

1 TORCHQL: A Programming Framework for Integrity 2 Constraints in Machine Learning

3 ANONYMOUS AUTHOR(S)

4 Finding errors in machine learning applications requires a thorough exploration of their behavior over data.
5 Existing approaches used by practitioners are often ad-hoc and lack the abstractions needed to scale this process.
6 We present TORCHQL, a programming framework to evaluate and improve the correctness of machine learning
7 applications. TORCHQL allows users to write queries to specify and check integrity constraints over machine
8 learning models and datasets. It seamlessly integrates relational algebra with functional programming to allow
9 for highly expressive queries using only eight intuitive operators. We evaluate TORCHQL on diverse use-cases
10 including finding critical temporal inconsistencies in objects detected across video frames in autonomous
11 driving, finding data imputation errors in time-series medical records, finding data labeling errors in real-
12 world images, and evaluating biases and constraining outputs of language models. Our experiments show that
13 TORCHQL enables up to 13x faster query executions than baselines like Pandas and MongoDB, and up to 40%
14 shorter queries than native Python. We also conduct a user study and find that TORCHQL is natural enough
15 for developers familiar with Python to specify complex integrity constraints.

18 1 INTRODUCTION

19 Machine learning models can fail in unexpected and harmful ways. Examples include fatalities
20 caused by self-driving vehicles [Wakabayashi 2018], vision models performing worse on people
21 with darker skin [Wilson et al. 2019], language models producing text containing offensive stereo-
22 types [Abid et al. 2021], and medical diagnosis models degrading in performance when used in
23 new hospitals [Zech et al. 2018]. Identifying and avoiding such behaviors is crucial to ensuring
24 performance, reliability, and trustworthiness of machine learning applications.

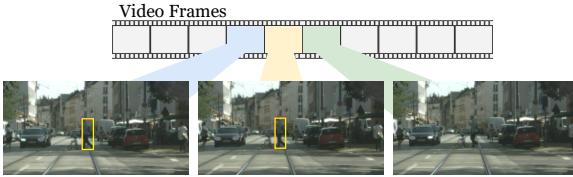
25 Finding and characterizing these behaviors is difficult. As a running example, consider Figure 1a,
26 which depicts three consecutive video frames from the CITYSCAPES self-driving dataset [Cordts et al.
27 2016]. A state-of-the-art vision model, OneFormer [Jain et al. 2022], fails to predict a pedestrian in
28 the third frame. This is a safety-critical issue since predicted objects should not suddenly disappear,
29 especially not pedestrians. This model failure, however, is difficult to find since the error occurs
30 sparsely within the 15K validation samples and it has minimal impact on most numeric evaluation
31 metrics, so even the highest-performing models exhibit such faults.

32 Similar to the above example, errors often exhibit sparsely in model predictions over vast datasets.
33 As such, in order to identify and characterize these errors, one must sift through large numbers of
34 correct predictions. Moreover, simply identifying the errors may not be enough; oftentimes we
35 must identify the underlying cause of these errors. In general, erroneous model predictions can be
36 cast as violations of certain properties. For instance, the error in Figure 1a violates the property of
37 object permanence across three adjacent frames in a video sequence.

38 We desire a general mechanism for specifying and checking such properties over data. *Integrity*
39 *Constraints*, originally studied in the context of databases [Godfrey et al. 1998], constitute such a
40 mechanism. As such, we can cast the problem of detecting and avoiding these errors as the problem
41 of specifying appropriate integrity constraints and evaluating them over model predictions. For
42 example, the error in Figure 1a can be cast as the violation of the constraint: *for every sequence of*
43 *three consecutive frames, a person detected in the first two frames should be detected in the third.*

44 Existing machine learning frameworks lack adequate abstractions to support integrity constraints.
45 Writing such constraints from scratch in Python typically also requires writing the infrastructure

46
47 2018. 2475-1421/2018/1-ART1 \$15.00
48 <https://doi.org/>



```
Query('temporal_consistency', base='preds')
    .join('preds', key=fid_plus1, fkey=fid)
    .join('preds', key=fid_plus2, fkey=fid)
    .project(match_three_bboxes)
    .project(zip_adj_preds).flatten()
    .filter(lambda sid, frames, lb, bx: (
        lb[0] == lb[1] and lb[2] == "No Match"
        and is_center(bx[0])))
```

(b) Integrity constraint as a query.

Fig. 1. (a) shows three consecutive frames in a video from the CITYSCAPES self-driving object-detection dataset. OneFormer, a state-of-the-art model, detects pedestrians in the center of the first two frames but not in the third. We formalize this fault as a violation of an integrity constraint. (b) shows a TORCHQL query to find all such violations over a dataset named ‘preds’ consisting of video frames along with the model’s predictions.

code needed to execute them. This places a burden on the user who may not know the constraint a priori and may want to test out multiple candidate constraints. Furthermore, once these constraints are specified, executing them over large scale data may require further optimizing the code, which is another burden on the user.

While frameworks such as Pandas and MongoDB can potentially abstract away some of the required infrastructure, they either have ad-hoc support for Python or come with complex interfaces and strict schema requirements. Furthermore, they may not naturally support the vast number of possible modalities that machine learning models operate over, from bounding boxes over video frames to unstructured text or audio data. This is also a challenge for systems such as LMQL [Beurer-Kellner et al. 2023] that seek to enforce integrity constraints over model outputs at runtime but are only applicable to the domain of language modeling.

We therefore observe that any effective framework for programming integrity constraints must address the following challenges:

- (1) *Scalability*: it must be able to check integrity constraints in large-scale machine learning settings,
- (2) *Interactivity*: it must support interactive testing and inspection of constraint violations, and
- (3) *Expressivity*: it must be able to support integrity constraints specified over a diverse landscape of models, datasets, and use-cases.

In this paper, we present TORCHQL, a framework that satisfies the above criteria. TORCHQL enables users to specify integrity constraints as *queries*. For this purpose, it introduces a query language that seamlessly integrates relational algebra with functional programming. Figure 1b shows an example query that specifies the constraint that detects the error in Figure 1a. TORCHQL executes the queries over a database representing datasets and model predictions to find violations of the corresponding integrity constraints. This allows developers to quickly and scalably test potential bugs on large datasets and models. Furthermore, the compositionality of TORCHQL queries enables the rapid refinement of previous iterations to interactively prototype new potential issues and reduce false alarms among discovered faults.

TORCHQL queries are easy to write yet highly expressive, capable of representing complex faults with only eight operators and simple user-defined functions. For example, the query in Figure 1b uses only four table operators—join, project, filter and flatten—along with three user-defined functions to search for violations of the object permanence constraint. This query finds only 627 frames containing such errors out of 15,000, allowing it to effectively detect errors in large-scale machine learning tasks, while systems like Python require considerable optimizations to do so without timing out. We discuss this example in more detail in Section 6.1.

We demonstrate how TORCHQL enables the use of integrity constraints for detecting and fixing correctness problems in a diverse set of machine learning tasks—object detection, data imputation,

image classification, text generation, and natural language reasoning—across domains involving self-driving videos, time-series medical records, and large language models (LLMs). We show the effectiveness of TORCHQL on these tasks and domains through five extensive case studies. Moreover, our evaluation shows that TORCHQL enables faster query executions (up to 13x best-case speed-ups) than baseline systems such as Pandas and MongoDB, and more concise queries (up to 40% shorter) than native Python. We also conduct a user study with 10 users to evaluate the usability of TORCHQL. We find that they were able to quickly learn TORCHQL to program integrity constraints and find model mispredictions.

We summarize the contributions of the paper:

- (1) We propose integrity constraints as a means for discovering and characterizing faults in general machine learning tasks.
- (2) We develop TORCHQL¹, a framework for programming and checking integrity constraints over models and datasets in a manner that is *scalable*, *interactive*, and *expressive*. At its core, TORCHQL provides abstractions to specify and evaluate integrity constraints as database queries.
- (3) We perform extensive experiments that demonstrate the efficiency and conciseness of TORCHQL for integrity constraint queries in a multitude of use-cases across a variety of domains.
- (4) We also complement our experiments with a user study and multiple case studies to validate the usability and expressiveness of TORCHQL.

2 ILLUSTRATIVE OVERVIEW

In this section, we illustrate how TORCHQL is used to iteratively discover and check integrity constraints to detect model faults, and contrast with programming the same constraints in native Python. This process is done offline, once the model is trained, but before it is deployed. In Section 6, we discuss other use-cases, such as using the constraints to detect model faults at runtime.

Consider the error in Figure 1 where a self-driving object detector fails to detect a pedestrian in the third frame of a video sequence. Since this is a critical bug, a machine learning practitioner would want to find other instances of this error, both within the training data as well as when the model is deployed. We therefore seek to *characterize* this error as an integrity constraint violation. Before we can accomplish this, however, we must first discover the integrity constraint itself.

We first initialize a TORCHQL database db as shown in Figure 2 and populate it with a table named ‘preds’ containing predictions of OneFormer over 15,000 frames. Any query q we write can then be evaluated over this table by invoking it over db.

Since we do not know the constraint beforehand, we start with a simple initial constraint: frames in which at least one pedestrian is predicted. This can be written as a simple TORCHQL query as shown in Figure 3a. Here, we specify the query name (people_count) and the table over which the query operates (preds). Each prediction consists of the bounding boxes and their corresponding labels for each frame. We supply a *user-defined function* to the filter operator, which executes it over each frame when the query is run, and keep the predictions that satisfy the condition.

We also show the same constraint implemented in native Python in Figure 3b. The lack of query abstractions in Python means that we need to write a for loop, an if condition, and manually add satisfying predictions to the data structure (people_count). As such, the lack of query abstractions often necessitates one to implement both *what* the constraint is, as well as *how* to evaluate it.

```
preds = ... # frames and predictions
db = Database()
db.register(preds, 'preds')
q = Query(...) # creating a query
result = q(db) # running the query
```

Fig. 2. Initializing a TORCHQL database and running queries over it.

¹The TORCHQL system is currently available in the supplementary material. We plan to release it publicly upon publication.

```
148 Query('people_count', base='preds')
149     .filter(lambda p:
150         count_people(p['labels']) != 0)
151
```

(a) TORCHQL

```
people_count = []
for p in preds:
    if count_people(p['labels']) != 0:
        people_count.append(p)
```

(b) Native Python

Fig. 3. Queries for finding individual frames in which at least one pedestrian is predicted.

```
155 Query('changing_people', base='preds')
156     .join('preds', key=fid_plus_1, fkey=fid)
157     .filter(lambda f1, f2:
158         count_people(f1['labels']) !=
159             count_people(f2['labels']))
```

(a) TORCHQL

```
changing_people = []
for p in preds:
    for q in preds:
        if fid_plus_1(p) == fid(q):
            if count_people(p['labels']) != count_people(q['labels']):
                changing_people.append([p, q])
```

(b) Native Python

Fig. 4. Queries for finding pairs of consecutive frames with different predicted pedestrian counts.

```
164 Query('temporal_consistency', base='preds')
165     .join('preds', key=fid_plus_1, fkey=fid)
166     .join('preds', key=fid_plus_2, fkey=fid)
167     .project(match_three_bbboxes)
168     .project(zip_adj_preds)
169     .flatten()
170     .filter(lambda sid, frames, lb, bx:
171         lb[0] == lb[1] and
172         lb[2] == "No Match" and
173             is_center(bx[0]))
```

(a) TORCHQL

```
temporal_consistency = []
for p in preds:
    for q in preds:
        for r in preds:
            if fid_plus_1(p) == fid(q) and
                fid_plus_2(p) == fid(r):
                bboxes = match_three_bbboxes(p, q, r)
                bboxes = zip_adj_preds(*bboxes)
                for box in bboxes:
                    if (sid, frames, lb, bx) in box:
                        if lb[0] == lb[1] and
                            lb[2] == "No Match" and
                            is_center(bx[0]):
                                temporal_consistency.append(box)
```

(b) Native Python

Fig. 5. Queries for finding sequences of three consecutive frames where a pedestrian is detected in the center of the first two frames but not in the third.

