

1 **A data-driven approach for deploying safety policies for schedule planning in industrial**
2 **construction projects: a case study**

3 Maedeh Taghaddos¹; Estacio Pereira²; Carlos Osorio-Sandoval³; Ulrich Hermann⁴; Simaan

4 AbouRizk^{5*}

5 ¹ PhD Graduate, Department of Civil and Environmental Engineering, University of Alberta, Donadeo
6 Innovation Centre for Engineering, 9211 116 Street NW, Edmonton, Alberta, Canada, T6G 1H9; Email:
7 taghaddo@ualberta.ca

8 ² Assistant Professor (Teaching), Department of Civil Engineering, University of Calgary, ENF 208, 2500
9 University Dr NW, Calgary, Alberta, Canada, T2N 1N4 Email: estacio.pereira@ucalgary.ca

10 ³ Assistant Professor, Department of Civil Engineering, University of Nottingham, B74 Coates Building,
11 University Park, NG7 2RD, Nottingham, United Kingdom Email: carlos.osorio@nottingham.ac.uk

12 ⁴ Manager, Construction Engineering, PCL Industrial Management Inc., 9925 56 Ave., Edmonton, AB,
13 Canada T6E 3P4; Email: rhhermann@pcl.com

14 ⁵ Professor, Department of Civil and Environmental Engineering, University of Alberta, 5-080 NREF,
15 9105 116 St NW, Edmonton, AB T6G 2W2, Email: abourizk@ualberta.ca

16 *Corresponding Author; 5-080 NREF, University of Alberta, Edmonton, Alberta, Canada T6G 2W2

17

18 **ABSTRACT**

19 Construction, by nature of its work, is more accident-prone than other industries despite advancements in
20 improving safety performance. Proactive mitigation and assessment of safety performance on construction
21 projects remain challenging due to the difficulty of acquiring, storing, and using data to produce accurate
22 predictive models. This research is focused on devising methods that allow decision-makers to leverage
23 existing data in the planning phase to streamline the development of predictive models. A data-driven
24 approach to predict the probability of a safety incident occurring in a given construction project and within
25 a novel discipline-level schedule is presented. By implementing the proposed model, decision-makers can
26 evaluate and mitigate the risk of a given project incident occurring by deploying discipline-level safety
27 policies in the planning phase and modifying the schedule accordingly. A predictive model is developed
28 based on selected safety-related metrics extracted from a dataset comprising daily payroll data and incident
29 reports, which represent 28 million working hours within eight different industrial construction projects in

30 Canada. The model was implemented in a case study based on an industrial project to demonstrate the
31 framework's functionality and practical utility during the project planning phase. The results show that the
32 revised safe plan can be achieved by incorporating safety considerations in the planning phase.

33 **PRACTICAL APPLICATIONS**

34 This research provides a practical solution for enhancing safety in the planning phase of
35 construction projects through a data-driven model. By leveraging existing historical data, decision-
36 makers can predict potential safety incidents within specific disciplines without the need for
37 detailed quantitative planning information. This approach also enables effective adjustments to be
38 made to the schedule in order to mitigate risks. Furthermore, the discipline-level approach
39 facilitates proactive safety planning by implementing discipline-specific safety policies that align
40 with the unique characteristics of each discipline.

41 Using a case study based on an industrial project, the proposed framework demonstrates its
42 functionality and practical utility by identifying suitable safety-related metrics that construction
43 enterprises typically record. These sources can include safety-related data, such as incident reports,
44 as well as data recorded for other purposes, such as payroll data. The results highlight that
45 incorporating safety considerations in the planning phase enables the development of a revised
46 safety plan. In conclusion, the key takeaway is that by considering safety-related metrics and
47 utilizing HR data available in all companies, organizations can proactively assess and improve
48 safety performance.

49 **INTRODUCTION**

50 According to the Association of Workers' Compensation Boards of Canada (2019), the number of
51 accepted lost-time claims per year in the construction industry has increased by 9% from 2016
52 (25,514) to 2018 (27,952). Additionally, the number of fatalities in the Canadian construction

53 industry is higher than in any other industry, with 202 deaths recorded in 2016 and 199 in 2018
54 (AWCBC 2019). According to estimates from the International Labor Organization (ILO), an
55 average of 4% of annual global Gross Domestic Product (GDP)—adding up to trillions of
56 dollars—is lost due to direct and indirect costs incurred as a result of occupational accidents and
57 diseases, which includes lost working time, workers' compensation, interruption of production,
58 and medical expenses (Takala et al., 2014). Although AFPM (2017) demonstrated that the incident
59 rates had decreased greatly in the heavy industrial sector when measured over the past 30 years
60 (Recordable Incident Rate reduced by more than 1000%), there is still great concern about events
61 with a serious injury or fatal consequences as they have remained consistent. Due to the
62 construction industry's impact on Canada's economy, the improved safety performance of
63 construction projects is critically needed.

64 Based on the literature, there are many causes for accidents at construction sites, and a holistic
65 approach has to be considered to mitigate hazards (Ahn et al. 2020; Mohammadi et al. 2018;
66 Pereira et al. 2020). Decision-makers should be aware of how they can impact project safety
67 performance. For instance, poor project scheduling can impact safety performance by 1)
68 generating delays and increasing the pressure on workers (Han et al. 2014; Mitropoulos et al.
69 2005); 2) developing site congestion — increasing the probability of workers being struck-by
70 heavy equipment (Ahn et al. 2020; Zhang et al. 2018); 3) increasing the number of inexperienced
71 people on, and 4) increasing the crew size and consequently reducing the number of inspections
72 and safety observations (Jiang et al. 2015).

73 Although there is no doubt about how different aspects of a project can impact the occurrence of
74 incidents, it is still challenging to integrate decisions from various project factors to mitigate
75 incidents and assess safety performance proactively. Issues may arise because construction

76 practitioners have difficulty identifying safety measures, collecting data, and proactively assessing
77 the impact of factors on safety performance (Pereira et al. 2018a).

