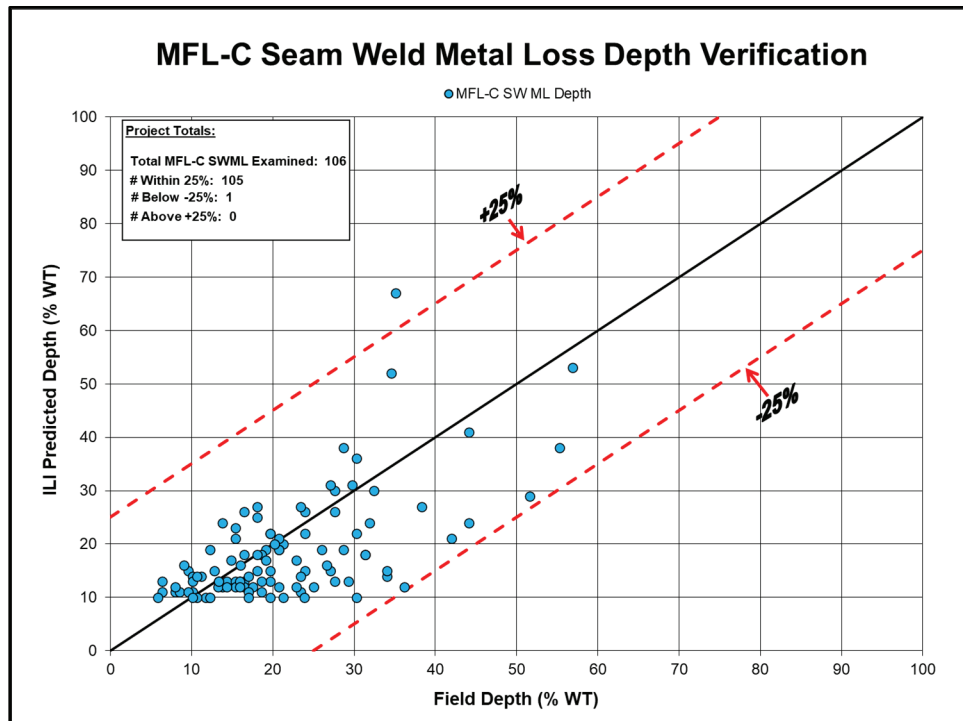**Figure 1.** Unity Plots for MFL-A and MFL-C Metal Loss



**Figure 2.** Unity Plot for MFL-C Seam Weld Metal Loss

The unity plots show that the majority of depths for the MFL-A metal loss features evaluated in-ditch are called within tool specifications (±10% WT). The MFL-C metal loss feature depths reported by ILI are generally under-called, and the MFL-C seam weld metal loss feature depths are within tool

specifications (±25% WT). The length unity plots show that the ILI tools under-call the length for metal loss features. The MFL-C tool shows a trend of reporting metal loss feature lengths more accurately than the MFL-A tool.

Given the number of correlated metal loss features available, an API 1163 Level 3 analysis was performed. There are various probabilistic analyses to determine if the tool performed within the specified confidence or certainty of the ILI vendor. Kiefner evaluated whether the probability met or exceeded the 80% certainty specified by the ILI vendor using a 95% confidence interval. Two different probabilistic approaches were considered for this analysis. These approaches are the Agresti-Coull and Clopper-Pearson methods as specified in the API 1163 2nd Edition. Clopper-Pearson results are shown here, but the final upper and lower bounds for the accept and reject criteria are similar for this sample size.

An initial statistical analysis was performed to determine the average and standard deviation and whether outliers or extreme values were present. **Table 1** shows the results from the statistical analysis; a negative value represents that the ILI tool has under-called the correlated features compared to the in-ditch data. Two MFL-A metal loss features were removed from the analysis due to extreme values. An example of an extreme value was a deep pinhole feature within general corrosion, as shown in **Figure 3**. This is a geometrically difficult feature for an MFL-A tool to report accurately.

