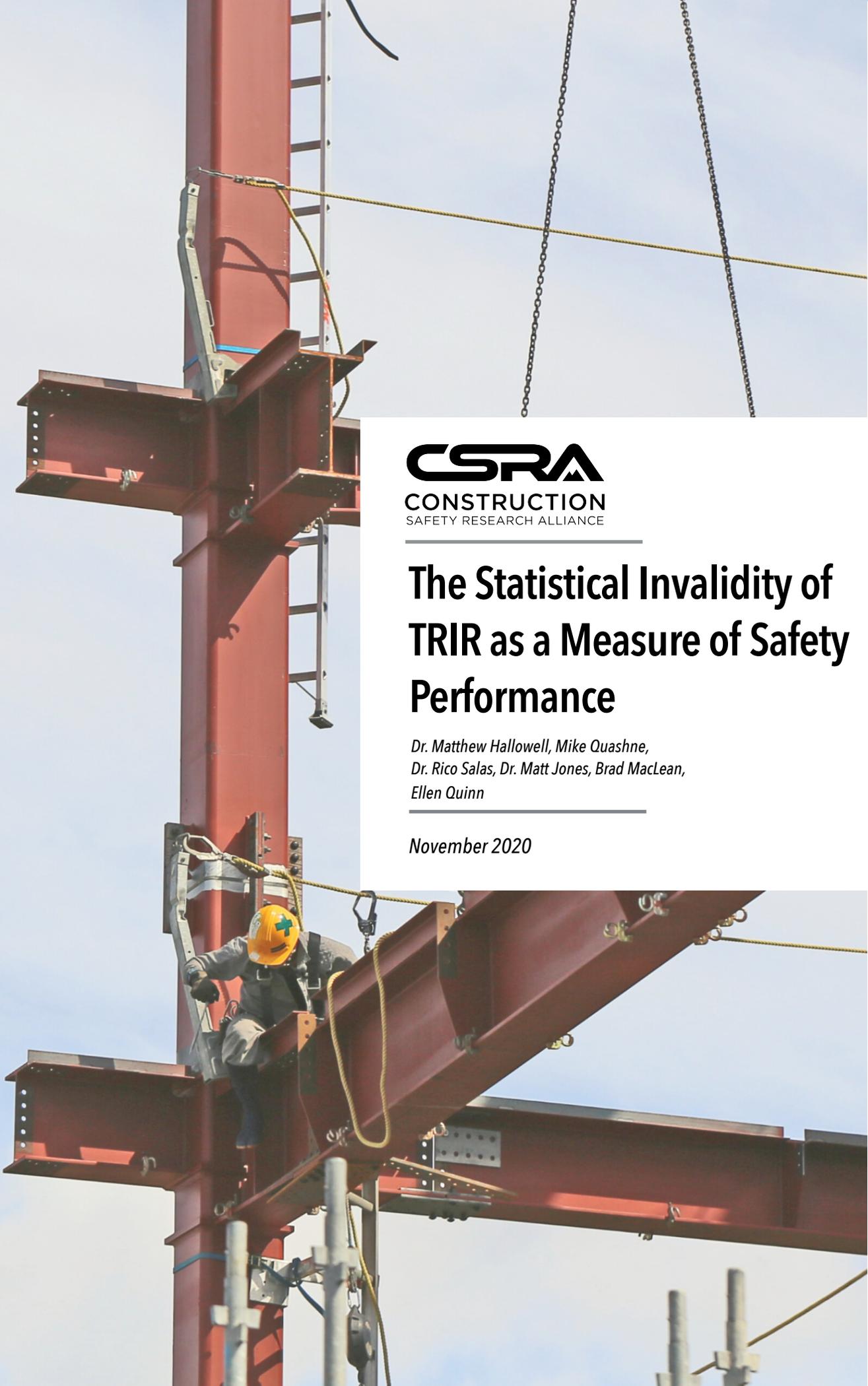




The Statistical Invalidity of TRIR as a Measure of Safety Performance

*Dr. Matthew Hallowell, Mike Quashne,
Dr. Rico Salas, Dr. Matt Jones, Brad MacLean,
Ellen Quinn*

November 2020



Key Take-Aways

Parametric and non-parametric statistical analysis data revealed that:

01.

There is no discernible association between Total Recordable Incident Rate (TRIR) and fatalities;

02.

The occurrence of recordable injuries is almost entirely random;

03.

TRIR is not precise and should not be communicated to multiple decimal points of precision; and

04.

In nearly every practical circumstance, it is statistically invalid to use TRIR to compare companies, business units, projects, or teams.

17 YEARS OF

DATA &



TRILLION WORKER HOURS



Introduction

“ *Given the way that it's used, to what extent is TRIR statistically valid?* ”

Total recordable incident rate (TRIR) has been used as the primary measure of safety performance for nearly 50 years. Simply, TRIR is the rate at which a company experiences an OSHA-recordable incident, scaled per 200,000 worker-hours. TRIR is based upon a standard definition of a “recordable” incident that was created and institutionalized in the recordkeeping requirements of the Occupational Safety and Health Act of 1970 (US Bureau of Labor Statistics 2019). According to the general criteria, an incident is recordable if it results in work-related injury or illness involving loss of consciousness or requiring medical treatment beyond first aid, days away from work, restricted work, or transfer to another job (US Department of Labor 2010). Since organizations conform to this same definition, the TRIR metric has been used to compare industries, sectors, companies, and even projects.

TRIR is used in many ways to measure safety performance, from the worksite to the board room. For example, organizations use TRIR to report results, benchmark against peers, prequalify and select contractors, evaluate the performance of managers, and track the impact of safety initiatives (Lofquist 2010; Manuele 2008; Wilbanks 2019). TRIR is also a primary safety metric of concern among executives because it may impact worker's compensation insurance premiums, influence public image, and be scrutinized by potential customers or investors (Karakhan et al. 2018; Lingard et al. 2017; Lofquist 2010; Ritchie 2013; Salas 2020; Truitt 2012). Although other lagging measures like the Days Away; Restricted; Transferred (DART) rates are also considered, no safety metric is as ubiquitous as TRIR.