78 From a theoretical perspective, three levels can be considered for safety assessment: the first is
79 project level, the second is discipline level, and the third is crew level. The project level has the
80 least details (broadest) to plan for safety, and the crew level has the most details (finest). The most
81 common approach is assessing safety performance at a project level, which may not be feasible in
82 practice as there is insufficient detailed information available for safety planning (Goh and Chua
83 2013; Karakhan et al. 2018; Lingard et al. 2017; Salas and Hallowell 2016; W Guo and Wing Yiu
84 2016). Therefore, practitioners find making decisions based on project-level models challenging
85 since their results may neither be economical nor produce the required output. For example, if a
86 project has a large crew size, reducing the crew size in a specific discipline may not necessarily
87 improve the project's safety performance as much as another measure (e.g., workers' age). In other
88 words, evaluating the crew size for each discipline is essential rather than the overall evaluation in
89 the entire project context. Hence, discipline-based evaluations are preferred rather than project-
90 based evaluations.

91 Another practical challenge is related to data acquisition and utilization. Small companies have
92 difficulty collecting safety-related data (e.g., available resources (Bavafa et al. 2018)). Some
93 companies may not have the knowledge or skills required to produce safety predictive assessment
94 models (Boon et al. 2019) using techniques like Machine Learning (ML) algorithms. Large
95 organizations may have issues with data silos — the data repository — as it is not easily accessible
96 to the entire organization. Although Pereira et al. (2020) demonstrated that integrating different
97 department databases may lead to the development of better models to predict safety, concerns
98 about privileges and sensitive information hinder database integration in practice. Moreover,

99 departments may collect information from different levels (e.g., project, discipline, and task),
100 which may make it impossible to integrate all into one single model. This study proposes to use
101 Human Resources (HR) data (e.g., payroll), a common database among all companies, to address
102 the data acquisition problem.

103 The HR dataset is one of the most reliable datasets, with considerable attention given to it because
104 of its monetary aspect. Moreover, it is very informative and is often made available due to
105 regulations. Several safety-related measures can be identified from this database, such as workers'
106 age and experience, job type, number of workers, new workers' rate (i.e., the rate at which new
107 workers are deployed to work on the project), the number of supervisors, etc. This manuscript
108 presents a case study that uses the payroll database in a data-driven approach for deploying safety
109 policies for discipline-level schedule planning. The approach described herein was able to identify
110 several safety-related measures that companies can use to proactively measure safety performance
111 at the discipline level without relying on subjective decisions.

112 Based on these measures, a predictive safety (or incident) assessment model is suggested, and a
113 case study on using the findings in practice is presented. The approach presented in this manuscript
114 may help organizations use their payroll database to derive safety models to proactively test
115 scenarios and control the safety performance on construction projects.

116 The novelty of this paper is to **proactively** improve safety in a **discipline-level schedule** (in the
117 planning phase). This prediction is through the **novel data acquisition** from a reliable, unbiased,
118 and objective payroll dataset, which is commonly available in all companies. Moreover, this
119 research can facilitate the safety forecast by planning based on qualitative data to evaluate the
120 safety of various possible what-if scenarios. Therefore, the practitioners can consider the safety of
121 the schedule in the planning phase.

122 LITERATURE REVIEW

123 Deficient project scheduling is considered a root cause of construction accidents as it leads to time
124 pressure on workers, with subsequent problems including trade overlap, crowded workspaces, and
125 reduced attention to detail (Haslam et al. 2005; Neale and Gurmu 2021). It is believed that safety
126 can be improved by considering activity information regarding the number of workers, including
127 their occupation types, in the schedule (Choe and Leite 2017), and by minimizing the number of
128 workers on congested sites (Anvari et al. 2016). However, empirical data — needed to estimate
129 risks accurately based on these constraints — is usually not captured, leading to a decision-making
130 process that relies on subjective opinions. In addition, attributes such as the fundamental
131 characteristics of the work site and environment (e.g.; weather, uneven surfaces, specific tools,
132 and equipment (Hallowell et al. 2020)), as well as features of the workforce (e.g.; age, experience,
133 crew size, and specific trade) can also contribute to safety incidents. The combination of these
134 attributes that define the overall work environment can be used to predict safety outcomes
135 (Hallowell et al. 2020; Tixier et al. 2016a).

136 Construction companies typically record some information on several safety indicators at the
137 project level (e.g., site inspection logs, hazard reports, injury reports, etc.) to meet the regulator's
138 requirements (Versteeg et al. 2019). Usually, this information is only provided to the safety
139 department and is not used for predictive purposes (Pereira et al 2020). However, companies still
140 need to be proactive and incorporate measurements of valid and reliable metrics that may be
141 causally related to the occurrence of incidents or injuries (Lingard et al. 2017; Versteeg et al. 2019).
142 New technologies, like wearable devices, automated data collection, bar codes, etc., enable
143 companies to capture more data related to safety issues (Ahn et al. 2019). Further research in the

144 area is still needed to access quality data. Nevertheless, there are still practical challenges on how
145 to use the data to improve safety management in practice.

146 To analyze the collected data stored in different data sources, statistical and ML models have been
147 applied widely to predict safety outcomes in construction projects at the early design stages. As
148 stated above, one practical challenge is the lack of consistency regarding the safety outcomes and
149 predictors measured across projects and organizations. This lack of consistency may result in the
150 need for different prediction approaches on a case-by-case basis. For example, Esmaeli et al.
151 (2015) tested the validity of generalized linear models to predict safety outcomes based on a large
152 volume of data on attributes that cause "struck-by" accidents. Poh et al. (2018) used an ML
153 approach to develop a model that forecasts accident occurrence and severity of construction
154 worksites based on project-related and safety-related input features. Kang and Ryu (2019)
155 proposed a Random Forest model that can predict occupational accidents based on accident and
156 weather data and suggest which management features significantly contribute to the forecast.
157 Sarkar et al. (2019) developed optimized ML-based models to predict incident outcomes at the
158 workplace using Support Vector Machine (SVM) and Artificial Neural Network (ANN)
159 algorithms on incident reports data. Baker et al. (2020) used a large dataset of over 90,000 incident
160 reports and various ML algorithms to develop a model that predicts injuries and their severity and
161 shows which attributes have high predictive power when the safety outcomes are external and
162 independent.