**Table 1.** Summary of Sizing and Population Density

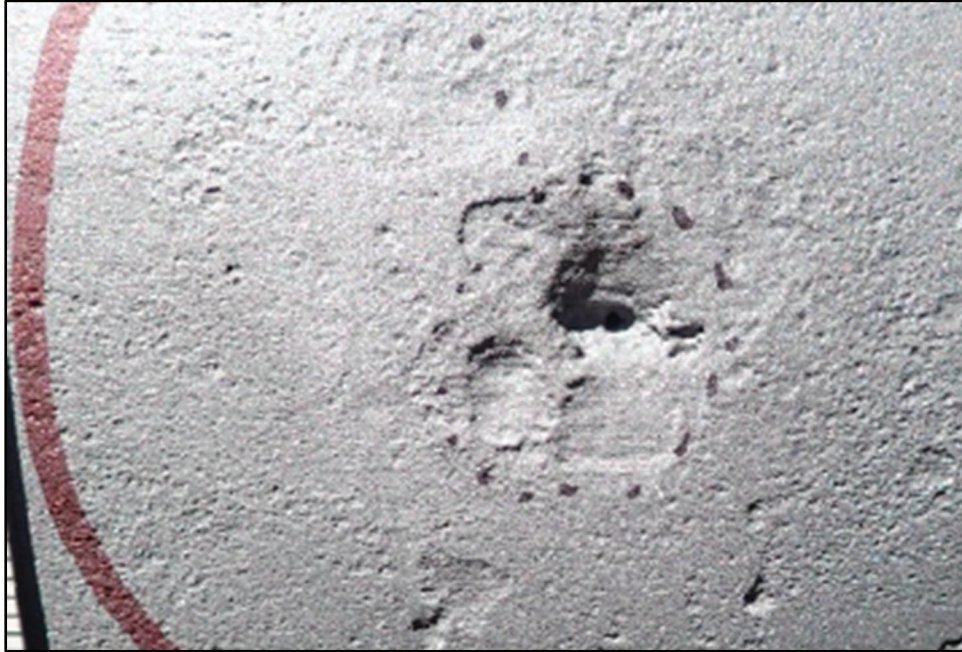|  | MFL-A | MFL-C | |
| --- | --- | --- | --- |
|  |  | Metal Loss | SW Metal Loss |
| Vendor Specified Depth Accuracy (%) | ±10 | ±15 | ±25 |
| Total Number of Matched Features | 460 | 218 | 106 |
| Number of Features Used in Analysis | 458 | 218 | 106 |
| Total Number of Features w/in Tool Specifications | 351 | 18 | 105 |
| Average Size Difference (%) | 2.6 | -7.0 | -2.9 |
| Standard Deviation (% WT) | 9.6 | 13.5 | 8.6 |
| 80% Random Error (% WT) | ±12.3 | ±17.3 | ±11.1 |
| Clopper-Pearson UB (%WT) | ±12.0 | ±22.9 | ±13.0 |

**Figure 3.** Example of a Pinhole (MFL-A Depth: 51% WT, Field Depth: Through Wall)

If the ILI Vendor's stated tool tolerance accuracies are applied, MFL-A metal loss features claim to be within ±10% WT, MFL-C metal loss features claim to be within ±15% WT, and seam weld metal loss features reported by MFL-C claim to be within ±25% WT, under a certainty of 0.80 at a confidence level of 95%.

Based on the Clopper-Pearson confidence interval analysis results, the following feature bias and tool tolerances were accepted for each tool technology.

- MFL-A metal loss features: bias of -2.6% WT and tool tolerance of ±12.0% WT
- MFL-C metal loss features: bias of +7.0% WT and tool tolerance of ±22.9% WT
- MFL-C seam weld metal loss features: bias of +2.9% WT and tool tolerance of ±13.0% WT

### Evaluating Reported Features by Different Technologies

The metal loss anomaly features reported in the original ILI pipeline listing were not correlated between the MFL-A and MFL-C tools. This means there were some reported features called by both tools, with differing sizing results. As a result, the process of correlating the MFL-A and MFL-C calls needed to be performed to ensure features were not multiple interacting defects but were, in fact, the same feature reported differently by each tool. Kiefner correlated the data again using spatial mapping internal proprietary software to visualize the correlation and ensure the accuracy of defect comparisons.

Although not discussed in depth within this case study, the lengths of features were also verified along with the depths, as shown previously. The dig results show that the MFL-A tool reports metal

loss depths better than the MFL-C tool depths, while the MFL-C tool results show better length than the MFL-A tool. This is expected since MFL-C technology can typically capture the axial changes along the pipeline and will detect the start and end of volumetric metal loss more accurately in most cases. This is expected since MFL-C technology can typically capture the axial changes in depth along the pipeline and will detect the start and end of volumetric metal loss more accurately in most cases. As a result, features with correlating calls between the ILI tools used the MFL-A depth and the MFL-C length to calculate failure pressures and determine a more optimal response time. Utilizing both data sets in this manner allows for improved prioritized response to the population of metal loss features as part of a whole dig program.