Despite the pervasive use of TRIR, its limitations are being recognized. For example, some argue that TRIR is a poor reflection of safety performance because it does not account for the actual or potential severity of an incident (Toellner 2001). For example, a four-stitch cut to the finger is counted in the same way as a fatality, and a near miss with the potential to be fatal is not in the TRIR metric at all. Others point out that TRIR is reactive in nature as it only counts incidents and does not consider the underlying safety program (Lingard et al. 2017; Lofquist 2010; Salas and Hallowell 2016).

More recently, some have begun to question the statistical validity of TRIR, suggesting that recordable injuries happen so infrequently that the metric is not stable or reliable. Since TRIR is typically reported over relatively short time frames (i.e., months, quarters, or years), the number of recordable injuries in each period can be exceedingly small. Therefore, it is suspected that the confidence interval of typical reporting periods is so wide that it renders the metric useless. This potential limitation is implicitly recognized by those who criticize TRIR as unfairly biased against small companies.

To better understand the validity of TRIR as a performance metric, this study attempted to answer the basic question: *Given the way that it's used, to what extent is TRIR statistically valid?*

More specifically, we aimed to test whether TRIR is statistically stable, precise, predictive, and indicative of high-severity events. The answer will help clarify if TRIR should be used to make important business decisions like comparing two contractors, evaluating the safety performance of managers, or concluding that a new safety intervention is effective.

Background

A statistical analysis of any variable requires its decomposition into underlying components. Since TRIR is a rate, the metric is comprised of two variables: ***the number of events and time***. Here, events are injuries and illnesses that conform to OSHA's definition of recordability. Time, on the other hand, is expressed as worker exposure hours. Because expressing TRIR as the number of incidents per worker-hour would yield an excessively small fraction, TRIR is scaled per 200,000 worker-hours. Since 200,000 worker-hours equates to approximately 100 employees working full time for one year, TRIR also reflects the percentage of workers who suffer a recordable incident in a year.

To compute a Total Recordable Incident Rate (TRIR) an organization applies Equation 1 for a specific period. For example, a company that accounts for 3 recordable incidents over 350,000 worker-hours in a month would have a TRIR of 1.71 per 200,000 worker-hours for that month.

Equation 1:

$$\text{TRIR} = \frac{\text{Number of Recordable Incidents} \times 200,000}{\text{Number of Worker-Hours}}$$

TRIR is typically reported as a single number. Statistically, this means that the company takes the single TRIR value to represent safety performance over the period as if it was the only possible outcome. This number is also often reported to one or more decimal points of precision, where subtle differences in TRIR are assumed to be meaningful (e.g., the difference between 1.2 and 1.6 from one month to the next). Although TRIR is used in this manner, the underlying assumptions have never been validated.

In addition, TRIR is often used as a dependent variable by academic researchers. TRIR has been applied frequently as an objective, empirical metric of safety performance and is often conceptualized as an ideal response variable. In fact, the authors of this paper have used TRIR on multiple occasions as a dependent variable when identifying which safety practices are more effective than others (Hallowell 2010; Hallowell and Gambatese 2009); validating safety leading indicators as predictive (Alruqi and Hallowell 2019; Hallowell et al. 2013; Hinze et al. 2013; Lingard et al. 2017; Salas and Hallowell 2016) and measuring intrinsic relationships between safety and other performance metrics like productivity and quality (Wanberg et al. 2013).



Example Cases

Case examples were created to illustrate the potential limitations of using TRIR as a comparative metric. The hypothetical cases were designed to vary greatly in the number of recordable injuries, the number of worker-hours accumulated, and the resulting TRIR. Case A represents a new contractor that has a recordable incident early in their company history. In contrast, Case B is a medium-sized company with fewer than 500 employees that accumulates just less than 1 million worker-hours in a year. Finally, Company C is a large company that amasses millions of worker-hours in a year.

Company A is an extreme example that underscores the limitation of reporting TRIR over very short timeframes or for small businesses. On the surface, the TRIR of 200 for Company A could be judged as over 50 times worse than average TRIR in the construction industry (US Department

of Labor 2016). However, the short timeframe makes it difficult to support this judgment. It also drives the general question: *how many worker-hours of exposure are needed before TRIR becomes statistically meaningful?*

Companies B and C offer interesting contrast as both accumulate relatively large numbers of worker-hours. Again, on the surface, it may appear that Company C is nearly twice as good as Company B. However, it is still unclear if Company B is *statistically* different from Company C given that injuries do not appear to occur at some regular, predictable interval.

Later in this paper, these three case companies are used to illustrate proper interpretation of TRIR and the implications of the results.

Company A

Has a recordable incident in the first 1,000 worker-hours that they are in business. At this point, their TRIR is 200 per 200,000 worker-hours.

Company B

Has 7 recordable incidents over 980,000 worker-hours in a given year. They report their yearly TRIR as 1.4 per 200,000 worker-hours.

Company C

Has 24 recordable incidents over 6,000,000 worker-hours in a given year. They report their yearly TRIR as 0.8 per 200,000 worker-hours.



Analytical Approach

This study was performed via a collaboration among senior leaders of ten construction companies and four academic researchers. The collaboration resulted in direct access to over 3 trillion worker-hours of internally reported incident data, which were analyzed by the academics using a variety of diagnostic and predictive analytics.

Both parametric and non-parametric analyses were used to study TRIR. A parametric statistical analysis is one that makes logical assumptions about the defining properties of the distributions (i.e., the metric follows a binomial distribution).

Parametric Approach

A parametric analysis starts by identifying the underlying distribution that governs the phenomenon to remove a layer of abstraction and create representative statistical equations. Most common metrics conform to well-known mathematical distributions like Normal, Binomial, or Poisson. When fitting a distribution, it is important to examine the assumptions about the metric to ensure that the chosen distribution is valid. Once a distribution is known, a mathematical function (equation) can be produced that allows us to interpret the precision of a given TRIR, or answer questions like: how much exposure time is required to produce a precise measure of TRIR? For simplicity, we use the term precision to refer to the width of the confidence intervals, where a wide interval is considered to be less precise than a narrow interval.

