

Execution Trace Based Feature Engineering To Enable Formative Feedback on Visual, Interactive Programs

Offering students immediate, formative feedback when they are programming can increase students' learning outcomes and self-efficacy. However, visual and interactive programs include dynamic user input and visual outputs that change over time, making it difficult to automatically assess students' code with traditional functional tests to offer this feedback. In this work, we introduce Execution Trace Based Feature Engineering (ETF), a feature engineering approach that extracts sequential patterns from execution traces, which capture the runtime behavior of students' code. We evaluated ETF on 162 students' code snapshots from a Pong game assignment in an introductory programming course, on a challenging task to predict students' success on fine-grained rubrics. We found that ETF achieves an average F_1 score of 0.93 over 10 grading rubrics, which is 0.1–0.2 higher than a high-performing syntax-based feature engineering approach from prior work. These results show that ETF has strong potential to be used for code classification, to enable formative feedback for students' visual, interactive programs.

ACM Reference Format:

. 2021. Execution Trace Based Feature Engineering To Enable Formative Feedback on Visual, Interactive Programs. 1, 1 (July 2021), 17 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Real-time, formative feedback promotes students' learning gains and self-efficacy [8, 23, 36], for it allows students to judge how well they are meeting assignment requirements, and make modifications to their code, both as they develop it, and before submitting it. To provide such formative feedback in real-time, CS instructors commonly write test cases, allowing students to run their code against these test cases when programming [10, 14, 19, 27]. However, visual, interactive programming projects, such as creating apps and games, include dynamic user interactions, and visual outputs that change over time, making it challenging to use test cases to assess these programs [37, 38].

In contrast to test case-based approaches, data-driven methods allow instructors to offer formative feedback by grading a smaller set of programs instead of writing test cases [29, 35, 42, 47]. These data-driven methods start with transforming code into input vectors using *feature engineering*, typically by extracting syntax elements from an abstract syntax tree (AST), where nodes and their children correspond to specific code elements (e.g., if statements). However, when applying these syntax-based AST feature extraction techniques to classify programs based on fine-grained assignment rubrics, prior work showed mixed results, which are often not high enough to ensure the quality of student feedback [1, 2, 35].

A known key limitation of such AST-based feature engineering approaches is that code is first and foremost “executable”, and that its execution traces include dynamic data about the program functionality, which cannot be directly captured by the syntax-based features [3, 26]. Some prior work has used execution traces to classify students' sorting programs based on their specific strategies, and has shown that execution-trace-based classification achieved higher accuracies than a syntax-based classification approach [26]. However, no prior work has conducted feature extraction

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2021 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

on the execution trace of visual, interactive programs, which include dynamic user interactions and various changing outputs.

In this work, we explore extracting useful features from execution traces that capture the runtime behavior of visual, interactive programs. We designed an execution trace-based feature engineering approach (ETF) to transform students' source code into feature vectors, for classification algorithms to build models based on rubric-based labels (e.g., the presence of a key-triggered movement). We evaluated ETF by classifying 162 students in-progress and submitted code snapshots. We found it to achieve high prediction performance with an average of 0.93 F_1 score over 10 grading rubric items, which is 0.1–0.2 higher than a high-performing syntax-based code classification approach. Our work has the following contributions: 1) We designed and implemented a novel, execution trace-based feature engineering (ETF) approach to extract temporal patterns in students' visual, interactive programs (Section 3); 2) We evaluated the ETF approach on students' code snapshots for a widely-used, representative visual, interactive program assignment (Section 4.1); 3) We investigated the effects of the amount of labeled data available on the performance of the ETF approach, using a cold start analysis (Section 4.2).

2 RELATED WORK

Students can easily get stuck during programming [17], and may need help in order to proceed. Such timely feedback that is provided in the middle of programming is called “**formative feedback**”. Unlike summative feedback, offered after a student completing a programming task, formative feedback aims to monitor students' on-going learning process and provide timely assessment and suggestions [13]. Receiving such formative feedback has been shown to improve students' learning outcomes [8] and self-efficacy [23]; and a systematic review by Shute has shown that real-time feedback is critical to learning [36]. In programming education, such formative feedback is built on an accurate assessment of students' current progress based on their code, which can take the form of analyzing students' progress against instructor rubrics [35]. With the growing size of CS classrooms [6], such assessment is challenging to offer manually.