## Evaluating Reported Features by Depth Error Distribution

Measurement errors should tend to be normally distributed, and if they are not, further investigation needs to be performed to determine possible causes of the non-normal distribution. The metal loss features can be evaluated using cumulative distribution plots to visualize the results, with a smooth s-shaped curve providing this quick indicator of a normal distribution. A more vertical s-shape indicates a lower tool tolerance, and a shift to the right or left of the data centered around zero is a visual representation of the mean bias shift. Kiefner evaluated the MFL-A data based on 10% WT depth bins, with results provided in **Figure 4**, separated into two plots for visual clarity. The data shows that the tool tolerance appears to increase as the depths increase. This is notable since responding to the deepest features is a common protocol for operators to reduce the risk of leaks or burst failures on the line. If the deeper features have more variability, it decreases the confidence of an accurate prioritized response. Conversely, this data indicates that shallow defects are being reported accurately. Therefore, we have confidence that we are not missing a deep and severe defect that has been inaccurately reported as shallow.
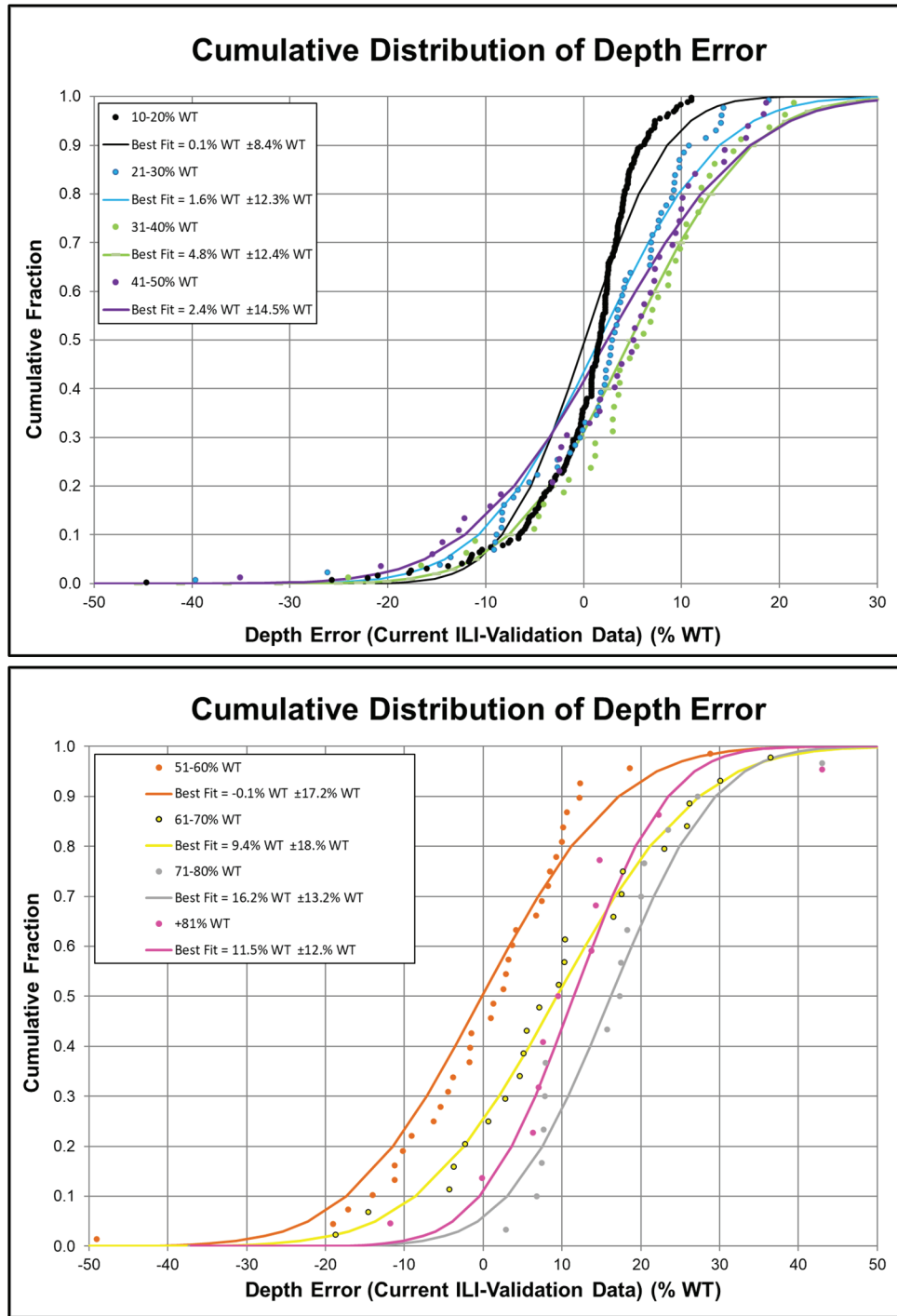
**Figure 4.** Cumulative Distribution Example of MFL-A Metal Loss Depth Bins