The parametric analysis revealed very important insight into how TRIR should be communicated. At present, most

TRIR as a Distribution

Two values are needed to compute TRIR: incidents and time. Incidents are discrete events; they happen, or they do not, and there is no such thing as a fractional or negative incident. Time, however, is a continuous value and there are infinite possibilities. To work with time as a variable it is often defined in ranges like days, hours, or minutes. If we look at each worker-hour (the base unit used in TRIR) as a discrete event and an incident as a binary possibility, this yields what is known as a Bernoulli trial.

Non-parametric statistical analyses make no assumptions about the underlying probability distributions, and instead estimate their distributions solely from the data. Both are important for statistical modeling because they help to answer different questions. For example, parametric analyses help us to interpret the precision of TRIR by considering confidence intervals and the non-parametric analyses help us to test whether past TRIR is predictive of future TRIR. Recognizing the significant implications of these results, both approaches were taken to gain a full understanding of when, if ever, TRIR can be used as a comparative or predictive metric.

organizations report TRIR as a single number, often to multiple decimal points (e.g., 1.84). However, as will be shown, TRIR is subject to random variation and should not be communicated as a single number. For example, with the same underlying safety system, an organization would not expect the exact same number of recordable injuries every reporting period. Rather, it may be equally likely that long periods may exist between incidents or that incidents occur in clusters. Additionally, although TRIR is widely reported as a concrete measurement, it is actually a sample taken over a discrete period of time (e.g., one month or one year). To make statements about the possible outcomes of a safety management system, it is important to understand TRIR as a distribution of potential values rather than a single point estimate.

A Bernoulli Trail is Built Upon 3 Assumptions:

Assumption 1 - Each trial results in one of two possible outcomes (occurrence or non-occurrence). Incidents satisfy this assumption because they either occur or do not in one worker-hour. Although it is possible to have more than one recordable incident in a worker-hour, it is so rare as to be mathematically negligible.

Assumption 2 - The probability (p) of an injury remains constant from one worker-hour to the next. It is assumed that a worker is equally likely to be injured in each worker-hour. Although certain situations have a higher probability of injury than others, over enough time exposure each worker-hour can be represented as approximately the same.

Assumption 3 - The trials (worker-hour) are independent. Independence in a dataset means that there is not a well-established connection in the outcomes among trials (worker-hours). For TRIR, there are no known patterns in occurrence, especially when the data are considered in a very high number of trials (e.g., hundreds of thousands of worker-hours). Although there is an argument to be made that an incident in one worker-hour makes an incident in the next more or less likely, independence remains a reasonable assumption.

Computing Confidence Intervals

Now that TRIR is understood to be logically represented as a series of Bernoulli trials, the distribution that best represents the context must be identified. Because recordable incidents can only occur or not occur, the distribution of TRIR must be discrete. In other words, the number of incidents observed over any time period must be a whole number. Although the binomial distribution represents the distribution of potential outcomes from a sample of Bernoulli trials, the Poisson distribution is an even more accurate representation of TRIR.

A Poisson distribution can be thought of as a unique case of the binomial distribution where the probability of occurrence is **very small**. The Poisson distribution is reasonable to use when there are at least 20 trials and the probability of occurrence is less than or equal to 5% (Prins 2012). In the case of TRIR at a

basic unit level (incidents and worker-hours), both are true because incidents are rare and TRIR is typically measured over at least thousands of worker-hours. Another key benefit of the Poisson distribution is that it expresses the probability of a given number of events occurring in a fixed interval of time. This means that the Poisson distribution can be scaled to virtually any time range.

When a value is observed from a distribution of possible observations, the result should be communicated as a confidence interval. A confidence interval is the range of values that is likely to contain the true value with some degree of confidence. The level of confidence is expressed as a probability that the true value lies between an upper and lower interval. For example, TRIR could be expressed as a value that is contained in a range between values X and Y with 95% confidence. This can also be practically interpreted as a 5% probability that a long-term TRIR would be outside this range. Confidence intervals are computed using equations that best represent the underlying distribution. For a Poisson distribution, the Wilson confidence interval is the most appropriate approximation (Wallis 2013). The upper and lower bounds of the Wilson confidence interval are represented by Equation 2 below (Wallis 2013). Although there are other approximations of Poisson confidence intervals that are more accurate (e.g., Garwood 1936; Ulm 1990); the Wilson confidence interval is simple to compute with a basic calculator and also provides an approximation within 1% of the more computationally demanding methods. For example, the 'exact' method from Ulm (1990) requires a statistical package to analyze and produces results that are practically indistinguishable from the Wilson Interval shown in Equation 2.

$$\text{Equation 2: } \frac{\hat{P} + \frac{z^2}{2n}}{1 + \frac{z^2}{n}} \pm \frac{z}{1 + \frac{z^2}{n}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n} + \frac{z^2}{4n^2}}$$

Where, p is the number of actual events (incidents) divided by the number of trials (worker-hours) and z is the critical value of a standard normal distribution corresponding to the target confidence interval (e.g., $\alpha = 0.05$ for a 95% confidence interval). In this analysis a significance level (α) of 0.05 was always selected so the corresponding z is 1.96.

The Wilson confidence interval allows one to judge the precision of an observed TRIR using only two pieces of information: the number of incidents and the number of worker-hours in the sample. Once this information is known, the confidence interval can be approximated using Equation 2. Multiplying the endpoints by 200,000 shows the range in terms of TRIR as it is normally reported.

An example provides some clarity for the layperson. If a theoretical company had 1 recordable injury in 200,000 worker-hours ($n=200,000$, $p=1/200,000$), we can calculate the 95% confidence interval using Equation 2 as 0.18 to 5.66 injuries per 200,000 worker-hours. Theoretically, this corresponds to the range of results that the company's safety system is designed to produce that month. A TRIR below 0.18 would be interpreted as unusually low and a TRIR above 5.66 as unusually high. Therefore, reporting a TRIR of 1.00 per 200,000 is not appropriate or meaningful. Rather, the TRIR should be reported as an interval

of 0.18 to 5.66 with the most likely true value of 1.00. We can also ask the complementary question: Assuming the company's true injury rate is 1 per 200,000 worker-hours, how many injuries are likely to occur over future periods of 200,000 worker-hours if the safety system remains the same? Figure 1 shows the range of potential results and their probabilities.