163 Although several ML algorithms have been proposed, their use in proactively assessing safety
164 performance in companies is not common. Several challenges exist since it is expected that a single
165 set of metrics may not be suitable for all construction industry sectors (Nasir et al. 2012). These
166 issues present a challenge regarding the selection of data that should be collected to predict safety

167 performance. Therefore, we propose a data-driven method that leverages HR data, commonly
168 recorded for reasons beyond performance metrics, and allows for an adequate model selection
169 based on the available data. This approach should be flexible enough to incorporate new variables
170 as their impact on safety is better understood.

171 **METHODOLOGY**

172 This paper reports on the experience obtained in a case study that implemented a data-driven
173 approach to predict the probability of a safety incident occurring within a month based on
174 preliminary schedule planning information. The predictive model used in this approach was
175 developed based on daily payroll information from previous projects and dated safety incident
176 reports. The model was implemented in a hypothetical case based on an industrial project to
177 demonstrate the framework's functionality during the project planning phase. Figure 1 depicts an
178 overview of the methodology followed in this study.

179 **Project Background**

180 The functionality of the proposed framework is illustrated through a hypothetical case study
181 inspired by a previous study (Taghaddos et al. 2021). This case study demonstrates that an
182 industrial construction enterprise can use the data available from various departments within the
183 enterprise to devise methods for predicting safety incidents based on discipline-specific
184 information. The case study also exemplifies the framework's functionality during the project
185 planning phase. The case study focuses on the planning phase of a small (6-month) in-situ drainage
186 oil sands project in Alberta, Canada, and mainly involves five disciplines: civil, operators,
187 pipefitters, electrical workers, and ironworkers.

188 The case study analyzes the risk of accident occurrence over the 6 months of the project at the
189 planning phase based on qualitative data. Data categorization was implemented based on 28
190 million working hours of historical data. Utilizing the discipline-based incident predictions
191 generated by the proposed model, the project plan was adjusted accordingly to mitigate safety
192 risks. This case study provides a concrete example of how a construction enterprise can put to use
193 existing data within the proposed framework and apply it to a real-world scenario. Moreover, it
194 demonstrates the manner in which the framework can be tailored to meet the unique needs of
195 various disciplines within the construction industry.

196 **Data Collection**

197 A large dataset of approximately 28 million working hours was collected across eight industrial
198 construction projects in Canada. The dataset reports different features obtained from daily payroll
199 information: working hours, age, work experience, time working on the project of workers and
200 foremen, crew sizes, number of operators, changes in the number of workers in the project, and
201 project progress. The data were categorized into five trade disciplines (electrical, ironworkers,
202 pipefitters, civil, and operators).

203 Safety reports were collected to determine the monthly occurrences of incidents related to each
204 discipline in each project throughout the data collection period. Incident reports, lost-time injuries,
205 medical aid injuries, and modified work injuries were considered safety incidents. The payroll
206 information dataset and the monthly safety incident occurrence information were collated and
207 integrated into a single dataset. Table 1 shows the features contained in the final dataset. A monthly
208 data record was collected for each discipline in each project.

209 The data is split into four qualitative categories: very high, high, moderate, and low. The value
210 categories range from 4 (very high) to 1(low). This qualitative categorization can ease the use of
211 the model in the planning phase. For example, the practitioners have some ideas about time-
212 dependent crew size during the project (i.e., how it varies during construction) but do not yet have
213 the exact numbers. For example, let's assume July is the peak month for the structural-steel
214 discipline's work in a particular project; hence, one can predict a large crew size will be required
215 in this case. Such high-level prediction can easily be combined with planning to improve safety
216 measurements.

217 Table 2 shows the number of data points collected for each discipline in this study. Each data point
218 represents information (reported monthly) on the features listed in Table 1, including the
219 occurrence of each reported safety incident (i.e., the target variable). The incident rates per
220 discipline were found to be as follows: ironworkers (34.93%), pipefitters (52.72%), civil (47.45%),
221 operators (3.78%), and electrical workers (20.74%).

222 **Model Development**

223 The final dataset was retrieved from eight industrial construction projects. Ten input features were
224 selected by applying the Boruta feature selection algorithm, which iteratively removes features
225 that prove to be less relevant than random probes (Kursa and Rudnicki 2010). These ten features
226 are categorized into five different trade disciplines. Five different ML models (Taghaddos and
227 Mohamed 2021): SVM, Decision Tree, Naïve Bayes, Naïve Bayes (Kernel), and Fast Large-
228 Margin, were developed from the resulting dataset containing the ten selected predictor features
229 and the target variable (whether or not a safety incident occurred). A cross-validation method was
230 implemented to select the model that worked best for the collected data by comparing the
231 prediction performance of the five developed models (Zhang and Yang 2015) using two measures:

232 accuracy and incident recall. Accuracy measures the percentage of correctly predicted records,
233 whereas incident recall measures the true positives recognition rate (Al-Turaiki et al. 2016). In the
234 context of construction safety, a false negative (i.e., predicting an incident as a "No incident") is
235 more expensive than a false positive (i.e., predicting a safe situation as an "incident") because false
236 positives only impose precautions to the system. For this study, an incident recall with a 75%
237 threshold was considered an important criterion for model selection. The framework was
238 developed using the educational version of the well-known data-mining tool, RapidMiner Studio
239 (Mierswa and Klinkenberg 2018). A screenshot of the developed model, as it appears in the data-
240 mining software, is provided in Figure 2.

241 **Model Implementation**

242 The proposed model was implemented in the hypothetical case study described above to
243 demonstrate its successful utilization in adjusting the project plan based on discipline-based
244 incident predictions. The model was applied to predict monthly occurrences of incidents on the
245 original program with planned information. A sensitivity analysis was performed on different
246 planning strategies considering modifications to the predictor features to reduce the number of
247 predicted incidents.