To enable formative feedback, many learning systems make use of functional tests; some prior work also researched data-driven code classification. We first review prior work on test-case-based assessment to understand its benefits and limitations. Next, we describe prior work on data-driven code classification, which inspired our work.

1) **Test-case-based Automated Assessment.** Many learning systems offer test-case-based automated assessments, where students are allowed to submit their solutions multiple times. With every submission, the system checks students' solutions through a set of test cases and sends the report to the students (e.g., [10, 14, 19, 27]). However, visual and interactive programs include dynamic user interactions and program properties that change over time [37, 38]. Prior work that has applied functional tests to analyze these programs identified shared challenges for instructors to author such test cases, which requires temporal logic-based specifications of program behaviors and time bounds, using a domain-specific language [12, 38].

2) **Data-driven Student Code Classification.** Prior work has made use of data-driven methods to classify student programs for enabling formative feedback. These methods take as training data a set of graded programs and predict students' success or failures on a given assignment. To do that, prior work used syntax-based feature engineering approaches to extract features from the code abstract syntax tree (AST) and feed them into machine learning algorithms. We describe their methods and applications below.

2.a) *Methods.* Syntax-based approaches extract patterns inside a code AST, use the presence or absence of a feature, or the count of the feature to generate input vectors. As an example, we explain a recently-applied AST n -Gram feature

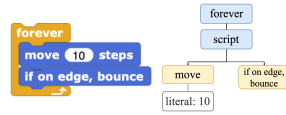


Fig. 1. A horizontal $n(2)$ -Gram (in yellow) and a vertical $n(2)$ -Gram (in blue).

extraction approach [2] by making an analogy to Natural Language Processing (NLP): In NLP, an n -Gram with $n = 1$ is a 1-Gram feature taken from each word; and an n -Gram feature takes a continuous sequence of n words to extract relationship between words (e.g., orders); in an AST, 1-Gram features represent each *node* in students' code. And to extract structural relationships, Akram et al. [2] designed the use of n -Grams to represent n -length sequences of code, where a *vertical n -Gram* is created by a depth-first search of leaves; a *horizontal n -Gram* is created by a breadth-first search of all direct children of each AST nodes (e.g., in Figure 1) [2].

Other AST feature engineering approaches were introduced with the usage of Neural Network models, which commonly apply an embedding layer that takes as input multi-dimensional code representations (e.g., a matrix) [3, 15, 25, 35]. In addition, many embedded approaches use sequential models, so that the embedding is trained based on relative locations of code elements [15, 25]. Depending on the Neural Network's architecture, these embedded approaches also vary in ways to represent code as vectors to feed into the models. For example, Alon et al. used leaf-to-leaf feature extraction, connecting the shortest path from each two leaf nodes [4, 5]. They used this approach to represent code and used it for an attention-based model (*code2vec*) to predict method names. The results show that it achieved relatively high accuracy (0.58 F_1 score), and that it learns meaningful vector embeddings that reveal semantic similarities [5].

2.b) *Applications on Student Code Classification.* 1-Gram and n -Gram-based AST feature extraction approaches have both been applied for student code classification tasks. For example, Azcona et al. applied 1-Gram feature extractions as a baseline to extract features in 591,707 snapshots of students' Python code, and achieved 0.598 F_1 in predicting students' overall success in completing programming problems, using a simple Naive Bayes model [7]. Compared to 1-Gram feature extractions, prior work has shown that n -Grams provide more predictive features for code analysis. For example, Akram et al. used n -Grams with n ranging from 1 to 4 to extract features, and used a Gaussian Process model to infer scores on 642 students' code pieces, in a block-based programming environment. This achieved an R-square of 0.94, higher than the 0.88 achieved by the baseline 1-Gram approach [1, 2].