### Evaluating Reported Features by Geometry POF Category Distribution

ILI Vendors provide tool performance parameters based on the geometry they can size and detect. The inherent limitations of MFL mean that the difficulties in accurately detecting the edges of volumetric metal losses will be apparent based on the feature's geometry. Therefore, it is sensible to evaluate the reported ILI features tool performance from verification digs similarly and divide the data into geometry subsets based on the POF category definitions, including General Corrosion,

Pitting, Axial Grooving, and Axial Slotting. Therefore, it is sensible to evaluate the reported ILI features tool performance from verification digs similarly and divide the data into geometry subsets based on the POF category definitions, in this case including General Corrosion, Pitting, Axial Grooving, and Axial Slotting. These were used in this case study for evaluating tool performance, and it was found that pitting defects were the least accurately sized by MFL-A, as shown in **Figure 5**. The implications would be that any pitting defects reported by ILI may be prioritized to ensure they are not being undercalled in depth and reduce the overall risk of a leak failure on the line.
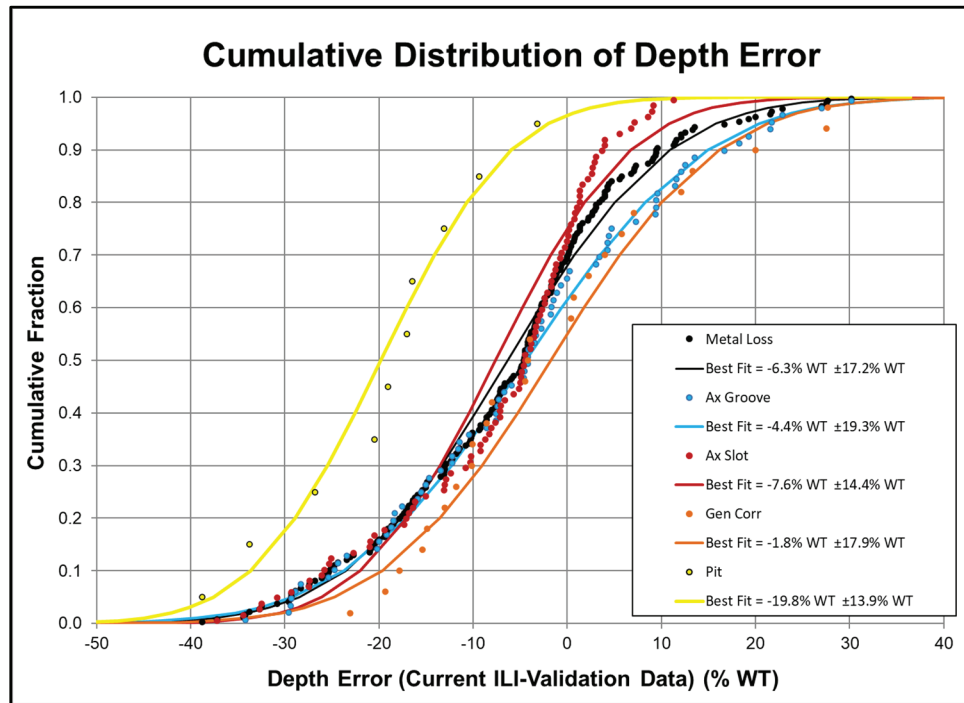


**Figure 5.** Cumulative Distribution Example of MFL-A Data Metal Loss Based on POF Category Geometry

### Evaluating Reported Features by Dig Location Distribution

Factoring human error into the dig data can be key in understanding why a data set does not fit with the whole group, or in other words, it can explain why the sample data is not normally distributed. If pit gauges are used, the error of this tool needs to be considered and included in the ILI tool performance calculations. It may be significant enough to affect the validation results, or it may compound the error and result in inconclusive validation of the ILI data. However, using only pit gauge measurements in modern dig programs is uncommon, particularly for complex corrosion geometry metal losses. In this case study, the non-destructive evaluation (NDE) vendor used pit gauges to roughly determine the depth of each metal loss feature to understand the overall severity, and then laser profilometry tools were used to determine the true depth. Laser profilometry tools have an accuracy of sizing metal losses of up to 0.0009 inches, far exceeding the accuracy found in ILI tools.