248 **RESULTS**

249 Table 3 shows the ten selected predictor features. These features are selected by applying the
250 Boruta feature selection algorithm to the dataset and using experts' opinions.

251 Five ML models were then developed and applied to the dataset to predict the target variable (i.e.,
252 whether or not an incident occurred) based on the selected predictor features. This process is
253 illustrated in Figure 2. In this figure, "CV" represents the cross-validation techniques in each

254 model. The accuracy and incident recall of the prediction performance of the models were
255 measured and compared after applying a cross-validation method. Table 4 shows the prediction
256 performance measures of the five models. Based on the results in Table 4, the Naïve Bayes
257 (Kernel) model was selected to demonstrate the framework's functionality in a hypothetical case
258 study.

259 Similar to the eight-project data stored in a database of industrial projects, the case study project
260 involves discipline-specific information. As such, the safety prediction incident results are also
261 discipline based. Moreover, the predictions are time-dependent since the prediction depends on
262 the features for the selected time period. Table 5 presents the original project schedule (Scenario
263 1). As demonstrated in the framework description section, each discipline is considered an
264 individual data point. The framework analyzed 28 entries for 6 months (civil, pipefitters, and
265 ironworkers) or 4 months (electrical/operators).

266 In Scenario 1 (Table 5), the framework assessed the entries based on the trained model and with
267 the original planning strategy, it predicted ten incidents (which is equivalent to 35% of the entries)
268 within the 6 months — Pipefitters (4), Civil (2), Ironworkers (3), and Electrical (1). In Scenario 2
269 (Table 6), a tentative planning strategy is used to reduce the number of accidents. Due to the
270 framework's ability to consider a variety of strategies capable of reducing the likelihood of an
271 accident, the following discipline-based mitigation strategies were adopted (individually or
272 grouped): decreasing the percentage of foremen older than 50 years, decreasing the rate at which
273 new workers are deployed to work on the project, reducing the crew size, increasing workers with
274 more than 3 years of experience, and reducing discipline working hours. The output results shown
275 are accident-free, reinforcing the framework's ability to test different planning strategies to
276 improve the project safety performance.

277 It is important to mention that just one discipline (Pipefitter – month 3) is required to reduce the
278 working hours. This strategy should be considered carefully since reducing working hours may
279 lead to further project completion delays and/or increase production pressure on workers in the
280 months ahead. The case study demonstrated that discipline-level data can produce better-tailored
281 strategies than those developed for the project level. As Pereira et al. (2020) suggested, crew size
282 is a factor that may impact project safety performance, and increasing the number of foremen may
283 reduce the likelihood of accidents. The framework presented in this research demonstrated that
284 reducing the crew size, not considering the discipline, may lead to increased project cost and would
285 not improve safety performance. Moreover, the framework also demonstrated that schedule
286 planning strategies such as new workers rate (i.e., the rate at which new workers are deployed to
287 work on the project), discipline working hours, discipline-based progress, and accumulative
288 working hours should be considered concomitantly with their impact on the project safety
289 performance.

290 While some features can be considered in the planning phase, others may assist the organization
291 in identifying flaws in the Safety Management System. Workers' and foremen's age features show
292 that a reinforcement of the organization's policies should account for this specific group to ensure
293 they are better prepared to identify hazards or follow safety procedures. The impact of workers'
294 experiences also demonstrated that safety induction and safety training should assess the workers'
295 abilities to retain the course knowledge and apply it in practice.

296 **DISCUSSION**

297 This case study validates the premise that the proposed predictor features in Table 3 can predict
298 the occurrence of safety incidents, enabling an informed decision-making process regarding
299 discipline-level schedule planning. The approach can simulate risks for any context, including new

300 work if the fundamental attributes remain stable (Hallowell et al. 2020). The proposed framework
301 predicts the likelihood of a safety incident occurring, overcoming a common limitation of
302 traditional attribute-based safety risk assessment, which predicts the outcome of an incident should
303 one occur (Choi et al. 2020; Esmaeili et al. 2015; Hallowell et al. 2020; Koc et al. 2021; Tixier et
304 al. 2016a; b).

305 The proposed framework uses historical data, both at the business and project management levels,
306 to discover data-driven knowledge and use it to support project management decisions, as
307 advocated by You and Wu (2019). The study emphasizes the importance of effective data
308 management in the context of construction safety. The method incorporates metrics extracted from
309 daily payroll data, which are typically not used in traditional safety planning. By leveraging
310 existing data, the method streamlines the development of predictive models for safety incidents,
311 equipping decision-makers with a useful tool for the proactive assessment and mitigation of risk.
312 Another novel feature of the developed framework is its ability to adjust the construction plan
313 based on safety incident predictions in order to mitigate risk. These adjustments could include
314 changing the crew composition, crew size, or time-dependent discipline working hours, all of
315 which can affect the project duration.

316 The case study presented in this paper demonstrates how a data-driven model can be incorporated
317 into the scheduling process contributing to safety planning. This finding agrees with Yi and
318 Langford (2006), who advocate for scheduling construction to reduce accident risks. The
319 suggested discipline-based scheduling is aligned with Hallowell and Gambatese (2009), who
320 recommend considering risks based on activities to target high-risk activities in safety programs.

321 The feature selection process results align with other studies predicting safety outcomes. For
322 example, Rivas et al. (2011) identified "task duration in hours", "length of time doing the job",

323 "job type" and "worker age" within the five most relevant features predicting accidents. Poh et al.
324 (2018) used the Boruta algorithm, and the features "percentage of project completion" and
325 "average monthly project manpower" were within the selected features to predict accident
326 occurrence and severity. Choi et al. (2020) also found that "age" and "service length" were among
327 the most important factors in predicting the likelihood of fatal accidents. These factors are
328 commonly associated with accident precursors (Pereira et al. 2018b). While the proposed model
329 does not attempt to uncover underlying causality, these alignments with knowledge of the
330 construction safety domain were essential during the development of the model and the subsequent
331 interpretation of its predictions (Mannering et al. 2020).