Evaluating this constraint gives us 4,591 frames out of the total 15,000 frames, many of them false positives. This means that the constraint is too broad and must be refined further. A subsequent hypothesis may look for two consecutive predictions which differ in the number of people detected. In such a case, it is possible that some people predicted in the first frame are missed in the second.

We can construct the query in Figure 4a to represent this hypothesis by chaining more table operators to the original query. First, we use a join operator to join the table of pedestrians with the same table shifted forward in time by one frame to get a new table containing pairs of consecutive predictions. Here too, we use user-defined functions to specify the key and foreign key of the join. This gives us a table of consecutive predictions. Observe that TORCHQL allows its queries to directly run on the objects being queried themselves. This allows us to seamlessly integrate relational algebra with lambda functions, as queries can contain and execute arbitrary Python functions (e.g., fid and fid_plus_1) over these objects without needing to define an explicit schema.

We can then apply another filter operator to select only pairs of predictions with a differing number of pedestrians. This now reduces the number of filtered frames to 3,638 from the 4,591 that were filtered out in the initial query. As was the case with the previous iteration, however, the Python implementation (Figure 4b) requires writing multiple looping and conditional statements to mimic the join and filter operators.

197 Our latest query produced a table of consecutive predictions with a differing number of people.
198 However, such predictions also include people that did not go missing, such as those exiting or
199 entering the second frame. We therefore need to identify the exact persons that have eluded the
200 object detector. This requires converting our latest table of consecutive predictions into a table of
201 consecutive persons. We also wish to do so over three frames instead of two to avoid even more
202 false positives, and restrict our constraint to predictions in the center of the frame to avoid issues
203 with people leaving or entering the frame.

204 The TORCHQL query identifying all such sequences of three consecutive frames is shown in
205 Figure 5a. Here, the sequences are such that the pedestrian is detected in the center of the first
206 frame, detected again in the second frame, but not in the third.

207 Our first step is to join an additional ‘preds’ table to get a set of three consecutive predictions
208 instead of two. We again use the `join` operator, similar to how we used it in the previous query.
209 Although we have a set of three consecutive predictions, each prediction contains an unordered
210 list of detected objects. As a result, it is not immediately obvious which bounding boxes refer to
211 the same person. We therefore use a well-developed tool from object detection for aligning the
212 bounding boxes between frames [Bolya et al. 2020]. We can apply this complex function, denoted as
213 `match_three_bboxes` to all of our prediction sequences using the `project` table operator, which
214 applies a function to every row in a table to create a new table.

215 We further use a combination of other user-defined functions, table operators, and a final filter to
216 extract consecutive objects where the object is detected as a person in the first two frames but not
217 in the third. This final constraint refines our search of critical errors from 3638 to 627 frames out
218 of 15,000. In contrast, the Python implementation shown in Figure 5b requires the additional join
219 to be explicitly set up as a nested loop, and more loops to mimic the `flatten` operator. Moreover,
220 the Python implementation now fails to scale to the 15,000 frames, timing out after 10 minutes,
221 while the TORCHQL query executes in around 43 seconds. While the Python implementation can be
222 optimized, doing so adds 24 more lines of code, imposing additional user burden.

223 We thus show that throughout this process, TORCHQL allowed us to quickly and easily execute
224 and iterate over prototype queries on not only the original dataset but also on intermediate tables
225 produced while trying to program the constraint to detect temporal inconsistencies. On the other
226 hand, iterating over constraints in Python is accompanied by setting up the infrastructure to scalably
227 execute the constraint. Our final query involves the composition of 4 unique table operators, each
228 of which is parameterized by a custom lambda function to handle various data types. This process
229 is representative of how the scalability, interactivity, and expressivity of TORCHQL allows us to
230 program integrity constraints to discover errors in machine learning applications.

231 3 THE TORCHQL FRAMEWORK

232 This section presents the TORCHQL framework. We first describe the underlying data model that
233 queries operate over and then present the syntax and semantics of queries.

234 3.1 Data Model

235 While integrity constraints can be programmed using existing systems like Pandas or native Python,
236 as we illustrated in Section 2, they are ill-equipped for our needs. If arbitrary Python objects are
237 needed for a debugging task, Python is preferable to Pandas, since the latter does not support
238 querying over arbitrary objects. Conversely, if the objects adhere to a rigid structure or can be easily
239 flattened into a table of primitive datatypes, it is more efficient to take advantage of Pandas querying
240 abstractions. Other candidates include NoSQL database systems like MongoDB and its in-memory
241 counterpart, ArangoDB. However, they are designed to store and query over key-value pairs and
242

246 247 248 249 250 251	Querying Systems	Data Representation	Flexible Schema	User-Defined Functions	Querying Abstractions
Pandas	numpy.array tables	X	✓	✓	
MongoDB [Banker et al. 2016]	key-value stores	✓	X	✓	
TORCHQL	object collections	✓	✓	✓	
Python, OMG [Kang et al. 2018]	arbitrary objects	✓	✓	X	

252
253 Table 1. Comparison of querying systems.

254 do not support user-defined functions (UDFs) that allow users to leverage external modules and
 255 express arbitrarily complex queries. We summarize these tradeoffs in Table 1.

256 In order to be usable, TORCHQL must support queries over arbitrary Python objects without
 257 requiring much data wrangling. TORCHQL therefore uses an object-relational representation to
 258 represent the data being queried. It inherits Python’s data model and considers an *object* (*o*) to
 259 be a fundamental abstraction of data. As such, each object may either be a Python primitive or
 260 instantiation of a Python class, a list or set of other objects, or a collection of key-value pairs.

261 These objects may be further organized into lists with other similar objects, referred to as
 262 *tables* (*T*). Objects in the same table are oftentimes related with each other, but need not conform
 263 to any schema, like frames from a video, or prompts for a large language model. Tables can be
 264 further assigned names and organized into a collection called a *database* (*D*). We show the semantic
 265 domains for objects, tables, and databases used by TORCHQL in Figure 6a.

266 An object-relational representation has several advantages over other common representations.
 267 First, while it requires objects being queried to be contained within a collection, it has no restrictions
 268 on the structure of the objects being queried. This allows for the flexibility to query over models and
 269 datasets without the associated overhead of converting the data into a compatible representation.

270 Moreover, storing these objects within collections allows for enough structure within each table
 271 to provide succinct yet powerful querying abstractions over them akin to their relational algebra
 272 counterparts. Operations such as joining tables or grouping rows turn into single-line specifications
 273 without the need for setting up low-level infrastructure that would otherwise be required.

274 An object-relational representation also allows for the support of executing UDFs over the objects
 275 within each table. TORCHQL supports this by allowing users to supplement each relational algebra
 276 abstraction with one or more UDFs if needed. This also allows the user to leverage external libraries,
 277 machine learning models, and other tools to write their queries within the TORCHQL framework.

279 3.2 The Query Language

280 We now describe TORCHQL’s query language. We first present its syntax, shown in Figure 6b,
 281 followed by the operational semantics, shown in Figure 7.

282 **3.2.1 Syntax.** A TORCHQL *program* comprises a sequence of statements each of which defines a
 283 named base table or a named query. Each query is a chain of table operators that define a sequence
 284 of transformations to apply to one or more tables in the database. TORCHQL provides eight unique
 285 table operators. These operators, with the exception of flatten and unique, are empowered by
 286 user-defined functions (UDFs). These UDFs enable powerful transformations by running arbitrary
 287 Python functions on each object of a given table, or in the case of reduce, over the entire table. We
 288 discuss these operators in more detail in Section 3.2.3, after describing the program semantics.

289 **3.2.2 Semantics.** A TORCHQL program is interpreted as a sequence of database transformations
 290 induced by its constituent statements. Statement *register*(*n*, *T*) adds base table *T* with the name
 291 *n* to the database whereas statement *n* ← *Q* runs query *Q* over tables in the database to result in a
 292 new table which is added to the database with the name *n*. In particular, *Q* is comprised of the name
 293

<pre> 295 (object) $o ::= p$ (name) n (operator) $a ::= \text{join}(n, f_1, f_2)$ 296 $[o_1, \dots, o_n]$ (function) f $\text{filter}(f) \mid \text{flatten}()$ 297 $\{o_1, \dots, o_n\}$ (query) $Q ::= n \mid Q.a$ $\text{project}(f)$ 298 $\{p_1 : o_1, \dots, p_n : o_n\}$ (statement) $S ::= \text{register}(n, T)$ $\text{order_by}(f)$ 299 (table) $T ::= [o_1, \dots, o_n]$ $n \leftarrow Q$ $\text{group_by}(f)$ (database) $D ::= [n_1 \rightarrow T_1, \dots, n_k \rightarrow T_k]$ (program) $P ::= \epsilon \mid S; P$ $\text{unique}() \mid \text{reduce}(f)$ </pre>	
(a) Semantic Domains.	(b) Abstract syntax.

Fig. 6. Core language of TORCHQL.

Program semantics

$$D[\cdot] : N \rightarrow T$$

$$\begin{array}{c}
\frac{}{D \vdash \epsilon \triangleright D} [\text{PROGRAM_E}] \quad \frac{D \vdash Q \triangleright T \quad D, T \vdash a \triangleright T'}{D \vdash Q.a \triangleright T'} [\text{QUERY_OP}] \quad \frac{}{D \vdash n \triangleright D[n]} [\text{QUERY_N}] \\
\\
\frac{D[n \mapsto T] \vdash P \triangleright D'}{D \vdash \text{register}(n, T); P \triangleright D'} [\text{PROGRAM_REG}] \quad \frac{D \vdash Q \triangleright T \quad D[n \mapsto T] \vdash P \triangleright D'}{D \vdash n \leftarrow Q; P \triangleright D'} [\text{PROGRAM_QUERY}]
\end{array}$$