As an extension of the above example, if the same company experienced 2 injuries in the next 200,000 worker-hours, the 95% interval would then be 0.55 to 7.29 per 200,000 worker-hours. The company might be concerned that they doubled their injuries from the previous interval from 1.00 to 2.00 per 200,000 worker-hours. However, a test for significance shows that there is no statistical difference between the months even though the number of injuries doubled. This result means that the difference in the count of injuries alone does not reveal anything significant about the difference in the safety system between the two periods.

Probability of Experiencing N Incidents in 200,000 Worker-Hours

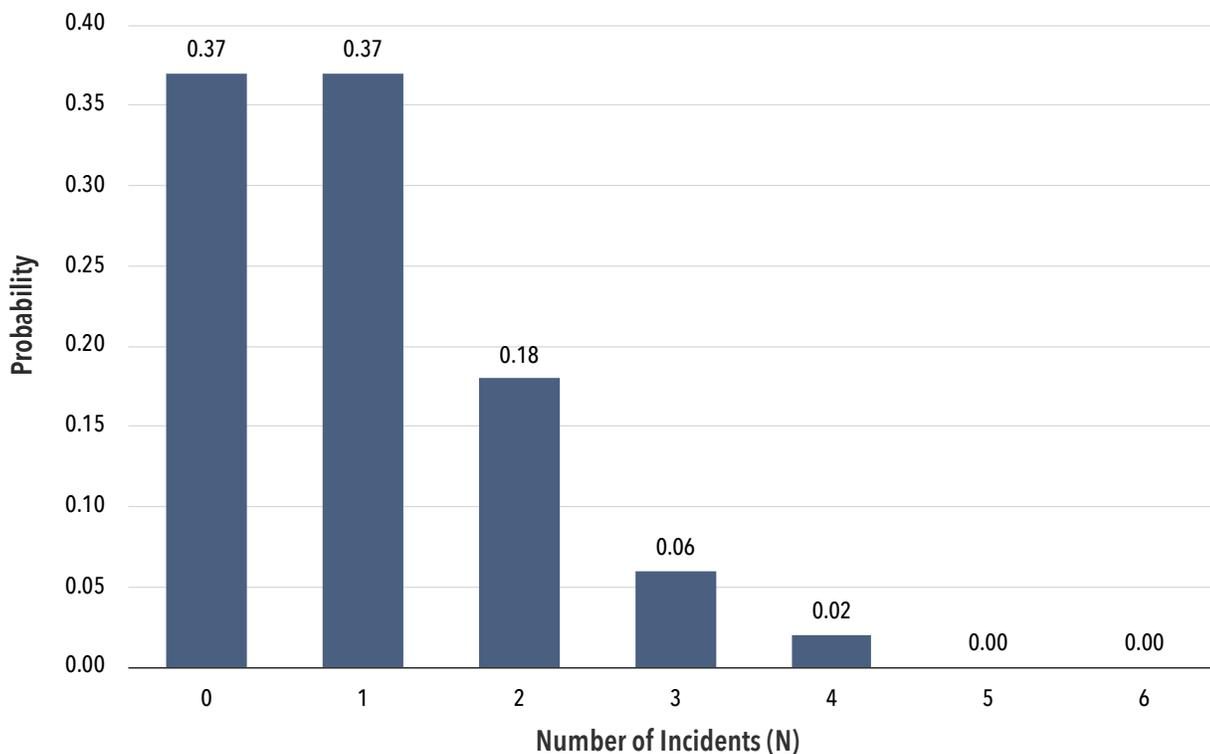


Figure 1 - Probability of experiencing N incidents in the next 200,000 worker-hours assuming a true TRIR of 1 recordable injury in 200,000 worker-hours.

Analysis of Case Examples

Using the three case examples from the background (Companies A, B, and C), we can explore how reporting TRIR as a confidence interval vastly changes its interpretation and meaning. For example, the TRIR of company A, which had a TRIR of 200 resulting from one recordable injury in 1,000 worker-hours, would be correctly reported as 35 to 1,128 with 95% confidence. This statistically confirms the logical interpretation that this new small company does not have enough exposure time to return a meaningful TRIR. The comparison of Companies B and C is also interesting. Since both companies reported their TRIR over many worker-hours, it may seem appropriate to recognize their difference in TRIR as meaningful. However, as shown in Figure 2, the distribution of Company C fits entirely within the distribution of Company B indicating that the two injury records are statistically indistinguishable. That said, another important conclusion can be made: the TRIR for Company C can be stated with a much higher precision (smaller interval) than Company B. Therefore, there is distinct advantage to having a greater number of worker-hours accumulated. However, it still remains unclear whether the historical TRIR has any bearing on future TRIR, which was the primary subject of the non-parametric empirical analysis.

Company A

Has a recordable incident in the first 1,000 worker-hours that they are in business.

TRIR range: 35.31 to 1,125.51 per 200,000 worker-hours.

Company B

Has 7 recordable incidents over 980,000 worker-hours in a given year.

TRIR range: 0.69 to 2.95 per 200,000 worker-hours.

Company C

Has 24 recordable incidents over 6,000,000 worker-hours in a given year.

TRIR range: 0.54 to 1.19 per 200,000 worker-hours.