332 **CONCLUSION**

333 This manuscript proposes a novel data-driven approach for deploying safety policies for discipline-
334 level schedule planning. This novel approach enables practitioners to account for safety
335 considerations in the planning phase and proactively make appropriate decisions without needing
336 detailed quantitative information. Five ML models were developed from payroll data collected in
337 eight large industrial construction projects. The accuracy and incident recall of the prediction
338 performance of the models were measured and compared to select the model that worked best on
339 the collected data. Subsequently, the practical utility of the model was demonstrated through a
340 case study.

341 The findings reveal that the predicted occurrence of safety incidents can be reduced by modifying
342 the predictor features during the project's planning phase, achieving a safer planning strategy. In
343 the case study project, the original plan and schedule were revised based on discipline-specific
344 tentative planning strategies—e.g., decreasing the rate of new workers and crew sizes (to varying
345 degrees depending on the discipline). Accordingly, the incident rate was reduced from 35% to 0%,

346 resulting in an incident-free plan. The discipline-based safety plan for the early planning stage as
347 proposed herein is beneficial to practitioners in that it provides the basis for expanding the plan to
348 the work-package level later in the project.

349 This study makes three important contributions to knowledge and practice in this domain. First, it
350 provides a framework for proactive safety improvement in the planning phase and for deploying
351 discipline-level safety policies by identifying suitable safety-related metrics that construction
352 enterprises typically record for other purposes. In this manner, it helps with tackling the large
353 volumes of project-level data to identify, capture, and analyze the features that affect safety
354 performance by leveraging existing data for model development. Moreover, the discipline-based
355 approach used in the case study demonstrates the adaptability of the proposed framework to meet
356 the discipline-specific needs, and align with the unique features, of the construction industry.

357 Second, it provides a novel data acquisition method. Specifically, this study demonstrates that
358 payroll data and incident reports, which tend to be more reliable and unbiased than other data
359 sources due to their monetary/regulatory nature, can be used to develop a model for safety-related
360 decision support.

361 Third, this study integrates discipline-level scheduling with safety prediction. A key consideration
362 in this regard is that the discipline-specific level of scheduling is not so high-level as to miss the
363 vital discipline-specific features that are important in decision-making (such as in the case of
364 project-level scheduling), and not so detailed as to make planning and decision-making
365 cumbersome (such as in the case of crew-level scheduling).

366 The authors believe that this research can help project practitioners identify which data should be
367 collected in their projects and define strategies to improve their construction plans in terms of

368 safety based on insights emerging from the data. By leveraging the data-driven discipline-level
369 safety prediction model, project teams can make informed decisions and implement proactive
370 measures to enhance safety performance. Furthermore, the model's flexibility allows for the
371 inclusion of additional factors specific to each project, ensuring a comprehensive and tailored
372 approach to safety management. For example, in a heavy construction project involving excavation
373 work, factors such as the type of equipment used (e.g., excavators, bulldozers, etc.), as well as the
374 competency of equipment operators, can significantly influence safety outcomes. Ultimately, this
375 research aims to contribute to the advancement of construction safety practices by promoting
376 evidence-based decision-making and proactive risk mitigation strategies.

377 **LIMITATIONS and FUTURE WORK**

378 While the developed framework has been found to be capable of predicting the probability of a
379 safety incident occurring in a construction project and within a novel discipline-level schedule,
380 this study is subject to certain limitations.

381 One notable limitation is the lack of worker-level data that could lead to more accurate predictor
382 features. For example, research has shown that the psychological status of workers on the site
383 directly influences unsafe behaviors (Guo et al. 2017); however, this feature is difficult to measure
384 or predict. Particularly as some types of worker-level data are typically subject to data protection
385 laws. In future work, site-related data at the project level should be collected (in the form of
386 incident days and incident-free days) in order to develop a more comprehensive model (Choi et al.
387 2020).

388 The example presented in this paper can be developed further by enhancing the dataset. The dataset
389 could be expanded to include near-miss incidents, which would add an additional metric on safety
390 performance (Shen and Marks 2015), provided that the incident reports specify the trade(s)

391 involved in the incident and therefore are aligned with the data structure followed in the present
392 study. In this manner, the prediction of near-miss incidents could be incorporated to build upon
393 the present work.

394 **DATA AVAILABILITY STATEMENT**

395 Data used in this study were provided by a third party. Direct requests for these materials may be
396 made to the provider indicated in the Acknowledgments.

397 **ACKNOWLEDGMENTS**

398 This project was supported by a Collaborative Research and Development Grant (CRDPJ 492657)
399 from the Natural Sciences and Engineering Council of Canada. The authors would also like to
400 thank PCL Industrial Management Inc. for their continued support, collaboration, in-depth
401 knowledge, and provision of historical project data.

402 **REFERENCES**

- 403 American Fuel and Petrochemical Manufacturers – AFPM (2017). 2017 AFPM Occupational
404 Injury and Illness Report.
405 <https://www.afpm.org/sites/default/files/issue_resources/AFPM_Charts.pdf> (accessed 5
406 December 2022). Ahn, C. R., S. Lee, C. Sun, H. Jebelli, K. Yang, and B. Choi. 2019. “Wearable
407 Sensing Technology Applications in Construction Safety and Health.” *J. Constr. Eng. Manag.*, 145
408 (11): 03119007. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001708)
409 [7862.0001708](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001708).
- 410 Ahn, S., L. Crouch, T. W. Kim, and R. Rameezdeen. 2020. “Comparison of Worker Safety Risks between
411 Onsite and Offsite Construction Methods: A Site Management Perspective.” *J. Comput. Civ. Eng.*,
412 146 (9): 05020010. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001890](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001890).

413 Al-Turaiki, I., M. Alshahrani, and T. Almutairi. 2016. "Building predictive models for MERS-CoV infections
414 using data mining techniques." *J. Infect. Public Health*, 9 (6): 744–748.
415 <https://doi.org/10.1016/j.jiph.2016.09.007>.