Even though laser profilometry is a highly accurate measurement technique, it can still be beneficial to evaluate the dig location when evaluating dig locations. Something as simple as poor tool calibration could be an issue that is easily remedied if observed at a single dig location. In this case study, there were no discernible differences from the dig locations, separated according to the distance on the pipe for every 50,000 feet. It should be noted that this assessment should be done with careful consideration, as an individual location may also simply be susceptible to a certain geometry of metal loss and lead to a bias shift in the data, not necessarily that the NDE technician incorrectly gathered the data.
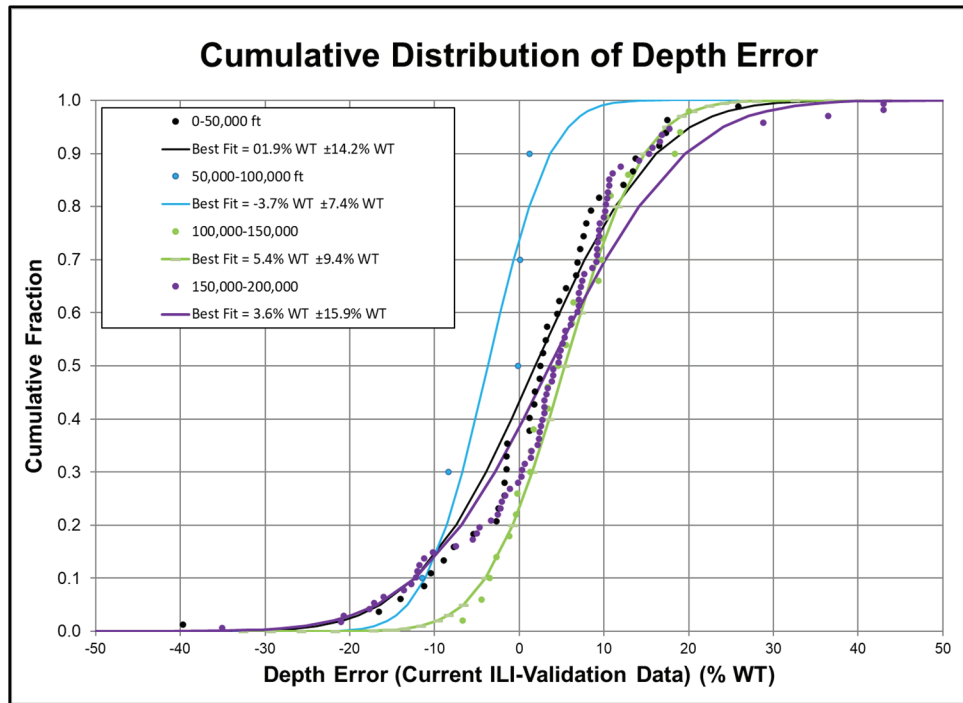


**Figure 6.** Cumulative Distribution Example of MFL-A Metal Loss based on Location

## Conclusions

As technologies for both ILI tools and NDE in-ditch techniques continue to improve, it is crucial to focus on how to use the data they provide effectively. This is key to making smart decisions without getting overwhelmed by the sheer volume of information or falling into the trap of simply adding more data to a unity plot, running calculations, and assuming everything checks out.

We must be careful about lumping all the data together without thinking critically about its context or quality. Blindly aggregating data can hide important details or lead to overly simplistic conclusions. Evaluating and questioning how the data is used helps ensure we make informed choices and get the most out of these advanced tools rather than just going through the motions.

# REFERENCES

[1] American Petroleum Institute. (2005). *API Standard 1163: In-line Inspection Systems Qualification.* 1st Edition.

[2] American Petroleum Institute. (2013). *API Standard 1163: In-line Inspection Systems Qualification.* 2nd Edition.

[3] American Petroleum Institute. (2021). *API Standard 1163: In-line Inspection Systems Qualification.* 3rd Edition.

[4] Razali, N. M., & Wah, Y. B. (2011). Power comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors, and Anderson-Darling tests. *Journal of Statistical Modeling and Analytics, 2*(1), 21-33.

[5] Ghasemi, A., & Zahediasl, S. (2012). Normality tests for statistical analysis: A guide for non-statisticians. *International Journal of Endocrinology and Metabolism, 10*(2), 486-489.

[6] Tabachnick, B. G., & Fidell, L. S. (2019). Using multivariate statistics. *Pearson.*

[7] Wickham, H., & Grolemund, G. (2016). R for data science: Import, tidy, transform, visualize, and model data. *O'Reilly Media, Inc.*

[8] Smart, L. Haines, H. (2014). Validating ILI Accuracy Using API 1163. Proceedings of the 2014 10th International Pipeline Conference. Calgary, AB, Canada. IPC2014. ASME.

[9] Pipeline Operators Forum. (2021). Specifications and requirements for in-line inspection of pipelines. POF 100. Standard Practice. www.pipelineoperators.org.