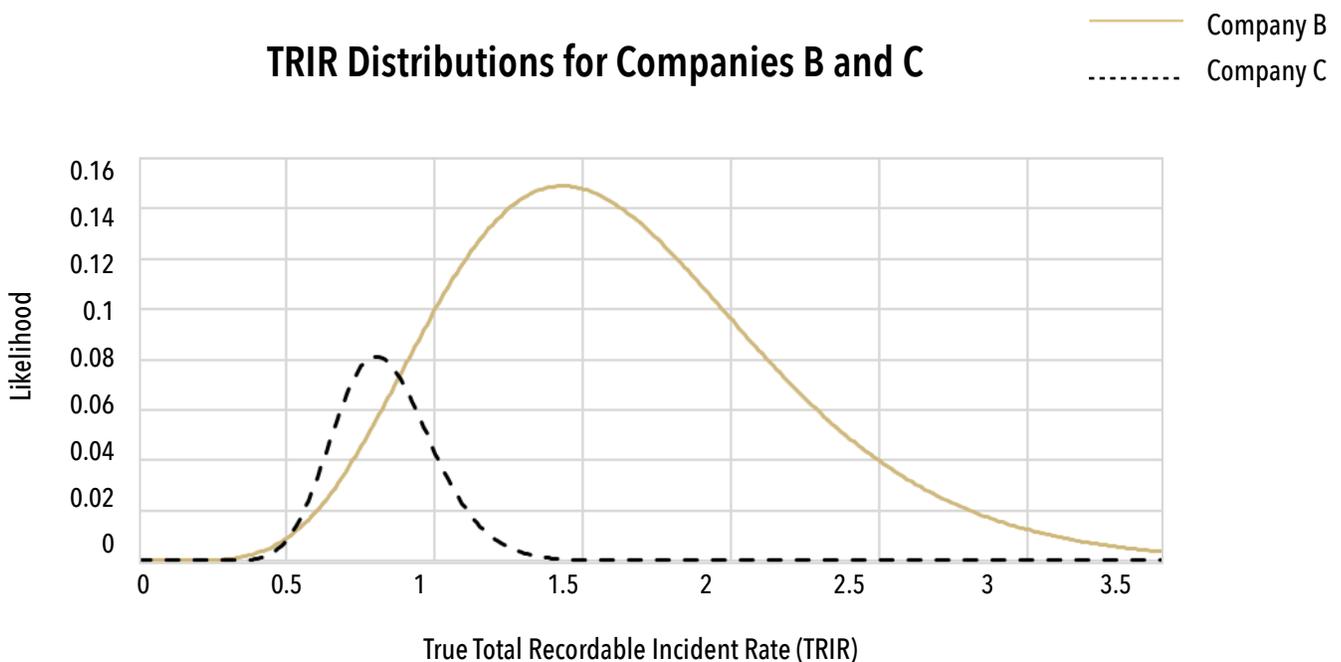


Figure 2 - Likelihood of the true injury rates for Companies B and C (per 200,000 hours), based on their observed histories.

To assist practitioners with the interpretation of TRIR as a range, Tables 1 and 2 were created. Table 1 shows the relationship between precision (width of the confidence interval), time (worker-hours), and TRIR. Simply, it shows how many worker-hours are needed to achieve a given precision level for a variety of TRIR scenarios. The number of hours needed for a given precision can be approximated by Equation 3 below, which is based upon the Wilson binomial formula (Krishnamoorthy and Peng 2007).



Equation 3:

$$n \cong \frac{z_a^2 \left[(Pq - 2d^2) + \sqrt{(Pq - 2d^2)^2 - d^2(4d^2 - 1)} \right]}{2d^2}$$

$$d = \frac{\text{Desired TRIR Precision}}{200,000 * 2}$$

Where n is the required number of worker-hours to achieve precision d given a probability of injury in each hour p, and q=1-p.

Table 1 reveals some potentially surprising results. For example, if an organization wishes to report a TRIR of 1.00 with a precision of 0.1 (e.g., 0.95 to 1.05 per 200,000 worker-hours), about 300 million worker-hours of exposure time is required. In other words, unless the TRIR of 1.0 is derived from 300 million worker-hours of exposure time, it should not be reported to even one decimal point. Furthermore, reporting two meaningful decimal places for TRIR requires approximately 30 billion worker-hours of data. The required number of worker-hours is so high because recordable injuries occur, on average, so infrequently that they do not produce statistical stability. This raises questions about reporting TRIR to two decimal points and making important business decisions with this level of granularity.

To further illustrate how TRIR should be interpreted, Table 2 shows the 95% confidence interval for a series of scenarios with varying TRIR and number of worker-hours. For example, Table 2 shows that if an organization measured their TRIR to be 1.0 over a period of 1,000,000 worker-hours (i.e., 5 recordable injuries occurred over the span of 1,00,000 worker-hours), the correct interpretation is that the safety system is designed to produce a TRIR between 0.43 and 2.34. Note that in Table 2 some scenarios are not possible because they would not correspond to a whole number of incidents; however, they are included to maintain continuity.



Non-Parametric Approach

To complement the parametric analysis and to add a degree of validation, an empirical analysis was performed with injury data provided by 10 construction organizations that are members of the Construction Safety Research Alliance (CSRA). The partner organizations represented infrastructure, power generation and delivery, oil and gas, industrial and energy, and commercial building sectors. The partners provided monthly counts of recordable injuries, fatalities, and worker-hours for a 15-year period. In total, the dataset included 3.26 trillion worker-hours of data. The empirical analysis of these data is summarized below, and the analytical approach is explained in fine detail in Salas (2020).

Random & Unpredictable Nature of TRIR

To assess the significance of TRIR, the data for each organization were fit to a generalized linear model (GLM) using different distribution functions to determine the best model (i.e., the one with the lowest Akaike Information Criterion). The best fit model was then subjected to 100 repeated sampling and 10-fold cross-validation to measure the strength, stability, and response of the model to the observed data. Then, a Monte Carlo Simulation was used to evaluate how much of the final result is due to random variation. Holding the data generating process constant, the Monte Carlo Simulation tested how different estimators fare in trying to uncover underlying parameters. Put simply, this method involved using the historical TRIR data for each organization to create a predictive equation and then testing how well that equation correctly estimated future TRIR.

The results indicated that TRIR is 96-98% random. This means that the best model was only able to predict the observed TRIR in 2-4% of the trials. Since safety is a chaotic system due to the interfaces between people, culture, policies, regulations, equipment and other external factors (e.g. economy, weather, natural events, etc.), this is not a surprising finding. However, it does have very important implications. First, it provides empirical evidence to support the logical assumptions made in the parametric analysis. Second, the fact that TRIR is random in nature provides further evidence that it must be reported as a

range (i.e., confidence interval) and absolutely should not be expressed as a single number.

In addition to measuring the significance, the models were also tested to assess the extent to which historical TRIR predicts future TRIR. The results showed that at least 100 months of data were required to achieve reasonable predictive power because of the high degree of random variation. Since TRIR is generally used to make comparisons or decisions on the order of months or years, this finding indicates that for all practical purposes, TRIR is not predictive. For example, a client that hires a contractor with a TRIR of 0.75 cannot reasonably expect that the contractor will achieve that same performance on their upcoming project.