Most prior work has used data-driven methods to predict students' success or failures of **the entire problem** [7, 21, 45]. However, to offer students informative feedback, it is important to understand what specific rubric items a student has succeeded or failed in. Some prior work has tried to predict students' success on **specific rubrics**, e.g., by 1) generating rubrics based on a set of correct programs [9, 48]; or by 2) predicting success/failures on rubric items based on training data of a set of graded programs. For example, Shi et al. applied *code2vec* on 207 students' program submissions to predict their success on six different rubric items, in an introductory programming assignment [35]. They found that *code2vec* achieved an average F_1 of 0.69 over all assignments, while a simple SVM model on a 1-Gram feature extraction achieves an average of 0.46 F_1 score. This work applied data-driven code classification in a real-world, practical setting: in practice, instructors do not have a large amount of training data at hand to train a data-hungry machine learning model (e.g., such as in [7] (591K) and [5] (14M)). [35] shows the potential of applying data-driven code classification not only on a relatively small dataset, but also to predict fine-grained, rubric-specific grades.

3) **Representing Programs by Execution Traces.** In a systematic literature review on machine-learning-based code analysis, Allamanis et al. critiques that **executability** contributes the key difference between code and natural language, and that **execution traces**, which log variable states during each step of program execution, add an important dimension of analyzing programs by presenting a dynamic viewpoint, and directly linking a piece of programming code to its functionality [3]. In the domain of student code analysis, prior work has identified similar limitations of syntax-based approaches, and pointed to the lack of work of execution-trace-based code analysis [26, 28]. Consequently, recent work that applied execution traces-based code analysis has revealed promising results [26, 44]. For example, PaaÅsen et al. used execution-trace-based distance measures to classify programs into different strategies (e.g., bubble sort v.s. Insertion sort), and found that execution-trace-based classification achieving 90% accuracy, higher than the 80% accuracy achieved by syntax-based approaches [26]. This shows the potential of using execution traces to classify students' programming code.

3 THE ETF APPROACH

The goal of ETF is to extract useful and relevant features from execution traces. It is a feature engineering approach that leverages n -Gram-based feature extraction, which has not previously been applied on execution traces for student code classification. ETF starts from collecting a set of students' programming code, along with a class label for each piece of code, generated from instructor gradings (i.e., positive or negative), and is specifically designed to conduct automated code classification on programs that have the following properties:

- (1) **Dynamic user interaction.** Programs respond to user interactions (e.g., mouse, keyboard). For example, in games, users use a mouse or keyboard to control actor movements; in an app, users need to navigate to different pages / options using mouse, or keyboard input.
- (2) **Object-specific program states.** Program states can be represented as properties of objects (e.g., sprites¹) such as positions and rotations, which correspond to visual output on the screen.
- (3) **Properties that can change over time.** Program behaviours can be *a function of time* (e.g., movement is described by the actor position that changes over time); in addition, a behaviour can also *change over time* (e.g., move and then stop [37].)

These properties are shared by many different types of visual, interactive program projects, such as games, simulations and apps, and can be easily created in many novice programming environments, such as SCRATCH [32], Snap! [24], and Greenfoot [18], as well as programming environments for creative practitioners, such as Processing [31]. ETF conducts feature engineering on these types of visual, interactive programs in four steps, the first three corresponding to the three features above: 1) generating execution traces that provide the inputs for dynamic user interaction (Section 3.1); 2) summarizing traces to include program state-based properties (Section 3.2); 3) extracting features that describe relevant properties that change over time (Section 3.3). The last step 4) is to filter features to generate a \vec{x} vector for each student's code snapshot (Section 3.4).

An Example Assignment. As an example, consider the Pong assignment, which we will use later in our experiments. Pong consists of a paddle sprite and a ball sprite. The ball moves around the stage², and a player can use the keyboard to control the up and down movement of the paddle to catch the running ball. If the paddle catches the ball, the player score increases; but if the paddle misses the ball and the ball hits the back wall, the game ends. The Pong

¹The word "sprite" is used in block-based programming environments, such as Snap! and SCRATCH, to represent an object that has its own code (called scripts) and costumes. For example, a sprite can be a game actor or an app widget.