Operator semantics

$$\begin{array}{c}
f : \mathbb{O} \rightarrow \mathbb{O}, \quad s : \mathbb{O} \rightarrow \text{Bool}, \quad g : T \rightarrow T, \\
\sigma_s : T \rightarrow T, \quad \cdot \bowtie_{f_K, f_{FK}} \cdot : T \rightarrow T, \quad \cdot \cdot \cdot \cdot \cdot : [\mathbb{O}] \rightarrow T \rightarrow T
\end{array}$$

$$\frac{}{D, T \vdash \text{join}(n, f_K, f_{FK}) \triangleright T \bowtie_{f_K, f_{FK}} D[n]} [\text{JOIN}]$$

$$\frac{D, T \vdash \text{project}(f) \triangleright T'}{D, t : T \vdash \text{project}(f) \triangleright f(t) : T'} [\text{PROJECT1}] \quad \frac{}{D, [] \vdash \text{project}(f) \triangleright []} [\text{PROJECT2}]$$

$$\frac{}{D, T \vdash \text{reduce}(g) \triangleright g(T)} [\text{REDUCE}] \quad \frac{}{D, T \vdash \text{filter}(s) \triangleright \sigma_s T} [\text{FILTER}]$$

$$\begin{array}{c}
D, T \vdash \text{filter}(\lambda x. f(x) \leq f(t)) \triangleright T_{\leq x} \\
D, T \vdash \text{filter}(\lambda x. f(x) > f(t)) \triangleright T_{>x} \\
D, T_{\leq x} \vdash \text{order_by}(f) \triangleright L \\
D, T_{>x} \vdash \text{order_by}(f) \triangleright R
\end{array}
\quad
\begin{array}{c}
D, T \vdash \text{filter}(\lambda x. f(x) = f(t)) \triangleright T_= \\
D, T \vdash \text{filter}(\lambda x. f(x) \neq f(t)) \triangleright T_\\neq \\
D, T_\\neq \vdash \text{group_by}(f) \triangleright T'
\end{array}$$

$$\frac{D, t : T \vdash \text{order_by}(f) \triangleright L++[t]++R}{D, t : T \vdash \text{order_by}(f) \triangleright L++[t]++R} [\text{ORDER1}] \quad \frac{D, t : T \vdash \text{group_by}(f) \triangleright (f(t), t : T_=) : T'}{D, t : T \vdash \text{group_by}(f) \triangleright (f(t), t : T_=) : T'} [\text{GROUP1}]$$

$$\frac{}{D, [], \vdash \text{order_by}(f) \triangleright []} [\text{ORDER2}] \quad \frac{}{D, [] \vdash \text{group_by}(f) \triangleright []} [\text{GROUP2}]$$

$$\frac{D, T \vdash \text{filter}(\lambda x. x \neq t) \triangleright T_\\neq \quad D \vdash \text{unique}() \triangleright T'}{D, t : T \vdash \text{unique}() \triangleright t : T'} [\text{UNIQUE1}] \quad \frac{}{D, [] \vdash \text{unique}() \triangleright []} [\text{UNIQUE2}]$$

$$\frac{D, T \vdash \text{flatten}() \triangleright T'}{D, t : T \vdash \text{flatten}() \triangleright t++T'} [\text{FLATTEN1}] \quad \frac{}{D, [] \vdash \text{flatten}() \triangleright []} [\text{FLATTEN2}]$$

Fig. 7. Operational semantics of TORCHQL.

n of an existing table and a chain of zero or more table operators $a_1.a_2 \dots a_k$. When executing Q over database D , table $T_0 = D[n]$ is first retrieved and then the table operations specified in Q sequentially transform the table. Starting from table T_0 , operation a_i transforms table T_{i-1} into T_i to eventually produce the output table T_k .

These semantics are vital for TORCHQL to be effective at interactively programming integrity constraints. Iterating over previous hypotheses requires the results of executed queries to be stored and usable across programming sessions. Moreover, since the TORCHQL language allows a single query to be arbitrarily complex, these semantics allow for simpler queries by decomposing complex operations into multiple queries. Queries can also act as preprocessors, allowing for features to be extracted from unstructured data before using those features to program integrity constraints. Sequentially executing and storing results of previous statements allows for these use cases.

344 3.2.3 *Operators*. We now describe the operator semantics of TORCHQL. The eight table operators
345 are largely drawn from relational algebra and natively support complex operations over tables. In
346 general, each operator takes as input a table, and produces a table as output. With the exception of
347 *flatten* and *unique*, users can supplement these operators with UDFs. However, the operators
348 differ in how these supplied UDFs are executed over the objects within each table.

349 (1) *Join*. This operator composes objects from two tables into a single table. The composition is
350 achieved through the supplied UDFs $f_K : \mathbb{O} \rightarrow \mathbb{O}$ and $f_{FK} : \mathbb{O} \rightarrow \mathbb{O}$, as shown in *Join* in Figure 7:

$$T_i \bowtie_{f_K, f_{FK}} T_j = [[o_i, o_j] | o_i \in T_i, o_j \in T_j, f_K(o_i) = f_{FK}(o_j)]$$

353 Here, \mathbb{O} denotes the domain of all objects. Since the key-foreign key pairs for the tables are defined
354 by the results of UDFs, one can join tables based on values that can be results of arbitrarily complex
355 operations not existing within the tables. An example of this is in the *join* of the query in Figure 1,
356 where the function `fid_plus_1` returns the frame ID of the next frame, while `fid` returns the frame
357 ID of the current frame. Using these functions allows one to join the table containing all the frames
358 with itself to produce a table containing pairs of consecutive frames.

359 This join is a form of equijoin, and can thus be implemented as a *hash join*, where we use a hash
360 table to join the objects from table T_i to those of T_j . Using a hash join algorithm allows us to
361 achieve this with a complexity of $O(m + n)$, where $m = |T_i|$ and $n = |T_j|$, as opposed to typical join
362 algorithms that perform with complexity $O(mn)$.

363 (2) *Project*. This operator transforms each object in a table according to the supplied UDF $f : \mathbb{O} \rightarrow \mathbb{O}$
364 as described in *PROJECT1* and *PROJECT2* in Figure 7. This allows us to perform powerful and
365 flexible transformations on the rows of the table, such as including image transformations using
366 image processing libraries, or analyzing the inferences made by external models as shown in the
367 query in Figure 10b where the *project* function is used to generate prompts for the LLMs.

368 (3) *Filter*. This operator uses the supplied UDF $s : \mathbb{O} \rightarrow \text{BOOL}$ to filter out rows from a table. This
369 is depicted in *FILTER* in Figure 7 using the selection operator σ from relational algebra.

371 (4) *Order By*. This operator reorders the objects of the table T according to the supplied UDF
372 $f : \mathbb{O} \rightarrow \mathbb{O}$ as shown in *ORDER1* and *ORDER2* in Figure 7. Again, the UDF allows the reordering of
373 objects without the need to explicitly populate the table with the values to order the objects by.

374 (5) *Group By*. This operator allows grouping objects of a table into subtables, where each subtable
375 contains the objects that produce the same value when passed to the UDF $f : \mathbb{O} \rightarrow \mathbb{O}$. This operator
376 is unique in the sense that its result contains nested tables that can be further manipulated using
377 table operators. This allows us to perform powerful group-by operations on the table, such as
378 grouping by the hour of the day or the sequence ID as shown in the query in Figure 14a.

379 (6) *Flatten*. This operator flattens each object $o_j \in T$ where o_j is a collection $[o_{j1}, o_{j2}, o_{j3}, \dots]$ to
380 produce table T' such that each element $o_{jm} \in T'$. If o_j is not a collection, then it is left unchanged.

381 (7) *Unique*. This operator returns a table T' where duplicate rows from T are removed.

383 (8) *Reduce*. This operator is different from the previously discussed operators in that it runs the
384 supplied UDF $g : \mathbb{T} \rightarrow \mathbb{T}$ over the *entire table* rather than over individual objects within that table.
385 Here, \mathbb{T} is the domain of tables. This not only allows for general aggregation functions like `count`,
386 `length`, `min`, `max`, and others, but also for recursive functions over tables to be run as a part of a
387 TORCHQL query without having to switch to native Python.

388 The use of UDFs within these operators enables a versatile querying system. First, since the
389 UDFs can be specified directly over objects, users do not need to wrangle the data to fit a schema
390 in order to efficiently query a model's predictions. Second, unlike in other querying systems,
391 extracting features necessary for building these queries now becomes integrated into the query

writing process, rather than being a separate component that practitioners independently refine. Finally, since UDFs can be arbitrary Python code, TORCHQL is Turing complete.

4 IMPLEMENTATION

We have implemented TORCHQL in Python. The underlying database system is designed to be in-memory. This design choice eliminates the overhead of storing and retrieving objects from disk while executing queries, thereby reducing the turnaround time for each query. As such, all tables being queried over are permanently loaded in the main memory for as long as the programming session lasts. While the data for the tables may be stored to disk and retrieved, they need to be loaded into memory before writing queries over them.

When writing queries, the tables that they access must exist in the database before they can be executed. Therefore, when a database is newly initialized, as is the case when starting a programming session for a new task, it must be populated with these tables using the `register` function. This is typically used for loading in the training, validation, or test datasets, after which model predictions can be obtained over them via TORCHQL queries.

In order to seamlessly integrate this process, as well as the interaction between TORCHQL queries and machine learning models, we design TORCHQL to support conventional machine learning frameworks, specifically PyTorch. To do so, TORCHQL tables inherit the base PyTorch Dataset class. In other words, TORCHQL tables are instantiations of PyTorch Datasets. This has a few advantages while working with PyTorch machine learning pipelines.

For one, users can directly register PyTorch datasets into TORCHQL databases without needing to cast them as TORCHQL tables explicitly, allowing them to directly query their data. For instance, we can download the training data for the MNIST dataset using PyTorch’s API and directly load it into a TORCHQL database as shown in Figure 8.

```
train_data = datasets.MNIST(
    root = 'data', train = True,
    transform = ToTensor(),
    download = True,
)
db = Database("mnist")
db.register(train_data, "train")
```

Fig. 8. Loading a PyTorch dataset.

This also allows TORCHQL’s database engine to leverage PyTorch’s in-built support for batch processing and vectorized tensor operations, speeding up the execution of queries whose UDFs work with batched tensors.

5 EVALUATION SETUP

We now discuss the setup for evaluating the effectiveness of TORCHQL as an integrity constraint programming framework. As discussed in Section 1, we require TORCHQL to be scalable, interactive, and expressive. We therefore evaluate it by answering the following research questions:

- RQ1. Expressivity:** Can TORCHQL be used to program integrity constraints in diverse settings?
- RQ2. Performance:** Are TORCHQL queries concise and efficient on large-scale data?
- RQ3. Usability:** Is TORCHQL intuitive and easy to use for users unfamiliar with the system?

For **RQ1**, we write queries for five machine learning tasks over domains such as self-driving videos, time-series healthcare data, trap-camera images, and natural language text. Out of 182 queries that were written, we present and analyze 11 queries in-depth across five case studies in Section 6. The tasks are summarized in Table 2 and the chosen queries are described in Table 3.

For **RQ2**, we implement seven of these queries in three baseline systems Python, Pandas, and ArangoDB to compare their conciseness and efficiency over their corresponding datasets and models. We present our findings in Section 7. For **RQ3**, we conduct a user study with 10 participants writing three queries of increasing complexity in TORCHQL. We present the results in Section 8.

Task	Dataset	Dataset Size	Model	Application domain
Object Detection	CITYSCAPES validation set [Cordts et al. 2016]	15K	OneFormer [Jain et al. 2022]	Self-Driving Videos
Data Imputation	Physionet-2012 Challenge [Silva et al. 2012]	4000	SAITS [Du et al. 2023]	Time-Series Healthcare Data
Image Classification	iWildCam training set [Beery et al. 2020]	121K	ResNet50 [He et al. 2016]	Trap-Camera Images
Text Generation	Alpaca [Taori et al. 2023]	52K	T5 [Raffel et al. 2020]	Natural Language
Natural Language Reasoning	GSM8K [Cobbe et al. 2021] Date Understanding [et al. 2023]	1319 369	Mistral 7B [Jiang et al. 2023]	Natural Language

Table 2. Summary of the datasets and models used in the experiments.

In the rest of this section, we describe each machine learning task from Table 3 including the chosen datasets and models, the correctness problems of interest, and the integrity constraints and queries to address them. We then proceed to the sections that answer each research question.