Lack of Relationship Between TRIR & Fatalities

Effect measurements revealed that variation in TRIR has no association with fatalities. That is, trends in TRIR do not associate statistically with fatality occurrence. Instead, fatalities appear to follow different patterns, suggesting that they occur for different reasons. This finding challenges the long-standing assumption, derived from the Heinrich safety pyramid, that injuries of different severity levels exist at fixed ratios and have the same underlying causes (Heinrich 1959). It also debunks the notion that reducing TRIR is a surrogate for mitigating the risk of high-impact events.



Table 1. Relationships among Precision (width of confidence interval), TRIR, and Exposure Time (worker-hours)

Precision	TRIR	Worker-Hours	Precision	TRIR	Worker-Hours	Precision	TRIR	Worker-Hours
0.1	0.20	62,409,083	0.25	0.20	10,715,491	0.5	0.20	3,197,046
0.1	0.40	123,404,751	0.25	0.40	20,137,222	0.5	0.40	5,357,737
0.1	0.60	184,709,031	0.25	0.60	29,819,030	0.5	0.60	7,682,889
0.1	0.80	246,092,232	0.25	0.80	39,575,018	0.5	0.80	10,068,588
0.1	1.00	307,507,116	0.25	1.00	49,361,749	0.5	1.00	12,481,762
0.1	1.25	384,297,072	0.25	1.25	61,616,226	0.5	1.25	15,517,884
0.1	1.50	461,099,608	0.25	1.50	73,883,277	0.5	1.50	18,566,029
0.1	1.75	537,909,256	0.25	1.75	86,157,534	0.5	1.75	21,621,176
0.1	2.00	614,723,280	0.25	2.00	98,436,299	0.5	2.00	24,680,748
0.1	3.00	922,000,301	0.25	3.00	147,573,796	0.5	3.00	36,941,358

Table 2. 95% Confidence Intervals for a Series of TRIR Scenarios

	Total Recordable Incident Rate								
	0.1	0.2	0.5	1.0	1.5	2.0	3.0	4.0	5.0
100K	0.00 - 7.68	0.00 - 7.68	0.00 - 7.68	0.00 - 7.68	0.00 - 7.68	0.35 - 11.33	0.35 - 11.33	1.10 - 14.59	1.10 - 14.59
250K	0.00 - 3.07	0.00 - 3.07	0.00 - 3.07	0.14 - 4.53	0.14 - 4.53	0.44 - 5.83	0.82 - 7.06	1.71 - 9.36	2.20 - 10.47
500K	0.00 - 1.54	0.00 - 1.54	0.07 - 2.27	0.22 - 2.92	0.41 - 3.53	0.85 - 4.68	1.36 - 5.78	2.17 - 7.36	2.75 - 8.39
1M	0.00 - 0.77	0.04 - 1.13	0.11 - 1.46	0.43 - 2.34	0.68 - 2.89	1.09 - 3.68	1.82 - 4.95	2.59 - 6.18	3.39 - 7.38
2.5M	0.01 - 0.45	0.04 - 0.58	0.22 - 1.05	0.55 - 1.68	0.91 - 2.28	1.35 - 2.95	2.15 - 4.08	3.03 - 5.27	3.87 - 6.36
5M	0.02 - 0.29	0.09 - 0.47	0.27 - 0.84	0.68 - 1.48	1.07 - 2.04	1.52 - 2.64	2.39 - 3.76	3.29 - 4.86	4.20 - 5.96
10M	0.04 - 0.23	0.11 - 0.37	0.34 - 0.74	0.76 - 1.32	1.20 - 1.88	1.64 - 2.43	2.56 - 3.52	3.48 - 4.59	4.42 - 5.66
20M	0.05 - 0.18	0.13 - 0.31	0.38 - 0.66	0.82 - 1.22	1.28 - 1.76	1.74 - 2.30	2.68 - 3.36	3.63 - 4.41	4.58 - 5.46
50M	0.07 - 0.15	0.15 - 0.26	0.42 - 0.60	0.88 - 1.13	1.36 - 1.66	1.83 - 2.18	2.79 - 3.22	3.76 - 4.26	4.73 - 5.28

Key Findings

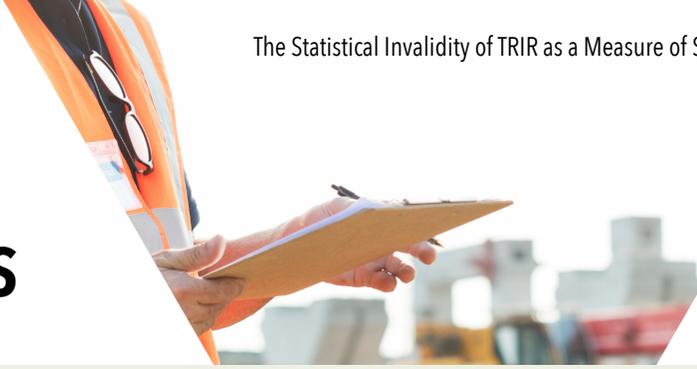


This study challenges conventional wisdom of safety measurement with an empirical analysis of 3.26 trillion worker-hours of TRIR data and a statistical demonstration of proof. The results reveal very strong evidence that TRIR is almost entirely random and is not indicative of future performance unless millions of worker-hours are amassed. The specific conclusions that follow logically and empirically are as follows:

1. **TRIR is not associated with fatalities.** The effect measurements revealed that there is no discernible association between fatalities and TRIR. Recordable injuries and fatalities follow different patterns and occur for different reasons. Thus, TRIR trends are not a proxy for high-impact incidents. Hence, it can be inferred that safety interventions, including policies, regulations and management systems, associated with the improvement in TRIR performance may not necessarily prevent fatalities.
2. **TRIR is almost entirely random.** Empirical analysis revealed that changes in TRIR are due to 96-98% random variation. This is logically confirmed by the fact that recordable injuries do not occur in predictable patterns or regular intervals. This is likely because safety is a complex phenomenon that is impacted by many factors.
3. **TRIR cannot be represented as a single point estimate.** Since TRIR is almost entirely random, a single number does not represent the true state of safety performance. Instead, TRIR is best expressed as a confidence interval and studied over extended periods of time. For example, a yearly TRIR value of 1.29 is statically meaningless for almost every organization. This finding was initially explored in the parametric analysis and was empirically validated in the non-parametric analysis.
4. **TRIR is not precise and should not be communicated to multiple decimal points.** Unless hundreds of millions of worker-hours are amassed, the confidence bands are so wide that TRIR cannot be accurately reported to even one decimal point. The implication is that the TRIR for almost all companies is virtually meaningless because they do not accumulate enough worker-hours.
5. **If an organization is using TRIR for performance evaluations, they are likely rewarding nothing more than random variation.** Because of the random nature of TRIR, it is unclear if a change in performance (positive or negative) is due to an underlying change in the safety system or if the organization is simply observing random variation.
6. **TRIR is predictive only over very long periods of time.** Previous researchers have postulated that TRIR can be used predictor of performance when taken over very long time periods (Alruqi and Hallowell 2019; Lingard et al. 2017; Salas and Hallowell 2016; Wilbanks 2019). The empirical results of this study confirmed that TRIR is only predictive when over 100 months of TRIR data are accumulated.