416 Anvari, B., P. Angeloudis, and W. Y. Ochieng. 2016. "A multi-objective GA-based optimisation for holistic
417 Manufacturing, transportation and Assembly of precast construction." *Autom. Constr.*, 71 (Part 2):
418 226–241. Elsevier. <https://doi.org/10.1016/J.AUTCON.2016.08.007>.

419 AWCBC. 2019. *Overview The Association of Workers' Compensation Boards of Canada produces this*
420 *annual report, which provides national statistics on work-related Fatalities and Lost Time Claim*
421 *compensation for injuries and diseases*.

422 Baker, H., M. R. Hallowell, and A. J. P. Tixier. 2020. "AI-based prediction of independent construction
423 safety outcomes from universal attributes." *Autom. Constr.*, 118: 103146. Elsevier B.V.
424 <https://doi.org/10.1016/j.autcon.2020.103146>.

425 Bavafa, A., A. Mahdiyar, and A. K. Marsono. 2018. "Identifying and assessing the critical factors for
426 effective implementation of safety programs in construction projects." *Saf. Sci.*, 106 (March): 47–
427 56. Elsevier. <https://doi.org/10.1016/j.ssci.2018.02.025>.

428 Boon, J., H. Yap, P. Eng, I. N. Chow, and K. Shavarebi. 2019. "Criticality of Construction Industry Problems
429 in Developing Countries: Analyzing Malaysian Projects." *J. Manag. Eng.*, 35 (5): 04019020–12.
430 [https://doi.org/10.1061/\(ASCE\)ME.1943-5479.0000709](https://doi.org/10.1061/(ASCE)ME.1943-5479.0000709).

431 Choe, S., and F. Leite. 2017. "Construction safety planning: Site-specific temporal and spatial information
432 integration." *Autom. Constr.*, 84: 335–344. Elsevier B.V.
433 <https://doi.org/10.1016/j.autcon.2017.09.007>.

434 Choi, J., B. Gu, S. Chin, and J. S. Lee. 2020. "Machine learning predictive model based on national data
435 for fatal accidents of construction workers." *Autom. Constr.*, 110: 102974. Elsevier.
436 <https://doi.org/10.1016/J.AUTCON.2019.102974>.

437 Esmaeili, B., M. R. Hallowell, and B. Rajagopalan. 2015. "Attribute-Based Safety Risk Assessment. II:
438 Predicting Safety Outcomes Using Generalized Linear Models." *J. Constr. Eng. Manag.*, 141 (8):
439 04015022. American Society of Civil Engineers (ASCE). [https://doi.org/10.1061/\(ASCE\)CO.1943-
440 7862.0000981](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000981).

441 Goh, Y. M., and D. Chua. 2013. "Neural network analysis of construction safety management systems: a
442 case study in Singapore." *Constr. Manag. Econ.*, 31 (5): 460–470.
443 <https://doi.org/10.1080/01446193.2013.797095>.

444 Guo, H., Y. Yu, T. Xiang, H. Li, and D. Zhang. 2017. "The availability of wearable-device-based physical
445 data for the measurement of construction workers' psychological status on site: From the
446 perspective of safety management." *Autom. Constr.*, 82: 207–217. Elsevier.
447 <https://doi.org/10.1016/J.AUTCON.2017.06.001>.

448 Hallowell, M. R., S. Bhandari, and W. Alruqi. 2020. "Methods of safety prediction: analysis and
449 integration of risk assessment, leading indicators, precursor analysis, and safety climate." *Constr.*
450 *Manag. Econ.*, 38 (4): 308–321. Routledge. <https://doi.org/10.1080/01446193.2019.1598566>.

451 Hallowell, M. R., and J. A. Gambatese. 2009. "Activity-Based Safety Risk Quantification for Concrete
452 Formwork Construction." *J. Constr. Eng. Manag.*, 135 (10): 990–998. American Society of Civil
453 Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000071](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000071).

454 Han, S., F. Saba, S. Lee, Y. Mohamed, and F. Pena-Mora. 2014. "Toward an understanding of the impact

455 of production pressure on safety performance in construction operations." *Accid. Anal. Prev.*, 68
456 (7): 106–116. <https://doi.org/10.1016/j.aap.2013.10.007>.

457 Haslam, R. A., S. A. Hide, A. G. F. Gibb, D. E. Gyi, T. Pavitt, S. Atkinson, and A. R. Duff. 2005. "Contributing
458 factors in construction accidents." *Appl. Ergon.*, 401–415. Elsevier Ltd.

459 Jiang, Z., D. Fang, and M. Zhang. 2015. "Understanding the causation of construction workers' unsafe
460 behaviors based on system dynamics modeling." *J. Manag. Eng.*, 31 (6): 10.1061/(ASCE)ME.1943-
461 5479.0000350, 04014099. [https://doi.org/10.1061/\(ASCE\)ME.1943-5479.0000350](https://doi.org/10.1061/(ASCE)ME.1943-5479.0000350).

462 Kang, K., and H. Ryu. 2019. "Predicting types of occupational accidents at construction sites in Korea
463 using random forest model." *Saf. Sci.*, 120: 226–236. Elsevier B.V.
464 <https://doi.org/10.1016/j.ssci.2019.06.034>.

465 Karakhan, A. A., S. Rajendran, J. Gambatese, and C. Nnaji. 2018. "Measuring and Evaluating Safety
466 Maturity of Construction Contractors: Multicriteria Decision-Making Approach." *J. Constr. Eng.
467 Manag.*, 144 (7): 1–13. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001503](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001503).

468 Koc, K., Ö. Ekmekcioğlu, and A. P. Gurgun. 2021. "Integrating feature engineering, genetic algorithm and
469 tree-based machine learning methods to predict the post-accident disability status of construction
470 workers." *Autom. Constr.*, 131: 103896. Elsevier. <https://doi.org/10.1016/J.AUTCON.2021.103896>.