5.1 Object Detection

Overview and Setup. In this task, the goal is to predict bounding boxes and their labels for individual frames in self-driving videos. We consider the validation set of the CITYSCAPES dataset, which contains 15,000 video frames, and use the OneFormer [Jain et al. 2022] model for predicting bounding boxes and their labels for each video frame.

Correctness Problems. Object detection models, in general, are evaluated using metrics like mean average precision (mAP). This metric considers a prediction correct if it has an Intersection over Union (IoU) ratio of at least 0.5. It then aggregates the precision of all predictions over all the frames. However, in cases where critical errors are rare and sparsely distributed over the predictions, this may cause the accurate predictions to overshadow the severe errors.

Integrity Constraints. Given this issue with the mAP metric, we seek to discover critical model faults by programming integrity constraints like the one shown in Figure 1b. We wrote a total of 57 queries to investigate various potential integrity constraints and analyze two such constraints in Section 6.1. The first, denoted S-Q1 in Table 2, discovers all objects from all the 15K frames that occur in the center of one frame, somewhere in the next consecutive frame, and are not detected in the third consecutive frame. A similar query was also investigated by Kang et al. [2018]. The second, S-Q2, aims to find vehicles with extremely high speeds across three frames. These objects with outlier speeds are likely to indicate incorrect bounding box predictions since most objects (including moving cars) have limited movement between frames.

5.2 Time-Series Data Imputation

Overview and Setup. In this task, the goal is to predict missing values within time-series healthcare data. We consider the Physionet-2012 Challenge’s [Silva et al. 2012] medical time-series dataset. Each sample in this dataset is composed of multiple univariate time series, each of which includes records of one feature across multiple time stamps, as illustrated in Figure 9. The dataset contains 4000 time-series samples and each sample consists of 35 lab values (features), collected hourly, within a 48-hour time window. Overall, up to 80% of the values in the samples are missing. We use the state-of-the-art model SAITS [Du et al. 2023] to impute these missing values.

491	492	493	Task	Query	Query Description	Query Operators					
						Join	project	filter	group	order	flatten
494	495	496	Object Detection	S-Q1	Retrieve all sequences of three continuous frames from the Cityscapes dataset in which an object appearing in the first two frames, and the center of the first, is not detected in the third frame by the OneFormer model.	✓	✓	✓	✗	✗	✓
497	498	499		S-Q2	Retrieve all vehicles predicted in three consecutive frames of the Cityscapes dataset by the OneFormer model and compute their speed across the frames.	✓	✓	✓	✗	✗	✗
500	501	502	Data Imputation	T-Q1	Compute the difference between all non-missing pairs of temporally consecutive feature values and compute the 99th percentile difference by following [Naik et al. 2023], for each feature.	✓	✓	✗	✗	✗	✗
503	504	505		T-Q2	Compute the interquartile range of non-missing ground-truth values across time and use the range to detect outliers in the SAITS model’s imputations, for each feature.	✓	✓	✗	✗	✗	✓
506	507	508		T-Q3	Compute the interquartile range of all values (including imputed ones) across time and use the range to detect outliers in the SAITS model’s imputations, for each feature.	✗	✓	✗	✗	✗	✗
509	510	511	Image Classification	I-Q1	Retrieve trap camera video sequences containing frames with more than one unique ground-truth animal.	✗	✓	✓	✓	✓	✗
512	513	514		I-Q2*	Find empty trap camera frames where a ResNet model predicts an animal.	✗	✓	✓	✓	✗	✓
515	516	517	Text Generation	B-Q1	Retrieve adjectives used to describe the nouns ‘farmer’ and ‘engineer’ in the T5 Alpaca dataset.	✗	✓	✗	✓	✗	✗
518	519	520		B-Q2*	From a set of adjectives biased toward farmers and a set biased toward engineers, find the adjectives with which GPT-3.5 consistently chooses to describe farmers or engineers.	✓	✓	✗	✗	✗	✓
521	522	523	Natural Language Reasoning	L-Q1*	Use the Mistral-7B LLM to answer and provide chain-of-thought reasoning for the GSM8K arithmetic reasoning questions, guaranteeing that the answer is an integer.	✗	✓	✓	✓	✗	✓
524	525	526		L-Q2*	Use the Mistral-7B LLM to answer and provide chain-of-thought reasoning for the Date understanding dataset, guaranteeing that the answer is a valid date format.	✗	✓	✓	✓	✗	✓

Table 3. Queries written for evaluating TORCHQL. Queries with an asterisk are only used in case studies.

530 *Correctness Problems.* Missing data in healthcare settings is prevalent due to irregular patient visits
 531 or clinical errors [Lee et al. 2017]. Imputing this data is necessary to train models for downstream
 532 use-cases. As a result, it is important that the imputed data is as accurate as possible. As discussed in
 533 [You et al. 2018], medical data such as the sample in Figure 9 should respect domain knowledge, such
 534 as the Glasgow Coma Scale’s range of 1 to 15, as well as common sense, such as the assumption of
 535 smooth variations in univariate time series data (henceforth called the *smoothness assumption*). We
 536 therefore aim to find imputed values that do not respect these assumptions. For instance, Figure 9
 537 shows an imputed temperature value (99) at the 4th hour, which significantly exceeds neighboring
 538 non-missing temperature values and causes an erratic spike in the temperature values.
 539

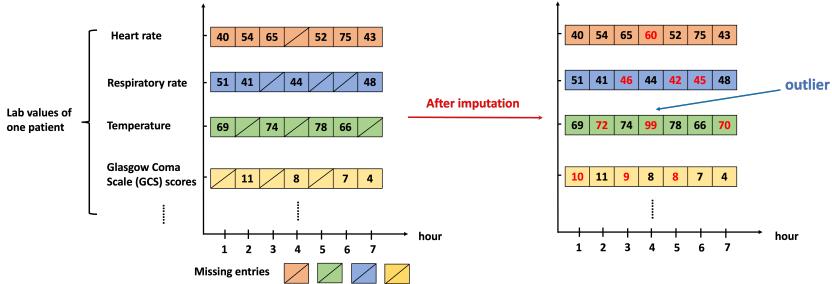


Fig. 9. A medical time series sample consisting of lab values of one patient over time. We show the original data with missing values on the left and the values imputed by SAITS on the right in red. The imputed value “99” for “Temperature” is an outlier since it deviates too far from other non-missing data for this feature.

Integrity Constraints. We wrote a total of around 30 queries to find different violations of this smoothness assumption. We investigate three of them T-Q1, T-Q2, and T-Q3 in Table 3, which describe two variations of the smoothness assumption.

The first such variant is that the difference between two continuous entries along the temporal dimension must be smooth and insignificant. One can use domain knowledge to determine a threshold to filter out those consecutive values with large differences, such as the outlier highlighted in Figure 9. In the absence of such domain knowledge, however, we formulate T-Q1 to collect all pairs of non-missing consecutive values from the entire dataset, and compute their differences. We then use the 99th percentile of all these gaps as the estimated threshold, in accordance with the method outlined in [Naik et al. 2023].

The other variant of the smoothness constraint that we specify captures the closeness between each imputed value and other entries within a time window. Specifically, we craft T-Q2 and T-Q3 to determine whether one imputed entry is an outlier or not with respect to other entries within a univariate time series. We discuss T-Q2 and T-Q3 in more detail in Section 6.2.

5.3 Image Classification

Overview and Setup. In this task, the goal is to classify the animal present in individual frames of videos captured by trap cameras. We evaluate the iWildCam dataset [Beery et al. 2020] and the predictions of a ResNet model [He et al. 2016] trained on the iWildCam training dataset.

Correctness Problems. The iWildCam dataset poses several challenges to machine learning models. We present two such problems that we investigate in more detail. First, models must have the ability to classify a frame as empty when no animal is present. As such, it is necessary for empty frames to be correctly labeled. However, there are several instances in iWildCam where empty frames are mislabeled as containing an animal. Second, the model tends to be inconsistent in its predictions. Over a sequence of frames, there is typically a single animal that moves around. However, the model sometimes predicts different animals in different frames within the same sequence.

Integrity Constraints. We wrote around 25 queries to capture various integrity constraints over the model predictions including whether nocturnal animals were predicted during the daytime, whether multiple animals were predicted across the same sequence of frames, and finding correlations between mispredicted classes and the ground truths. We analyze two of these queries.

The first, query I-Q1, aims to find sequences where the predictions by the model vary across frames, but the ground truth remains the same. The second, I-Q2, aims to find empty frames that are mislabeled to contain animals. Since the ground truth cannot be relied on, we compute the

589 pixel differences between pairs of consecutive frames to determine frames that have no movement
 590 and thus are likely empty. We investigate I-Q2 in more detail in Section 6.3.

591 5.4 Text Generation

593 *Overview and Setup.* In this task, the goal is to generate text (as a sequence of tokens) given a
 594 prompt. In particular, we consider the Alpaca instruction fine-tuning dataset [Taori et al. 2023]
 595 (Alpaca for short), which consists of 52K prompts and their responses from LLMs. In this experiment,
 596 we focus on the model responses from the GPT-3.5 (gpt-3.5-turbo) and T5² models.

597 *Correctness Problems.* Large language models (LLMs) are known to exhibit biases that can be
 598 potentially harmful [Havaldar et al. 2023; Liang et al. 2021; Zhao et al. 2017]. This can occur when
 599 the data they are trained over contain biases that the LLM may pick up on. However, typical metrics
 600 for language generation, such as the Bleu score, do not have the ability to check for them.
 601

602 *Integrity Constraints.* We investigate how integrity constraints can be programmed using TORCHQL
 603 to detect and study biases in LLMs. We wrote a total of around 50 queries to investigate various
 604 biases, and analyze two such queries in detail.

605 First, we write a query B-Q1 to discover the biases in the prompts of the Alpaca dataset by
 606 identifying the adjectives that are highly correlated with the occupations of farmer and engineer.
 607 This covers the class of adjective-profession biases [Kurita et al. 2019; Li et al. 2020; Nangia et al.
 608 2020; Smith et al. 2022], where certain adjectives are biased towards certain professions (e.g. brilliant
 609 scientists). Further details and examples of this bias are provided in Section 6.4. Second, we write a
 610 query B-Q2 to quantify bias in the model response. Here, we take adjectives associated with farmers
 611 and adjectives associated with engineers to construct prompts to determine how frequently an
 612 LLM uses the farmer adjective to describe farmers and the engineer adjective to describe engineers.
 613 More details are included in Section 6.4.

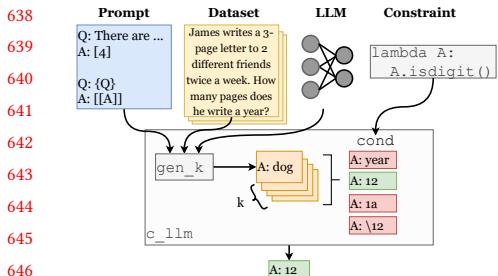
614 5.5 Natural Language Reasoning

615 *Overview and Setup.* In this task, the goal is to answer reasoning questions in natural language
 616 text. We consider two reasoning tasks: *date understanding*, which involves reasoning over dates
 617 and relative durations to determine the described date, and *arithmetic reasoning*, which involves
 618 solving word problems using basic arithmetic to compute desired quantities. For date under-
 619 standing, we use the BigBench benchmark’s date understanding task (Date) [et al. 2023] and for
 620 arithmetic reasoning we use GSM8K [Cobbe et al. 2021]. We use the top performing 7B parameter
 621 pretrained model from the Open LLM leaderboard which at the time of writing was Mistral-7B
 622 (`mistralai/Mistral-7B-v0.1`) [Jiang et al. 2023]. The prompt used for both tasks is the same
 623 prompt used by Lyu et al. [2023]. To decode the output of the LLM, we use a temperature of 0.4
 624 along with the default parameters to the Huggingface generation API.
 625

626 *Correctness Problems.* One of the major issues with LLMs is that the output generated is not
 627 guaranteed to conform to the conditions of the ground truths. For example, LLMs should generate
 628 valid dates in response to prompts in the date understanding task, and numbers for the arithmetic
 629 reasoning task. There have been elaborate prompting mechanisms devised to improve the reliability
 630 of LLMs such as chain-of-thought [Wei et al. 2022] and ReACT [Yao et al. 2023] and systems for
 631 programming these prompting mechanisms [Chase 2022].