²In Snap! and SCRATCH, a stage is a screen on which the program shows its sprite actions.

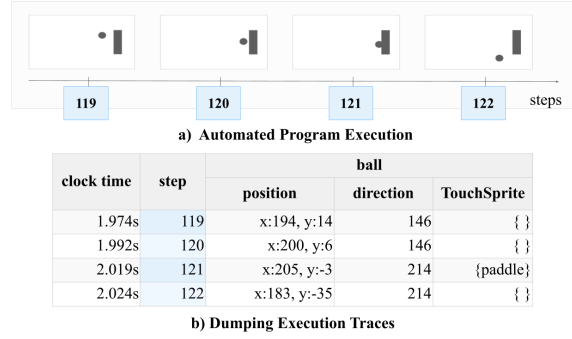


Fig. 2. Step 1: Generating Execution Traces.

assignment is commonly used in a variety of introductory programming courses [46] and camps [20, 34], using various code editors, such as Snap! [46], NetsBlox [20], and App Inventor [34], and is therefore representative of many learning contexts and environments. This example assignment includes the three properties of visual, interactive programs that we designed ETF to analyze. We describe the ETF approach below.

3.1 Step 1: Generating Execution Traces

Visual, interactive programs include program states that can be represented as *properties* that change over time [37]. For example, in Snap!, these properties can include:

- (1) time: #milliseconds spent from the start of program execution;
- (2) global variables: the names and values of global variables;
- (3) sprite-specific properties: properties that are related to specific sprites, such as x, y coordinates, directions, which sprites the current sprite is touching, which stage edge the sprite is touching, the sprite size, as well as the names and values of the sprite's local variables. The dumped execution trace tables use sprite names to label sprite-specific properties (e.g., to distinguish `ball.x` from `paddle.x`³).
- (4) inputs: automated execution of Snap! programs requires user inputs such as mouse and key pressing. ETF uses an existing test harness tool called SNAPCHECK [40], which instrument a program editor to simulate user inputs and to execute programs automatically [38]. SNAPCHECK allows manually authorizing input generation rules to automatically execute student programs [40], such inputs can be responsive to program states. For example, the rule “when `ball.y < paddle.y`, press down key” allows the paddle to follow the ball movement. Such inputs are also encoded as program states in a Snap! program, including properties such as key pressed, mouse clicked, and mouse positions on stage.

Systems such as Snap! and SCRATCH make use of *step functions* to update the above properties based on the current properties and the current user inputs [38, 40] in the intervals of milliseconds. We instrumented the step function in Snap! with a trace logging tool, so that with each step, it adds a row in an *execution trace table* with the properties listed above, and dumps the trace table at the end of the execution. Figure 2 gives an example of a part of the *execution trace table*, in which one row logs one discrete step created by the step function, with each entry maps to a property described above (i.e., a concrete program state).

³paddle.x denotes the x coordinate of the paddle sprite. ETF uses these properties to summarize trace and generate features (Section 3.3). To allow comparison across students, sprites need to have consistent labels across student programs.

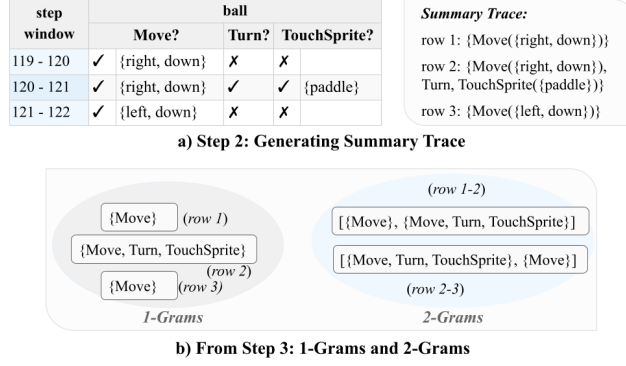


Fig. 3. Step 2&3: summarizing traces & generating features.