The results of this study may not be surprising to some professionals who have made these postulations for years. However, this is the first scientific evidence that explains why TRIR is not a valid comparative measure of safety performance.

Conclusion & Recommendations



The conclusions of this research may disrupt how the profession approaches safety measurement and reporting. Therefore, along with the senior executives from our industry partners, we offer the following practical interpretation of the results.

- 1. TRIR should not be used as a proxy for serious injuries and fatalities.** For years, many practitioners have used declining trends in TRIR to indicate the mitigation of fatality risk. However, the lack of statistical association between TRIR and fatalities suggests that the assumption holds no scientific merit. Practitioners should consider creating and testing targeted measurement, learning, and prevention efforts for serious injuries and fatalities.
- 2. TRIR should not be used to track internal performance or compare companies, business units, projects, or teams.** Since the average company requires tens of millions of worker-hours to return a confidence interval with one decimal point of precision, organizations should be very careful making *any* comparisons using TRIR. For example, most companies do not even have enough worker-hours to detect statistically significant changes in TRIR from one year to the next. Therefore, we challenge practices where TRIR is used to compare companies, projects, or teams. At best, TRIR is only useful for comparing industries or sectors of the US economy over long

periods of time. For the same reason, TRIR should not be used as the primary safety metric when incentivizing organizational performance or comparing or prequalifying contractors.

- 3. The safety profession must change how it communicates TRIR.** TRIR is almost always communicated as a precise number as if it was the only possible outcome (Stricoff 2000). Since recordable injuries are so infrequent and are a product of so much random variation, a single precise number is meaningless. Instead, TRIR should be accompanied by the range of potential outcomes that the safety system could have reasonably produced. The implication is that small companies would report large confidence intervals (high uncertainty) and large companies would report smaller confidence intervals (lower uncertainty). However, almost no company would be able to appropriately report TRIR to the level of precision most commonly used today.
- 4. Companies should not place much emphasis on short term changes in TRIR.** It is tempting to track TRIR over time to identify when performance is improving or degrading. However, as observed in the empirical analysis, changes that occur from month-to-month are mostly random and do not necessarily reflect any actual change in the safety system.

5. TRIR should not be used to measure the impact of safety interventions. Managers are often compelled to show reductions in TRIR resulting from safety initiatives and investments. However, controlled experiments and longitudinal data are needed to establish causal inference between safety interventions and TRIR trends. This renders TRIR entirely inadequate for attributing change because most companies drastically shift their management approaches, safety programs, and even their business models over the long timeframes needed to produce a precise and stable TRIR.



6. New approaches to safety measurement are needed. In addition to statistical invalidity, the use of TRIR also does not describe *why* the performance – good or bad – was achieved and what can be done to improve. This leaves organizations wondering, ‘Are we truly good, or simply lucky?’ or worse, ‘Are we truly bad, or do we simply need to log more worker-hours?’ The academic and professional community should consider alternative measures of safety performance that assess the actual safety system at high frequency. Increasing the number of reliable measurements could drastically improve the stability, precision, and predictive nature of safety metrics. To be comparative, however, these metrics must be standardized and consistently reported.

Since TRIR has remained the most pervasive measure of safety for nearly fifty years, this study underscores the need to scientifically test even the most basic assumptions of the safety profession. In the spirit of scientific inquiry, we recommend that other researchers propose alternative hypotheses about TRIR, conduct independent tests, and challenge the assumptions made in this paper. Although we stand by our conclusions, we recognize that other perspectives may generate different models and results.

Acknowledgements



This research was made possible by the generous funding provided by the members of the Construction Safety Research Alliance (CSRA) and through the participation, data, and leadership of the following members of the CSRA Board of Advisors:

- Brad MacLean: Wolfcreek Group
- Ellen Quinn: Otis Elevator
- Jenny Bailey: Xcel Energy
- Ester Brawley: California Resources Corporation
- Bruno Cagnart: TechnipFMC
- Len Colvin: Tennessee Valley Authority
- Matt Compher: Quanta Services
- Mike Court: Graham Group Ltd
- Nick DiMartino: MasTec
- Phillip Hodge: Chevron
- Greg Kelly: Enbridge
- Joe Parham: Alabama Power
- Mike Quashne: Baltimore Gas & Electric
- Gregg Slintak: Consolidated Edison Company