632 *Integrity Constraints.* TORCHQL is useful for orchestrating language model prompting due to
 633 its highly expressive user-defined functions. Placing language model inference calls inside user-
 634 defined functions allows for queries which prompt language models and process the output. We

635
 636 ²<https://huggingface.co/lmsys/fastchat-t5-3b-v1.0>



(a) Constrained generation overview.

```

def gen_k(n, dataset, prompt, k=10):
    return Query(n, base=dataset)
        .project(parse_and_prepare_prompt)
        .project(lambda prompts, token: (llm(prompts, k),
                                         token), bs=10)
        .project(parse_outputs)
        .flatten()

def c_llm(n, dataset, prompt, cond, ans, k=10):
    return gen_k(n, dataset, prompt, ans, k)
        .filter(lambda prompt, gen: cond(gen))
        .group_by(lambda prompt, gen: prompt)
        .project(lambda prompt, gens: ans(gens))

```

(b) Constrained generation with TORCHQL.

Fig. 10. Using TORCHQL to orchestrate prompting of language models with output constraints.

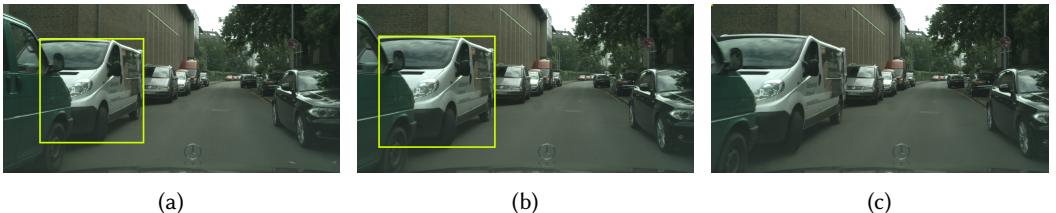


Fig. 11. One example of violating object consistency constraint where Figure 11a to Figure 11c are three continuous frames from CITYSCAPES dataset. The large truck detected in the first two frames is suddenly undetected in the third consecutive frame.

demonstrate the usefulness of this interface with a query to constrain the output of a language model similar to LMQL [Beurer-Kellner et al. 2023]. We then evaluate our prompting technique and show that TORCHQL can be used to constrain LLM output while maintaining accuracy as well as LMQL on two benchmark reasoning tasks.

In order to achieve this, we build a helper query, shown in Figure 10b, whose operation is shown in Figure 10a. This query samples multiple responses from an LLM and filters out the responses which violate the given constraints. The query uses four TORCHQL operations and many of the UDFs are single-lined lambda functions defined within the query specification, except for the implementation of the prompting semantics which include some use of regular expressions.

We use this helper query to write around 20 queries to enforce various constraints over the Mistral-7B LLM. We analyze two of them in this paper. First, we write the query L-Q1, which performs Chain-of-Thought (CoT) [Wei et al. 2022] prompting with constraints for the arithmetic reasoning task. The second query L-Q2 similarly attempts to constrain the LLM for the date understanding task. Both queries are explored in more detail in Section 6.5.

6 CASE STUDIES

In this section, we evaluate the expressivity of TORCHQL through a series of case studies for each of the above tasks. In each case study, we discuss the integrity constraints and queries from Table 3 and provide a detailed analysis of the issues discovered with those queries.

6.1 Case Study: Prediction Error Patterns in Self-Driving Object Detectors

Object Consistency Constraint. Section 2 describes a version of the object consistency constraint (S-Q1) specialized over pedestrians. S-Q1 finds 6,264 objects that are detected in two consecutive frames and disappear in the third. To evaluate the effectiveness of this query, we look at how well

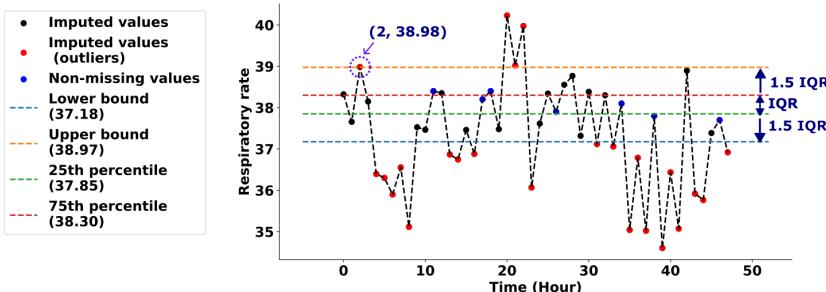


Fig. 13. Finding imputation errors as outliers to the range of expected respiratory rates for a patient. Using the non-missing values of respiratory rate, upper and lower bounds are calculated and all imputed values outside the lower-to-upper bound range are selected as imputation errors (shown as red dots).

6.2 Case Study: Imputed Value Monitoring for Time-Series Healthcare Data

We craft queries T-Q2 and T-Q3 to specify constraints on the smoothness assumptions in a univariate time series. We only focus on T-Q2 in this section as T-Q3 just differs in the way the threshold for detecting outliers is determined. T-Q2 joins the table containing imputation results with the ground truth table to obtain all the non-missing values and imputed values. We then employ a user-defined function to filter out imputation outliers over each univariate time series of every sample. This is accomplished by leveraging a threshold determined through established statistical methods for outlier detection [Tukey et al. 1977].

We illustrate this outlier detection method on a patient’s respiratory rate from the Physionet-2012 Challenge dataset shown in Figure 13. In this univariate time series, the non-missing respiratory rate values are denoted by blue dots while the imputed values (including outliers) are denoted by black or red dots. The outliers (denoted by red dots) are those imputed values outside the lower or upper bound shown in Figure 13. To determine these two bounds, we first calculate the interquartile range (IQR) of all the non-missing values, which is the range between the 25th percentile (denoted by LQ) and 75th percentile (denoted by UQ) of those values. They are then derived with the following formula:

$$\text{Lower bound} = LQ - 1.5 \cdot \text{IQR}, \text{Upper bound} = UQ + 1.5 \cdot \text{IQR},$$

As Figure 13 suggests, 24 outliers are identified from the imputed values with the constraint specified by T-Q2, most of which are all visually far away from their nearest non-missing values.

We further report the *recall* of T-Q2 over 20% randomly held-out non-missing entries, i.e., the portion of the entries with imputation errors covered by T-Q2. Ideally, T-Q2 covers a significant proportion of the total imputation errors allowing us to find where the model makes its worst errors and potentially allow for highly effective solutions. Hence, over the 20% held-out non-missing entries, we evaluate the ratio of the imputation errors covered by T-Q2 to the total imputation errors, denoted as *imputation error*. Six other queries were written to detect outliers over different variables using similar smoothness assumptions, and so we do the same to evaluate the recall and imputation error for all seven queries.

Ideally, both *recall* and *imputation error* approach 100%. The recall and imputation error are 40.15% and 54.27% respectively for all seven smoothness assumption queries, while they are 26.14% and 34.23% respectively just for T-Q2. The intermediate-level recall and imputation error of these queries justifies their validity and usefulness.

It is also worth noting that queries written in TorchSQL can reveal model prediction issues that cannot be captured by standard model performance metrics such as the Mean Square Error (MSE).

```

785 db = Database()
786 db.register(train_data, 'wilds_train')
787
788 imvar = Query('imvar', base='wilds_train')
789     .group_by(lambda im, lb, m: m[1].item())
790     .project(lambda seqid, r: [(seqid,
791         i-1, r[i-1][0], r[i-1][1], r[i-1][2],
792         i, r[i][0], r[i-1][1], r[i][2],
793         diff(r[i-1][0], r[i][0], 0.1, 0.15)
794     ) for i in range(1, len(r))])
795     .flatten()
796
797 imvar(db)

```

```

def is_empty(*row):
    diff = row[-1]
    mean = get_non_zero_mean(diff).item()
    return mean <= 0.13 and mean > 0

nonempty = Query('nonempty', base='imvar')
    .filter(lambda *args: not (
        args[3] == 0 and args[7] == 0))
non_empty(db)

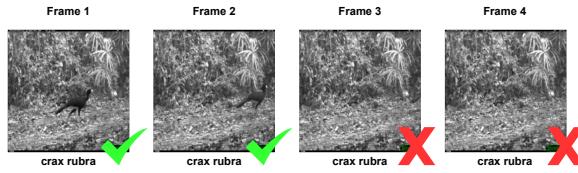
mislabeled = Query('mislabeled',
    base='nonempty').filter(is_empty)
mislabeled(db)

```

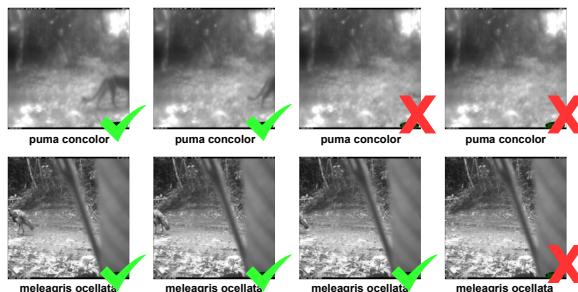
(a) Finding pixel differences between consecutive frames.

(b) Finding mislabeled frames.

Fig. 14. Query to find mislabeled frames in the iWildCam dataset.



(a) An example of mislabeled data found while manually exploring the iWildCam training data.



(b) Mislabeled data identified by the mislabeled query from Figure 14b. There are 280 such instances.

Fig. 15. Examples of mislabeled data from the iWildCam training dataset. Frames marked with a check are correctly labeled. Frames marked with a cross are those without an animal but still labeled non-empty. The correct label in such cases should be empty.

For example, for the entry at hour 2 in Figure 13 whose ground truth is 38.95, the imputed value at this entry is 38.98, which is marked as an outlier by T-Q2 since it is above the upper bound (38.97). However, the MSE metric cannot differentiate the cases where the imputed value is 38.98 or 38.92 for this entry since the difference between either value and the ground truth is the same. On the other hand, the imputed value 38.92 won't be flagged as a violation of T-Q2. One can therefore use T-Q2 to prefer an imputed value of 38.92 over 38.98, while the MSE is not fine-grained enough to determine which value is better.

6.3 Case Study: Mislabeled Data in iWildCam

We notice that in the iWildCam dataset, there are several cases where animals are present in the first few frames of a trap camera video, but then exit the frame. Once the animal exits, the subsequent frames should be labeled as empty. However, in many cases, these frames are labeled in

834 the ground truth as if the animal is still present. We show an example of this in Figure 15a, where
 835 the label for each frame is *leopardis pardalis*, despite the animal not being present in frame (d).