471 Kursu, M. B., and W. R. Rudnicki. 2010. "Feature selection with the boruta package." *J. Stat. Softw.*, 36
472 (11): 1–13. <https://doi.org/10.18637/jss.v036.i11>.

473 Lingard, H., M. Hallowell, R. Salas, and P. Pirzadeh. 2017. "Leading or lagging? Temporal analysis of
474 safety indicators on a large infrastructure construction project." *Saf. Sci.*, 91: 206–220. Elsevier.
475 <https://doi.org/10.1016/j.ssci.2016.08.020>.

476 Mannering, F., C. R. Bhat, V. Shankar, and M. Abdel-Aty. 2020. "Big data, traditional data and the
477 tradeoffs between prediction and causality in highway-safety analysis." *Anal. Methods Accid. Res.*,
478 25: 100113. Elsevier. <https://doi.org/10.1016/J.AMAR.2020.100113>.

479 Mierswa, I., and R. Klinkenberg. 2018. "RapidMiner Studio (9.1)."

480 Mitropoulos, P., T. S. Abdelhamid, and G. A. Howell. 2005. "Systems model of construction accident
481 causation." *J. Constr. Eng. Manag.*, 131 (7): 816–825.

482 Mohammadi, A., M. Tavakolan, and Y. Khosravi. 2018. "Factors influencing safety performance on
483 construction projects: A review." *Saf. Sci.*, 109 (2018): 382–397.
484 <https://doi.org/10.1016/j.ssci.2018.06.017>.

485 Nasir, H., C. T. Haas, J. H. Rankin, A. R. Fayek, D. Forgues, and J. Ruwanpura. 2012. "Development and
486 implementation of a benchmarking and metrics program for construction performance and
487 productivity improvement 1 This paper is one of a selection of papers in this Special Issue on
488 Construction Engineering and Management." *Can. J. Civ. Eng.*, 39 (9): 957–967.
489 <https://doi.org/10.1139/l2012-030>.

490 Neale, J., and A. Gurmu. 2021. "Production pressures in the building sector of the construction industry:
491 a systematic review of literature." *J. Eng. Des. Technol.* Emerald Publishing Limited.
492 <https://doi.org/10.1108/JEDT-12-2020-0529>.

493 Pereira, E., M. Ali, L. Wu, and S. Abourizk. 2020. "Distributed Simulation-Based Analytics Approach for
494 Enhancing Safety Management Systems in Industrial Construction." *J. Constr. Eng. Manag.*, 146 (1):
495 1–12. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001732](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001732).

496 Pereira, E., S. Han, and S. AbouRizk. 2018a. "Integrating Case-Based Reasoning and Simulation Modeling

497 for Testing Strategies to Control Safety Performance.” *J. Comput. Civ. Eng.*, 32 (6): 04018047.
498 [https://doi.org/10.1061/\(ASCE\)CP.1943-5487.0000792](https://doi.org/10.1061/(ASCE)CP.1943-5487.0000792).

499 Pereira, E., U. Hermann, S. Han, and S. AbouRizk. 2018b. “Case-Based Reasoning Approach for Assessing
500 Safety Performance Using Safety-Related Measures.” *J. Constr. Eng. Manag.*, 144 (9): 04018088.
501 American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001546](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001546).

502 Poh, C. Q. X., C. U. Ubeynarayana, and Y. M. Goh. 2018. “Safety leading indicators for construction sites:
503 A machine learning approach.” *Autom. Constr.*, 93: 375–386. Elsevier B.V.
504 <https://doi.org/10.1016/j.autcon.2018.03.022>.

505 Rivas, T., M. Paz, J. E. Martín, J. M. Matías, J. F. García, and J. Taboada. 2011. “Explaining and predicting
506 workplace accidents using data-mining techniques.” *Reliab. Eng. Syst. Saf.*, 96 (7): 739–747.
507 Elsevier. <https://doi.org/10.1016/J.RESS.2011.03.006>.

508 Salas, R., and M. Hallowell. 2016. “Predictive validity of safety leading indicators: empirical assessment
509 in the oil and gas sector.” *J. Constr. Eng. Manag.*, 142 (10): 10.1061/(ASCE)CO.1943-7862.0001167,
510 04016052. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001167](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001167).

511 Sarkar, S., S. Vinay, R. Raj, J. Maiti, and P. Mitra. 2019. “Application of optimized machine learning
512 techniques for prediction of occupational accidents.” *Comput. Oper. Res.*, 106: 210–224. Elsevier
513 Ltd. <https://doi.org/10.1016/j.cor.2018.02.021>.

514 Shen, X., and E. Marks. 2015. “Near-Miss Information Visualization Tool in BIM for Construction Safety.”
515 *J. Constr. Eng. Manag.*, 142 (4): 04015100. American Society of Civil Engineers.
516 [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001100](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001100).

517 Taghaddos, M., H. Taghaddos, U. Hermann, Y. Mohamed, and S. AbouRizk. 2021. “Hybrid multi-mode

518 simulation and optimization for subarea scheduling in heavy industrial construction.” *Autom.*
519 *Constr.*, 125 (2021): 1–18. <https://doi.org/10.1016/j.autcon.2021.103616>.

520 Tixier, A. J. P., M. R. Hallowell, B. Rajagopalan, and D. Bowman. 2016a. “Application of machine learning
521 to construction injury prediction.” *Autom. Constr.*, 69: 102–114. Elsevier B.V.
522 <https://doi.org/10.1016/j.autcon.2016.05.016>.

523 Tixier, A. J. P., M. R. Hallowell, B. Rajagopalan, and D. Bowman. 2016b. “Automated content analysis for
524 construction safety: A natural language processing system to extract precursors and outcomes
525 from unstructured injury reports.” *Autom. Constr.*, 62: 45–56. Elsevier.
526 <https://doi.org/10.1016/J.AUTCON.2015.11.001>.