CSRA

References

- Alruqi, W. M., and Hollowell, M. R. (2019). "Critical Success Factors for Construction Safety: Review and Meta-Analysis of Safety Leading Indicators." *Journal of Construction Engineering and Management*, 145(3), 04019005.
- Fagan, K. M., and Hodgson, M. J. (2017). "Under-recording of Work-related Injuries and Illnesses: An OSHA Priority." *Journal of Safety Research*, Elsevier Ltd, United States, 60, 79-83.
- Garwood, F. (1936). Fiducial limits for the Poisson distribution, *Biometrika*, 28, 437-442.
- Hollowell, M. (2010). "Cost Effectiveness of Construction Safety Programme Elements." *Construction Management and Economics*, 28(1), 25-34.
- Hollowell, M. R., and Gambatese, J. A. (2009). "Construction Safety Risk Mitigation." *Journal of Construction Engineering and Management*, 135(12), 1316-1323.
- Hollowell, M. R., Hinze, J. W., Baud, K. C., and Wehle, A. (2013). "Proactive Construction Safety Control: Measuring, Monitoring, and Responding to Safety Leading Indicators." *Journal of Construction Engineering and Management*, 139(10), 4013010.
- Heinrich, H. W. (1959). "Industrial Accident Prevention: A Scientific Approach." McGraw-Hill, New York, NY, USA.
- Hinze, J., Thurman, S., and Wehle, A. (2013). "Leading Indicators of Construction Safety Performance." *Safety Science*, Elsevier Ltd, 51(1), 23-28.
- Krishnamoorthy, K., and Peng, J. (2007). "Some Properties of the Exact and Score Methods for Binomial Proportion and Sample Size Calculation." *Communications in Statistics - Simulation and Computation*. 36: 1171-1186.
- Lingard, H., Hollowell, M., Salas, R., and Pirzadeh, P. (2017). "Leading or lagging? Temporal analysis of safety indicators on a large infrastructure construction project." *Safety Science*, 91.
- Lofquist, E. A. (2010). "The Art of Measuring Nothing: The Paradox of Measuring Safety in a Changing Civil Aviation Industry using Traditional Safety Metrics." *Safety Science*, 48(10), 1520-1529.
- Prins, J. (2012), "6.3.3.1. Counts Control Charts." e-Handbook of Statistical Methods, NIST/SEMATECH, <<http://www.itl.nist.gov/div898/handbook/pmc/section3/pmc331.htm>> (Sept. 20, 2019).
- Ritchie, N. W. (2013). "Journey to Zero: Aspiration Versus Reality." SPE European HSE Conference and Exhibition, 16-18.
- Salas, R. (2020). "Empirical Evaluation of Construction Safety: Insights from Economic Trends, National Cultures and Workplace Injury Rates." ProQuest Dissertations Publishing, Boulder, CO.
- Salas, R., and Hollowell, M. (2016). "Predictive Validity of Safety Leading Indicators: Empirical Assessment in the Oil and Gas Sector." *Journal of Construction Engineering and Management*, 142(10).
- Stricoff, R. S. (2000). "Safety performance measurement: Identifying prospective indicators with high validity." *Professional Safety*, 45(January), 36-39.
- Toellner, J. (2001). "Improving Safety and Health Performance: Identifying and Measuring Leading Indicators." *Professional Safety*, 42-47.
- Truitt, D. (2012). "Contractor Prequalification." *Professional Safety*, (March), 34-36.
- Ulm, K. (1990). Simple method to calculate the confidence interval of a standardized mortality ratio (SMR). *American Journal of Epidemiology*, 131(2), 373-375.
- US Bureau of Labor Statistics. (2019). "Injuries, Illnesses and Fatalities." <<https://www.bls.gov/iif/>>.
- US Department of Labor. (2010). "Clarification of assigned working hours when recording work-related injuries/illnesses." *Standard Interpretations*, <<https://www.osha.gov/laws-regs/standardinterpretations/2010-02-16>> (Nov. 19, 2019).
- US Department of Labor. (2016). "Clarification on how the formula is used by OSHA to calculate incident rates." *Standard Interpretations*, <<https://www.osha.gov/laws-regs/standardinterpretations/2016-08-23>> (Nov. 24, 2019).
- Wanberg, J., Harper, C., Hollowell, M. R., and Rajendran, S. (2013). Relationship between construction safety and quality performance. *Journal of Construction Engineering and Management*, 139(10), 04013003.
- Wardrop, R. (2010). "Bernoulli Trials." Course Notes, Wisconsin University Statistics 371, <<http://pages.stat.wisc.edu/~wardrop/courses/371achapter2a.pdf>> (Sept. 20, 2019).
- Wallis, S. A. (2013). "Binomial confidence intervals and contingency tests: mathematical fundamentals and the evaluation of alternative methods." *Journal of Quantitative Linguistics*. 20 (3): 178-208.
- Wilbanks, D. W. (2019). "Contractor Safety Prequalification." *Professional Safety*, (July), 22-26.

Authors



Dr. Hallowell is the Executive Director of the Construction Safety Research Alliance (CSRA). He is also a President's Teaching Scholar and Endowed Professor of Construction Engineering at the University of Colorado.

Dr. Matthew Hallowell



Mike is a member of the CSRA Board of Advisors and a Manager of Business Transformation and Performance Assessment for Baltimore Gas and Electric. Mike has been a safety and human performance practitioner since 2016 in the utility sector. He holds a Bachelor of Science in Aeronautical Engineering from the United States Air Force Academy, a Master of Science in Operations Research from the Air Force Institute of Technology, and is a Certified Utility Safety Professional (CUSP).

Mike Quashne



Dr. Salas is a Research Associate with the CSRA. He received his MS degree from the University of Cambridge and a PhD from the University of Colorado at Boulder. He has an extensive technical and work experience in design, construction and process safety in the oil and gas industry. He is currently the Asset Integrity Engineering Manager in one of the leading US oil and gas companies.

Dr. Rico Salas



Dr. Matt Jones

Dr. Jones is an Associate Professor of Psychology and Neuroscience and the Director of the Center for Research on Training. His research develops mathematical and statistical methods for understanding human behavior.



Brad MacLean

Brad MacLean is Chair of the CSRA and the Senior Vice President of Safety for the Wolfcreek Group. Brad has been a Health, Safety, and Environment practitioner in the Energy Sector since 1990, working in Canada and the United States. His experience extends to construction, operations, corporate governance and Board service. Brad is a founding member and past Chair of the Interstate Natural Gas Association of America (INGAA) Foundation Safety Committee.



Ellen Quinn

Ellen Quinn is the Vice Chair of the CSRA and was Vice President of Environmental, Health, and Safety for United Technologies Corporation (UTC) from 2016 to 2020. She received a BA from Franklin Marshall College, and MA from Harvard University, and an MS from Rensselaer Polytechnic University. She received the international Robert W. Campbell award for EH&S excellence in 2011. She has served on non-profit boards in the environmental and arts fields, including Trust for Public Land and Pilobolus.



© 2020, CSRA-The Construction Safety Research Alliance. All rights reserved. CSRA, UCB 428, 1111 Engineering Drive, Boulder, CO, 80309-0428
Production of this publication was supported by APlus Design, LLC, 20070 SW Newcastle Drive, Aloha, OR, 97078