836 In most cases when animals are present in the frame, they tend to move around, creating a
 837 significant difference between the pixel values of consecutive frames. We therefore program the
 838 integrity constraint that any pair of consecutive frames with a small difference in their pixel values
 839 are most likely empty, and should therefore be labeled as such.

840 We construct this integrity constraint over a sequence of three queries shown in Figure 14. We
 841 begin by finding the difference in pixels over all pairs of two consecutive frames using the `imvar`
 842 query (Figure 14a). The nonempty query is then used to filter to frames not labeled as empty. We
 843 then determine a threshold to distinguish the pixel differences between frames with empty labels
 844 and the ones with non-empty labels. We use this threshold (0.13) to build the `mislabeled` query
 845 (Figure 14b) to detect visually empty frames incorrectly labeled as nonempty. More details about
 846 the queries are provided in Appendix A.2.

847 Overall, we detected 280 instances in the training data out of 69,972 pairs of consecutive frames.
 848 We show some examples of sequences containing such mislabeled data in Figure 15b. Since these
 849 mislabels occur within the ground truth, we manually inspect 20 random instances and find three
 850 false positives. In other words, there are only three cases out of the 20 where the frame satisfied
 851 this condition but actually had an animal in it. This suggests the constraint is accurate in capturing
 852 empty frames that are mislabeled as non-empty.

854 6.4 Case Study: Bias Discovery in LLMs

855 We write integrity constraints in TORCHQL to discover adjective-profession biases in the Alpaca
 856 instruction-tuning dataset [Taori et al. 2023].

857 *Detecting Biases in Datasets.* We first discover biases in the Alpaca dataset with query B-Q1:

```
858
859     Query('farmer_adj', base='alpaca')
860         .project(lambda instr, inp, outp: [instr, inp, outp, extract_noun_adj_pairs(outp)])
861         .filter(lambda instr, inp, outp, pairs: len(pairs) > 0)
862         .project(lambda instr, inp, outp, pairs: pairs)
863         .flatten()
864         .project(lambda noun, adj: [noun.lower(), adj.lower()])
865         .group_by(lambda noun, adj: noun)
866         .filter(lambda noun, adj: noun == "farmer")
```

867 Intuitively, this query finds all the adjectives used to describe the noun *farmer* in the Alpaca
 868 dataset. We find nouns and adjectives using a natural language parsing model from the Spacy library
 869 [Honnibal and Montani 2017]. This library is used by the `extract_noun_adj_pairs` function called
 870 within the query. We then write a similar query to find adjective-profession biases for “engineer”.
 871 We visualize adjectives associated with *farmer* and *engineer* in the Alpaca dataset in Figure 16.

872 The figure shows a clear difference between adjectives associated with farmers and engineers.
 873 We see undesirable differences tied to factors like age, race, and economic status (*old*, *African*, and
 874 *poor* for farmers), and variations in perception, such as engineers associated with *responsible*.

875 *Monitoring Biases in Model Outputs.* Now that we know the biases that exist in the Alpaca dataset,
 876 we program integrity constraints for monitoring the biases in the outputs of LLMs. We do so by
 877 building on top of the ‘farmer_adj’ and ‘engineer_adj’ queries. We specify that a language model
 878 should not favor exclusively using adjectives associated with farmers when describing farmers,
 879 or adjectives associated with engineers when describing engineers. To do this, we prompt LLMs
 880 like GPT-3.5 with the question “Answer with one word. The person is a farmer. Are they more
 881 likely *farmer_adj* or *engineer_adj*?", in which *farmer_adj* and *engineer_adj* are a pair of farmer and



Fig. 16. Visualizations of the adjectives associated with *farmer* and *engineer* from the Alpaca dataset are shown on the left. Biased responses of GPT-3.5 when queried with the circled adjectives are on the right.

engineer adjectives. We automate this process with the following TORCHQL queries which prompt GPT-3.5 in the above manner and then compute how often the biased answer is chosen:

```

98     Query('gpt35_bias_answers', base='farmer_adj')
99         .join('engineer_adj', key=lambda *args: 1, fkey=lambda *args: 1)
100        .project(lambda farmer_adj, engineer_adj: [
101            f"Answer with one word. The person is a farmer. \
102            Are they more likely {farmer_adj} or {engineer_adj}?" ... ])
103        .flatten()
104        .project(lambda prompt: gpt35(prompt))
105        .project(lambda prompt, answer: [*parse_prompt(prompt), parse_response(answer)])
106        .group_by(lambda target_job, adj1, adj2, ans: (target_job, frozenset({adj1, adj2})))
107        .project(lambda key, rows: [key[0], key[1], [row[3] for row in rows]])
108    Query('gpt35_bias_chosen', base='gpt35_bias_answers')
109        .project(lambda job, adjs, res: [[job, adjs, w] for w in res if w in adjs])
110        .flatten()

```

We collectively refer to the above queries as B-Q2. Note that the above queries were querying the GPT-3.5 LLM within the query itself using the `project` function. We run a similar query for the T5 language model as well, and then compare the biases exhibited by both models. Out of all the model responses that contain an adjective correlated with farmers or engineers, T5 selects the biased adjective 58.0% of the time and GPT-3.5 selects the biased adjective 59.6% of the time. For GPT-3.5, these biased responses cover 23 of the 27 farmer-associated adjectives and 34 of the 40 engineer-associated adjectives. We show an example of this on the right of Figure 16, where GPT-3.5 consistently generates biased adjectives. Similarly, the biased responses of T5 cover 11 of the 14 farmer adjectives and 15 of the 20 engineer adjectives. Note that a perfectly unbiased model would select the biased adjective 50% of the time. This indicates that this query is able to identify biased responses from LLMs like GPT-3.5 and T5.

6.5 Case Study: Constraining the Output of Language Models

The query described and shown in Figure 10 implements a general constrained prompting interface. Within the UDF `parse_and_prepare_prompt`, we implement a simple prompting interface based on that of LMQL where text in a prompt surrounded by curly braces, such as “{question}”, is replaced by elements from a table, and text appearing in double square brackets, such as “[[answer]]” is generated by the language model. A benefit of using TORCHQL in this case is that it provides in-house support for batching. Since the language model takes batched inputs, as do most machine learning models, we use a batch size of 10 shown in Figure 10b by passing the `bs=10` argument to TORCHQL’s `project` operation for automatic batch-wise parallelism. Performing similar batching with LMQL requires using Python’s asynchronous functionality and setting up a semaphore to limit the number of concurrent processes.

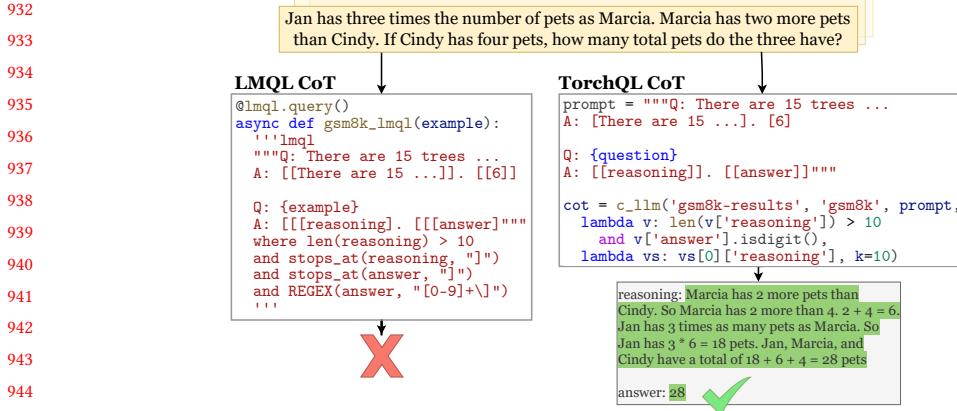


Fig. 17. Use of TORCHQL for constrained generation on the GSM8K arithmetic reasoning task compared to LMQL. The two methods can result in different responses on the same prompt. LMQL fails to produce any constraint-satisfying output while TORCHQL results in output which is a valid integer and the correct answer.

950	Method	Constraints	Date Understanding			Arithmetic Reasoning		
			Valid	Accuracy	Time (m)	Valid	Accuracy	Time (m)
951	CoT	No	96.75	52.85	2.6	81.88	37.83	9.35
952	LMQL CoT	Yes	<u>98.37</u>	52.85	68.03	<u>89.92</u>	<u>38.67</u>	226.12
953	TORCHQL CoT	Yes	99.73	53.39	<u>5.72</u>	97.80	44.28	<u>20.68</u>

Table 4. Performance results for using TORCHQL for language model prompting. Validity is the percent of all samples with a constraint-satisfying answer and accuracy is the task accuracy.

Unlike LMQL, we cannot perform any constraint-specific optimizations, since TORCHQL allows constraints to be arbitrary Python functions. Despite this, we show that for the tasks of arithmetic problem solving and date understanding, the TORCHQL queries still improve constraint satisfaction without hurting accuracy. Another benefit of our approach is that we can use black-box functions, including other models or even the Python interpreter. Using the `c_llm` function, we define 2 queries which perform constrained CoT prompting on the two tasks. The query performing constrained CoT prompting on arithmetic reasoning is L-Q1 and the query performing constrained CoT on date understanding is L-Q2 from Table 2. Query L-Q1 is shown on the right of Figure 17. L-Q2 is the same query with a modified prompt for the date understanding task and the constraint:

```
lambda args: len(args['reasoning']) > 10 and re.match(r"\d\d/\d\d/\d\d\d\d\d\d", args['answer']).
```

The full query is provided in Appendix A.2. We compare the TORCHQL query L-Q1 with its LMQL CoT counterpart in Figure 17. For the question shown in Figure 17, LMQL fails to produce any valid response while our method with TORCHQL produces a valid and correct response. This failure of LMQL is potentially caused by the decoding algorithm which may sometimes generate long streams of text without returning a valid answer.

The results for using TORCHQL to prompt an LLM for the GSM8K and Date dataset are in Table 4. The ‘Valid’ column shows the percent of samples where a constraint-satisfying answer is produced, and the ‘Accuracy’ column shows the task accuracy of the results. We see that this LLM constraint method with TORCHQL successfully constrains the LLM, resulting in better generation validity than LMQL for both tasks, as well as higher task accuracy. Finally, using TORCHQL for constrained

981	Query	Running time (seconds)				
		982	Python	ArangoDB	Pandas	TORCHQL
983	S-Q1	43.32±0.52	-	TO	44.51±0.42	
984	S-Q2	<u>11.27±0.87</u>	TO	32.89±0.20	<u>16.83±0.97</u>	
985	T-Q1	<u>0.11±0.02</u>	TO	0.13±0.03	0.11±0.04	
986	T-Q2	2.59±0.05	9.32±1.66	5.76±0.35	<u>2.65±0.38</u>	
987	T-Q3	0.90±0.05	57.38±1.89	1.93±0.55	<u>0.95±0.08</u>	
988	I-Q1	0.11±0.05	0.24±0.13	0.95±0.08	<u>0.18±0.01</u>	
989	B-Q1	0.012±0.02	-	0.085±0.03	<u>0.016±0.01</u>	

990 Table 5. Running time (seconds) for each system
991 (smaller is better). For each query, the quickest is
992 marked in bold, while the next quickest is underlined.
993 All queries had a timeout set to 10 minutes.