527 Versteeg, K., P. Bigelow, A. M. Dale, and A. Chaurasia. 2019. “Utilizing construction safety leading and
528 lagging indicators to measure project safety performance: A case study.” *Saf. Sci.*, 120: 411–421.
529 Elsevier. <https://doi.org/10.1016/j.ssci.2019.06.035>.

530 W Guo, B. H., and T. Wing Yiu. 2016. “Developing Leading Indicators to Monitor the Safety Conditions of
531 Construction Projects.” *J. Manag. Eng.*, 32 (1): 04015016–14. [https://doi.org/10.1061/\(ASCE\)ME](https://doi.org/10.1061/(ASCE)ME).

532 Yi, K.-J., and D. Langford. 2006. “Scheduling-Based Risk Estimation and Safety Planning for Construction
533 Projects.” *J. Constr. Eng. Manag.*, 132 (6): 626–635. [https://doi.org/10.1061/\(ASCE\)0733-](https://doi.org/10.1061/(ASCE)0733-)
534 [9364\(2006\)132:6\(626\)](https://doi.org/10.1061/(ASCE)0733-9364(2006)132:6(626)).

535 You, Z., and C. Wu. 2019. “A framework for data-driven informatization of the construction company.”
536 *Adv. Eng. Informatics*, 39: 269–277. Elsevier. <https://doi.org/10.1016/J.AEI.2019.02.002>.

537 Zhang, P., N. Li, Z. Jiang, D. Fang, and C. J. Anumba. 2018. “An agent-based modeling approach for
538 understanding the effect of worker-management interactions on construction workers’ safety-

539 related behaviors.” <https://doi.org/10.1016/j.autcon.2018.10.015>.

540 Zhang, Y., and Y. Yang. 2015. “Cross-validation for selecting a model selection procedure.” *J. Econom.*,
541 187 (1): 95–112. <https://doi.org/10.1016/j.jeconom.2015.02.006>.

542 **List of Figure Captions**

543 Figure 1: Method Overview

544 Figure 2: Developed model using RapidMiner Studio

545

546

TABLE 1. Features collected from daily payroll information

feature name	Description
Proj_id*	Project identifier
WH_month	Monthly working hours
WH_cml	Monthly cumulative working hours
WH_cml-pct	Percentage of cumulative working hours
WH_diff	Increase/decrease in working hours compared to the previous month
Age_wrks-30-less	Percentage of workers aged 30 or less
Age_frmn-30-less	Percentage of foremen aged 30 or less
Age_wrks-50-more	Percentage of workers aged 50 or more
Age_frmn-50-more	Percentage of foremen aged 50 or more
WEx_wrks-new	Percentage of new workers compared to the previous month
WEx_wrks-3-less	Percentage of workers with up to 3 years of experience
DS_wrks	Workers' average number of days on the site
DS_frmn	Foremen's average number of days on the site
Crew-size	Crew size
Proj_s-curve-inc	Monthly S-Curve increase
Proj_pct-cplt	Percentage of project completion
Proj_ramp	Increase/decrease of workers on the project
S_incident	Monthly safety incidents occurring

*All features, except Proj_id, were collected by discipline

547

548

TABLE 2. Number of collected data points by discipline

Trade discipline	Total data points	Incident data points
Ironworkers	146	51
Pipefitters	165	87
Civil	177	84
Operators	132	5
Electrical	135	28

549

550

TABLE 3. Selected predictor features

feature name	Description
WH_month	Monthly working hours
WEx_wrks-3-less	Percentage of workers with up to 3 years of experience
Age_frmn-30-less	Percentage of foremen aged 30 or less
Age_wrks-30-less	Percentage of workers aged 30 or less
WEx_wrks-new	Percentage of new workers compared to the previous month
Proj_s-curve-inc	Monthly S-Curve increase

Crew-size	Crew size
DS_frmn	Foremen average number of days on the site
Age_frmn-50-more	Percentage of foremen aged 50 or more
WH_cml-pct	Percentage of cumulative working hours

551

552

TABLE 4. Prediction performance measures

Model	Accuracy	Standard deviation	Incident recall
SVM	72.2	3.83	61.39
Naïve Bayes (Kernel)	70.47	3.65	75.25
Decision Tree	67.97	3.31	57.43
Naïve Bayes	62.90	3.89	74.26
Fast Large-Margin	62.90	3.89	74.26

553

Table 5. Original project planning - Scenario 1

Month #	1	1	1	2	2	2	2	2	3	3	3	3	3	4	4	4	4	4	5	5	5	5	5	6	6	6	6	6
Discipline-Civil (Yes-1/No-0)	1	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0
Discipline-Ironworkers (Yes-1/No-0)	0	1	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0
Discipline-Pipefitters (Yes-1/No-0)	0	0	1	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0
Discipline-Electrical (Yes-1/No-0)	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1
Discipline-Operators (Yes-1/No-0)	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0
Age_frmn-30-less	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
DS_frmn	1	1	1	1	1	1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
Age_wrks-50-more	1	1	1	1	1	1	1	1	2	2	2	2	2	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
Age_wrks-30-less	4	4	4	4	4	4	4	4	4	3	3	3	3	3	3	2	2	2	2	2	2	2	2	2	2	2	2	2
NewWorkersRate	4	4	4	4	4	4	4	4	3	3	3	3	3	3	3	3	3	3	1	1	1	1	1	1	1	1	1	1
Crew_Size	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	3	3	3	3	3	3	3	3	3	3
WEx_wrks-3-less	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2	1	1	2	1	2	1	1	2	1	2	2	1
WH_month	2	1	1	2	3	1	1	1	3	3	2	3	3	3	3	3	4	3	3	3	3	4	3	3	3	3	4	3
Proj_s-curve_inc	1	1	1	1	2	1	1	1	2	2	1	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
WH_cml-pct	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2	1	2	1	2	2	2	2	2	2	2	2	2	2
Incident (Yes/No)	No	No	No	No	No	No	No	No	No	Yes	No	Yes	Yes	No	Yes	Yes	Yes	No	No	No	Yes	Yes	No	No	No	Yes	Yes	No