994	Query	Conciseness (number of tokens)				
		995	Python	ArangoDB	Pandas	TORCHQL
996	S-Q1	183	-	<u>175</u>	113	
997	S-Q2	<u>109</u>	207	113	97	
998	T-Q1	<u>123</u>	171	148	109	
999	T-Q2	<u>111</u>	174	120	99	
1000	T-Q3	<u>149</u>	167	183	140	
1001	I-Q1	75	64	43	<u>55</u>	
1002	B-Q1	104	-	<u>98</u>	89	

1002 Table 6. Conciseness comparison (smaller is better).
1003 For each query, we show the count of tokens of the
1004 query in each system. The smallest number is marked
1005 in bold and the next smallest is underlined.

1006 generation is up to 11x faster than using LMQL even though we used the same batch size and
1007 model, but only about 2x slower than unconstrained CoT.

7 QUANTITATIVE EVALUATION

1008 The previous section studied the expressiveness of TORCHQL across different use-cases. We now
1009 conduct quantitative experiments to evaluate the efficiency and conciseness of TORCHQL queries
1010 compared to their counterparts in baseline systems.

7.1 Setup

1011 We consider three baselines for this set of experiments: native Python, the standard language used for
1012 machine learning applications; ArangoDB [triAGENS GmbH [n. d.]], which is a document-oriented
1013 database and effectively an in-memory version of MongoDB; and Pandas [pandas development team
1014 2020], an in-memory data-analysis Python library. The OMG [Kang et al. 2018] model assertion
1015 system is another potential baseline, but it uses Python as its query language, so our Python
1016 comparison also serves as an OMG comparison.

1017 For each query in Table 2, with the exceptions of I-Q2, B-Q2, L-Q1, and L-Q2, we implement
1018 the query in TORCHQL as well as in each baseline system. We include all queries in Appendix
1019 A.3. In order to ensure a fair evaluation, we made the best effort to optimize the code within the
1020 style typical of each framework. For instance, we optimize the Python code so that its algorithmic
1021 complexity matches that of TORCHQL. We also attempt to use as few DataFrame operations in the
1022 Pandas implementations and as few database operations in the ArangoDB versions, as possible. We
1023 compare the different implementations of each query to evaluate its efficiency (by measuring the
1024 running time) in Table 5, and the conciseness (by measuring the number of tokens) in Table 6.

7.2 Results

1025 *Efficiency.* In all cases, TORCHQL queries are comparable or faster than their Python counterparts.
1026 This is expected since each Python query is optimized to match the algorithmic complexity of the
1027 TORCHQL version. The running time overhead introduced by TORCHQL’s querying abstractions
1028 is negligible for the most part, being less than one second for all queries except S-Q2, where it is
1029 around five seconds. This is still a small difference given the scale of the data for S-Q2. There is
1030 a trade-off for optimizing the Python code to such an extent: doing so drastically increases the
1031 required number of tokens, as evident from Table 6.

1032 Pandas times out on S-Q1 after 10 minutes due to the scale of the data, and is several seconds
1033 slower than TORCHQL in most other cases. The primary factor contributing to the inefficiency is the

1030 rigorous schema requirement of Pandas. Pandas is optimized for cases where the DataFrames have
 1031 a fixed structure with primitive datatypes in cells. As a result, additional operations are needed to
 1032 explicitly produce new attributes essential for the downstream operations but not in the schema.
 1033 Since these operations tend to modify the DataFrame itself, they introduce substantial overhead.

1034 ArangoDB is even slower and times out on S-Q2 and T-Q1. The reason is that for queries like S-Q2
 1035 which join two datasets on attributes that are not indexed, ArangoDB has to perform full nested
 1036 scans on the datasets for the join operations (see Figure 19d where the join is conditioned over
 1037 item1.frames[0] == item2.frames[0] but frames[0] is not indexed). On the other hand, TORCHQL
 1038 can construct indexes over such newly created attributes on the fly without extra operations.

1039 Overall, TORCHQL achieves at least a 13x and three order-of-magnitude speed-ups in with respect
 1040 to Pandas and ArangoDB respectively in the best case (S-Q1 for Pandas and T-Q1 for ArangoDB).
 1041 Note that Pandas timed out for S-Q1, which is why the 13x speed up in this case is a lower estimate.

1042 *Conciseness.* Table 6 shows the conciseness metric evaluated on each query written in different
 1043 systems. TORCHQL significantly reduces the number of tokens needed for writing a query (by up to
 1044 40%) with respect to native Python. This can substantially reduce users’ effort in the process of
 1045 iteratively writing and refining queries for specifying appropriate integrity constraints. Furthermore,
 1046 TORCHQL is more concise than Pandas for all cases except I-Q1, and more concise than ArangoDB
 1047 in general. This is because these baselines require substantial data wrangling so that their rigid
 1048 schemas can support arbitrary Python objects, which is not needed for TORCHQL.
 1049

1050 8 USER STUDY

1051 We conduct a small-scale user study to validate the usability of TORCHQL for programming integrity
 1052 constraints over complex forms of data. The study consists of 10 participants, including 6 graduate
 1053 students, 3 undergraduate students, and a systems engineer.
 1054

1055 8.1 Setup

1056 For the user study, each participant was provided with a tutorial of TORCHQL in a Jupyter notebook
 1057 (an interactive Python environment) and then asked to write three TORCHQL queries of increasing
 1058 complexity over the iWildCam dataset. All three tasks, detailed in Appendix A.5, focused on finding
 1059 model prediction errors. The first two tasks allowed them to make use of the supplied ground truth
 1060 labels. The third task, on the other hand, was more open-ended in nature, and required them to
 1061 program an integrity constraint that flagged erroneous predictions without using the ground truth.
 1062 No time limit was enforced on any task, though we measured the time it took users to complete
 1063 each task. We also measure the precision and recall of the responses of the users for the third task.
 1064

1065 8.2 Results

1066 Overall, users completed the first two tasks in an average of 3.41 and 20.45 minutes respectively. A
 1067 total of seven of the 10 users completed the third and most complex task, where they had to find
 1068 model mispredictions without using ground-truth labels. Their queries had an average precision
 1069 of 88.0% and recall of 71.0% for finding model mispredictions, compared to our solution with a
 1070 precision and recall of 91.0% and 78.0% respectively, and were written in an average of 54.33 minutes.
 1071 The large difference in average time taken for the three tasks reflects their differing complexity.
 1072 These aggregate numbers show the usability of TORCHQL despite the users’ lack of prior familiarity.
 1073

1074 Looking at individual users’ queries reveals where people had difficulty and general trends in the
 1075 use of TORCHQL. The first task required a single filter operation and all users correctly utilized
 1076 this operation to perform the first task. The second task required the use of a group_by followed by
 1077 an order_by operation, and seven users successfully used these two operations and composed them
 1078

1079 in the correct way. The other three users all correctly used the `group_by` operation, but resorted to
1080 manually iterating through the table with a `for` loop to determine the largest and smallest element
1081 of the table. Of the seven responses for the third task, five were exactly equivalent to our solution,
1082 one achieved within 0.01 F1 score of our solution, and one was significantly different from our
1083 solution and performed much worse. All seven solutions matched the overall structure in terms of
1084 TORCHQL operations of our query, but differed in their definition of UDFs. The five solutions that
1085 matched ours used UDFs that were expressed differently but resulted in the same overall operation,
1086 and the other two queries solved the problem differently by choosing different UDFs.

1087 9 RELATED WORK

1088 *Databases for machine learning.* Techniques from databases have been used to enhance various
1089 aspects of machine learning. Ideas from operation scheduling have been used to select appropriate
1090 hardware [Mirhoseini et al. 2017] or find semantically equivalent but more efficient operations to
1091 optimize deep learning computations [Jia et al. 2019]. Recomputing and swapping techniques from
1092 databases have been adopted to manage memory in deep learning frameworks [Pleiss et al. 2017;
1093 Wang et al. 2018]. Compression techniques and communication frameworks have been proposed to
1094 accelerate parameter servers in distributed training of large models [Goyal et al. 2017; Jiang et al.
1095 2017]. These approaches focus on the speed and efficiency of deploying models.

1096 *Integrity constraints in machine learning.* Integrity constraints have been used in machine learning
1097 for data analysis and cleaning, model debugging, verification, and data generation. For instance,
1098 “unit testing” data [Schelter et al. 2019] is proposed to assess data quality in data analysis and native
1099 Python assertions have been employed to capture model output bugs [Kang et al. 2018]. In addition,
1100 there has been extensive research in neural model verification [Liu et al. 2021], which focuses
1101 on verifying whether a given input-output constraint holds for a neural network. Finally, data
1102 generation languages such as Scenic [Fremont et al. 2019] construct data to test the integrity of a
1103 machine learning system or to describe correct or incorrect model behaviors [Kim et al. 2020].

1104 *Finding errors in machine learning models.* Various tools have been developed to identify errors
1105 in the machine learning systems, which can be operated in either passive or active mode. In passive
1106 mode, tools can either generate test samples to trigger model errors [Tian et al. 2018] or produce
1107 explanations to assist users in understanding the model’s prediction process. These explanations
1108 can be in various forms, such as the coefficients of the sparsified neural network layers [Wong
1109 et al. 2021], extracted robust features [Singla and Feizi 2021; Singla et al. 2021], and influence
1110 functions [Ilyas et al. 2022; Koh and Liang 2017; Salman et al. 2022], which can then be compared
1111 against domain knowledge to discover model prediction errors. They can even be integrated into
1112 the subsequent model debugging loop [Anders et al. 2022; Bhadra and Hein 2015; Cadamuro et al.
1113 2016; Kulesza et al. 2015]. Apart from that, users can actively provide rules [Kang et al. 2018] based
1114 on their knowledge to verify the correctness of model predictions. Except for the tools for finding
1115 errors in general machine learning systems, some tools have emerged for specific settings such as
1116 federated learning systems [Augenstein et al. 2019].

1117 10 CONCLUSION

1118 We introduced a new framework TORCHQL for programming and checking integrity constraints
1119 in machine learning applications. We showed that TORCHQL is able to help find errors in several
1120 tasks across various domains up to 13x faster than other querying systems while being up to 40%
1121 more concise than regular Python programs. We also validated the usability of TORCHQL via a user
1122 study. In the future, we intend to extend TORCHQL to further improve performance by parallelizing
1123 operations and to support synthesizing queries from natural language descriptions.

1128 DATA-AVAILABILITY STATEMENT

1129 The code for TORCHQL is available as a part of the supplementary material along with the data that
 1130 was used to evaluate it.

1132 REFERENCES

- 1133 Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings*
 1134 *of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 298–306.
- 1135 Christopher J Anders, Leander Weber, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian Lapuschkin.
 1136 2022. Finding and removing Clever Hans: using explanation methods to debug and improve deep models. *Information*
 1137 *Fusion* 77 (2022), 261–295.
- 1138 Sean Augenstein, H Brendan McMahan, Daniel Ramage, Swaroop Ramaswamy, Peter Kairouz, Mingqing Chen, Rajiv
 1139 Mathews, et al. 2019. Generative models for effective ML on private, decentralized datasets. *arXiv preprint arXiv:1911.06679*
 1140 (2019).
- 1141 Kyle Bunker, Douglas Garrett, Peter Bakkum, and Shaun Verch. 2016. *MongoDB in action: covers MongoDB version 3.0*. Simon
 1142 and Schuster.
- 1143 Sara Beery, Elijah Cole, and Arvi Gjoka. 2020. The iWildCam 2020 Competition Dataset. *arXiv preprint arXiv:2004.10340*
 1144 (2020).
- 1145 Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2023. Prompting is programming: A query language for large
 1146 language models. *Proceedings of the ACM on Programming Languages* 7, PLDI (2023), 1946–1969.
- 1147 Sahely Bhadra and Matthias Hein. 2015. Correction of noisy labels via mutual consistency check. *Neurocomputing* 160
 1148 (2015), 34–52.
- 1149 Daniel Bolya, Sean Foley, James Hays, and Judy Hoffman. 2020. Tide: A general toolbox for identifying object detection
 1150 errors. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III*
 1151 16. Springer, 558–573.
- 1152 Gabriel Cadamuro, Ran Gilad-Bachrach, and Xiaojin Zhu. 2016. Debugging machine learning models. In *ICML Workshop on*
 1153 *Reliable Machine Learning in the Wild*, Vol. 103.
- 1154 Harrison Chase. 2022. *LangChain*. <https://github.com/langchain-ai/langchain>
- 1155 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry
 1156 Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint*
 1157 *arXiv:2110.14168* (2021).
- 1158 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke,
 1159 Stefan Roth, and Bernt Schiele. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of*
 1160 *the IEEE conference on computer vision and pattern recognition*. 3213–3223.
- 1161 Wenjie Du, David Côté, and Yan Liu. 2023. Saits: Self-attention-based imputation for time series. *Expert Systems with*
 1162 *Applications* 219 (2023), 119619.
- 1163 Aarohi Srivastava et al. 2023. Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models.
 1164 *Transactions on Machine Learning Research* (2023). <https://openreview.net/forum?id=uyTL5Bvosj>
- 1165 Daniel J Fremont, Tommaso Dreossi, Shromona Ghosh, Xiangyu Yue, Alberto L Sangiovanni-Vincentelli, and Sanjit A
 1166 Seshia. 2019. Scenic: a language for scenario specification and scene generation. In *Proceedings of the 40th ACM SIGPLAN*
 1167 *Conference on Programming Language Design and Implementation*. 63–78.
- 1168 Parke Godfrey, John Grant, Jarek Gryz, and Jack Minker. 1998. Integrity constraints: Semantics and applications. In *Logics*
 1169 *for databases and information systems*. Springer.
- 1170 Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing
 1171 Jia, and Kaiming He. 2017. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*
 1172 (2017).
- 1173 Shreya Havaldar, Bhumika Singhal, Sunny Rai, Langchen Liu, Sharath Chandra Guntuku, and Lyle Ungar. 2023. Multilingual
 1174 Language Models are not Multicultural: A Case Study in Emotion. In *Proceedings of the 13th Workshop on Computational*
 1175 *Approaches to Subjectivity, Sentiment, & Social Media Analysis*. 202–214.
- 1176 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings*
 1177 *of the IEEE conference on computer vision and pattern recognition*. 770–778.
- 1178 Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional
 1179 neural networks and incremental parsing. (2017). To appear.
- 1180 Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. 2022. Datamodels: Predicting
 1181 predictions from training data. *arXiv preprint arXiv:2202.00622* (2022).
- 1182 Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. 2022. OneFormer: One Transformer to
 1183 Rule Universal Image Segmentation. *arXiv* (2022).

- 1177 Zhihao Jia, James Thomas, Todd Warszawski, Mingyu Gao, Matei Zaharia, and Alex Aiken. 2019. Optimizing DNN
1178 computation with relaxed graph substitutions. *Proceedings of Machine Learning and Systems* 1 (2019), 27–39.
- 1179 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian
1180 Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825*
(2023).
- 1181 Jiawei Jiang, Bin Cui, Ce Zhang, and Lele Yu. 2017. Heterogeneity-aware distributed parameter servers. In *Proceedings of the*
1182 *2017 ACM International Conference on Management of Data*. 463–478.
- 1183 Daniel Kang, Deepti Raghavan, Peter Bailis, and Matei Zaharia. 2018. Model assertions for debugging machine learning. In
1184 *NeurIPS MLSys Workshop*, Vol. 3. 10.
- 1185 Edward Kim, Divya Gopinath, Corina Pasareanu, and Sanjit A Seshia. 2020. A programmatic and semantic approach to
1186 explaining and debugging neural network based object detectors. In *Proceedings of the IEEE/CVF Conference on Computer*
Vision and Pattern Recognition. 11128–11137.
- 1187 Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *International conference*
1188 *on machine learning*. PMLR, 1885–1894.
- 1189 Todd Kulesza, Margaret Burnett, Weng-Keen Wong, and Simone Stumpf. 2015. Principles of explanatory debugging to
1190 personalize interactive machine learning. In *Proceedings of the 20th international conference on intelligent user interfaces*.
126–137.
- 1191 Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring Bias in Contextualized Word
1192 Representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*. 166–172.
- 1193 Chonho Lee, Zhaojing Luo, Kee Yuan Ngiam, Meihui Zhang, Kaiping Zheng, Gang Chen, Beng Chin Ooi, and Wei Luen James
1194 Yip. 2017. Big healthcare data analytics: Challenges and applications. *Handbook of large-scale distributed computing in*
smart healthcare (2017), 11–41.
- 1195 Tao Li, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Vivek Srikumar. 2020. UNQOVERing Stereotyping Biases via
1196 Underspecified Questions. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. 3475–3489.
- 1197 Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating
1198 social biases in language models. In *International Conference on Machine Learning*. PMLR, 6565–6576.
- 1199 Changliu Liu, Tomer Arnon, Christopher Lazarus, Christopher Strong, Clark Barrett, and Mykel J Kochenderfer. 2021.
Algorithms for Verifying Deep Neural Networks. *Foundations and Trends® in Optimization* 4, 3-4 (2021), 244–404.
- 1200 Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch.
2023. Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2301.13379* (2023).
- 1201 Azalia Mirhoseini, Hieu Pham, Quoc V Le, Benoit Steiner, Rasmus Larsen, Yuefeng Zhou, Naveen Kumar, Mohammad
1202 Norouzi, Samy Bengio, and Jeff Dean. 2017. Device placement optimization with reinforcement learning. In *International*
Conference on Machine Learning. PMLR, 2430–2439.
- 1203 Aaditya Naik, Yinjun Wu, Mayur Naik, and Eric Wong. 2023. Do Machine Learning Models Learn Common Sense? *arXiv*
1204 *preprint arXiv:2303.01433* (2023).
- 1205 Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. CrowS-Pairs: A Challenge Dataset for Measuring
1206 Social Biases in Masked Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*
1207 *Processing (EMNLP)*. 1953–1967.
- 1208 The pandas development team. 2020. *pandas-dev/pandas: Pandas*. <https://doi.org/10.5281/zenodo.3509134>
- 1209 Geoff Pleiss, Danlu Chen, Gao Huang, Tongcheng Li, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Memory-
1210 efficient implementation of densenets. *arXiv preprint arXiv:1707.06990* (2017).
- 1211 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J
1212 Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine*
1213 *Learning Research* 21, 1 (2020), 5485–5551.
- 1214 Hadi Salman, Saachi Jain, Andrew Ilyas, Logan Engstrom, Eric Wong, and Aleksander Madry. 2022. When does Bias Transfer
1215 in Transfer Learning? *arXiv preprint arXiv:2207.02842* (2022).
- 1216 Sebastian Schelter, Felix Biessmann, Dustin Lange, Tammo Rukat, Philipp Schmidt, Stephan Seufert, Pierre Brunelle, and
1217 Andrey Tapuntsov. 2019. Unit testing data with deeque. In *Proceedings of the 2019 International Conference on Management*
of Data. 1993–1996.
- 1218 Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. 2012. Predicting in-hospital mortality of icu
1219 patients: The physionet/computing in cardiology challenge 2012. In *2012 Computing in Cardiology*. IEEE, 245–248.
- 1220 Sahil Singla and Soheil Feizi. 2021. Salient ImageNet: How to discover spurious features in Deep Learning? *arXiv preprint*
arXiv:2110.04301 (2021).
- 1221 Sahil Singla, Besmira Nushi, Shital Shah, Ece Kamar, and Eric Horvitz. 2021. Understanding failures of deep networks
1222 via robust feature extraction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
1223 12853–12862.
- 1224
- 1225

- 1226 Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. 2022. “I’m sorry to hear that”:
1227 Finding New Biases in Language Models with a Holistic Descriptor Dataset. In *Proceedings of the 2022 Conference on*
1228 *Empirical Methods in Natural Language Processing*. 9180–9211.
- 1229 Rohan Taori, Ishaaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B.
1230 Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.
- 1231 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. 2021. CSDI: Conditional score-based diffusion models for
1232 probabilistic time series imputation. *Advances in Neural Information Processing Systems* 34 (2021), 24804–24816.
- 1233 Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. 2018. DeepTest: Automated Testing of Deep-Neural-Network-Driven
1234 Autonomous Cars. In *Proceedings of the 40th International Conference on Software Engineering* (Gothenburg, Sweden)
1235 (*ICSE ’18*). Association for Computing Machinery, New York, NY, USA, 303–314. <https://doi.org/10.1145/3180155.3180220>
- 1236 triAGENS GmbH. [n. d.]. ArangoDB Query Language (AQL) Introduction: ArangoDB Documentation. www.arangodb.com.
- 1237 John W Tukey et al. 1977. *Exploratory data analysis*. Vol. 2. Reading, MA.
- 1238 Daisuke Wakabayashi. 2018. Self-driving uber car kills pedestrian in Arizona, where Robots Roam. <https://www.nytimes.com/2018/03/19/technology/uber-driverless-fatality.html>
- 1239 Linnan Wang, Jinmian Ye, Yiyang Zhao, Wei Wu, Ang Li, Shuaiwen Leon Song, Zenglin Xu, and Tim Kraska. 2018.
1240 Superneurons: Dynamic GPU memory management for training deep neural networks. In *Proceedings of the 23rd ACM*
1241 *SIGPLAN symposium on principles and practice of parallel programming*. 41–53.
- 1242 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022.
1243 Self-Instruct: Aligning Language Model with Self Generated Instructions. *arXiv preprint arXiv:2212.10560* (2022).
- 1244 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-
1245 thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* 35
1246 (2022), 24824–24837.
- 1247 Benjamin Wilson, Judy Hoffman, and Jamie Morgenstern. 2019. Predictive inequity in object detection. *arXiv preprint*
1248 *arXiv:1902.11097* (2019).
- 1249 Eric Wong, Shibani Santurkar, and Aleksander Madry. 2021. Leveraging sparse linear layers for debuggable deep networks.
1250 In *International Conference on Machine Learning*. PMLR, 11205–11216.
- 1251 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing
1252 Reasoning and Acting in Language Models. In *The Eleventh International Conference on Learning Representations*.
- 1253 Cheng You, Dennis KJ Lin, and S Stanley Young. 2018. Time series smoother for effect detection. *PLoS one* 13, 4 (2018),
1254 e0195360.
- 1255 John R Zech, Marcus A Badgeley, Manway Liu, Anthony B Costa, Joseph J Titano, and Eric Karl Oermann. 2018. Variable
1256 generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study.
1257 *PLoS medicine* 15, 11 (2018), e1002683.
- 1258 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men Also Like Shopping: Reducing
1259 Gender Bias Amplification using Corpus-level Constraints. In *Proceedings of the 2017 Conference on Empirical Methods in*
1260 *Natural Language Processing*. 2979–2989.
- 1261
- 1262
- 1263
- 1264
- 1265
- 1266
- 1267
- 1268
- 1269
- 1270
- 1271
- 1272
- 1273
- 1274