3.2 Step 2: Summarizing Traces

Some key features may not be directly captured by execution traces. For example, in Pong, the exact x-coordinate of a ball is less important than what direction it is moves towards. In this step, ETF infers a set of summarized properties in an execution trace, and connects them sequentially to generate *summary trace*. To do this, we use the ETF approach to scan through the execution trace table in a 2-step sliding window, apply a *Trace Abstraction Function (TAF)* that takes in a sliding window as input, and map this input to a set of summarized properties, which is next logged to the *summary trace*, described below.

The TAF Function. TAF maps a trace sliding window to a set of properties in two different ways. First, it infers changes between the first and next step in the 2-step sliding window, this generates properties such as movement and variable changes. For example, in the step window 120-121 of Figure 3, TAF identified a change in position (Move), a change in direction (Turn). Second, in addition to these inferred properties, some properties such as TouchSprite are useful properties that are already recorded in the execution trace. For these properties, TAF include them into the summary trace when it appears at the last step of the sliding window. For example, the step window 120-121 of Figure 3 included a property called TouchSprite, as the TouchSprite property was non-empty at Step 121 in Figure 2.

Parameters. Properties also has parameters, to describe informations such as left (\leftarrow)/right (\rightarrow) direction. For example, in Figure 3, TAF finds a Move(\rightarrow, \downarrow) property from the sliding window of Steps 119-120, as the ball is moving towards the lower-right direction.

By default, ETF uses 9 types of properties. Except 2 program state properties: (KeyDown and ChangeVar), the rest 7 are sprite-specific properties and are labeled with the sprite names (e.g., Paddle.Move). Among the 9 properties, 4 were original trace properties, and were directly returned when the corresponding property in the last step of the sliding window is non-empty: KeyDown, TouchSprite, TouchEdge, and OffStage, derived from fields in execution trace table, using the same parameters with the corresponding execution trace field at the last step of the sliding window. Other than the original trace *properties*, we list below the 5 inferred properties summarized by the TAF function:

- Move($\leftarrow, \rightarrow, \uparrow, \downarrow$) Returned when a sprite position changes towards a certain direction.
- Turn. Returned when a sprite changes direction.
- ChangeSize(+, -) Returned when a sprite changes its size to bigger (+) or smaller (-).
- ChangeVar(variable name, +, -). Returned when a global variable changes value to bigger (+) or smaller (-). To enable comparison across student programs, the variable names are standardized (e.g., var1 and var2).

- `ChangeLocalVar(variable name, <+, ->)`. Returned when a sprite-specific variable changes value to bigger (+) or smaller (-).

We designed the above candidate properties to be generalizable, and expect them to be useful across assignments. But these properties could be extended to additional functions very easily (e.g., to include user mouse inputs). A researcher only needs to define a new property and how it can be inferred in ETF.

3.3 Step 3: Generating n -Gram-based Features.

ETF next transforms the summary traces created in Step 2 into a set of features using n -Gram-based approach (discussed in Section 2), where an n -Gram takes a contiguous sequence of n items in data. ETF extracts features of 4 types:

- **1-Grams.** Extracting **simultaneous** behaviors, taken from each row of the summary trace. For example, the 1-Gram {Move, Turn, TouchSprite} in Figure 3(b) means “*Touching the paddle while turning to a different direction simultaneously*”
- **2-Grams.** Connecting adjacent 1-Grams sequentially. For example, in Figure 3, a 2-Gram from rows 1-2 in the property set sequence is [{Move}, {Move, Turn, TouchSprite}], denoting two consecutive timestamps, the first having only the Move property, but not turning or touching any other sprite, and the next includes touching the paddle and turning to a different direction.
- **Power Sets.** Many program behaviors (e.g., ball bounce when touching paddle) will only be concerned with 1 or 2 properties, e.g., Move and Turn. In these cases the other properties (e.g. a KeyDown) are irrelevant, and would add noise that makes it difficult for a classifier to see patterns across students’ traces. To address this, we extract n -Grams of not only the full property set in each row of the *summary trace*, but also of subsets of the property set, such as the 2-set of just Move and Turn. When constructing power sets for 2-Grams, we apply the power set on the types of properties that are possible in this 2-Gram. For example, the 2-Gram generated by rows 2-3 in the property sequence [{Move}, {Move, Turn, TouchSprite}] may be divided into 6 subsets as its power sets: [{Move}, {Move}], [{}, {Turn}], [{}, {TouchSprite}], [{Move}, {Move, Turn}], [{}, {Turn, TouchSprite}], [{}, {Move, TouchSprite}], and the 2-Gram itself.
- **Parameterized.** Parameters from a summary trace does not always add useful information to the extracted n -Grams. For example, when capturing a random movement behavior, it may be unnecessary to record the direction of the movement; but when we capturing an up arrow key pressing event, the parameter of KeyDown provides important information. Therefore, for all extracted n -Grams, we keep both *non-parameterized n -Grams*, where we do not record the parameters, as well as *parameterized n -Grams*, where each property would include its parameters when they were logged in the summary trace.

This step generates one feature set for each program snapshot have, including distinct features from all four feature types above. ETF next collects a *full feature set*, which consists of distinct features from all students. To retain all possibilities of features, features are considered distinct from one another if they have different properties, or different parameters of the same property. For example, a *parameterized n -Gram* is distinct from its *non-parameterized* counterpart.

3.4 Step 4: Filtering Features

Merging duplicate features and removing rare features. Based on the *full feature set* generated from Step 3, if features have the exact same distribution among programs, the ETF algorithm then merges these features as one

Table 1. Prevalence of correct programs on each rubric item.

rubric item	key_up	key_down	upper_bound	lower_bound	space_start
prevalence	0.76	0.73	0.55	0.56	0.50

rubric item	edge_bounce	paddle_bounce	paddle_score	reset_score	reset_ball
prevalence	0.49	0.42	0.35	0.23	0.31

feature; and it calculates the support of each feature based on the proportion of student programs that include this feature, and remove features that have support lower than a certain threshold, determined by a hyperparameter tuning process, described in Section 4.1.2.

Generating \vec{x} vectors. After generating, merging duplicates, and removing rare features, we use the resulting feature set as the independent variables, and for each student program, we use 1 as representing the presence of a feature in the student program (i.e. the n -Gram appeared at least once in their abstracted execution trace), 0 as the absence of the feature, and generate 0-1 digitized \vec{x} vector for each student’s code snapshot, which is used as vector input for a classification model.

4 EVALUATION

We aim to evaluate how ETF achieves its goal of automatically assessing students’ code on specific, fine-grained rubric items. In future work, this will enable offering automated, formative feedback. Our evaluation of ETF focused on the following research questions:

- **RQ1:** How accurately does ETF perform rubric-based code classification of students in-progress and submitted code, and how does this compare to syntax-based approaches?
- **RQ2:** How does the amount of labeled data available affect the performance of ETF-based classification?

We evaluated ETF on a dataset collected from an authentic classroom assignment, described below. We next investigate RQ1 by comparing ETF features and features extracted in a syntax-based approach *across models*, and *across rubrics*, described in Section 4.1. Then, we investigate RQ2 by understanding *how much training data* ETF and syntax-based approaches require to achieve their classification performance, through a cold start analysis, described in Section 4.2.

Dataset. We evaluated ETF on 42 students’ 162 code snapshots for a Pong game assignment, sampling student code snapshots at 10 minutes (42), 20 minutes (40) and 30 minutes (38) of work, as well as their final submissions (42). This creates a representative sample of when students might ask for help—both during programming, and before submitting their code. To create rubric-based labels on all in-progress and submitted code snapshots, one researcher manually graded all 162 programs based on the 10 rubric-based assignment requirements specified by the course instructor: 1) & 2) *key_up / down*: Paddle moves up / down with the up / down arrow key. 3) & 4) *upper / lower_bound*: When touching the upper / lower bound, the paddle does not move upwards / downwards even when the up / down key is pressed. 5) *space_start*: The ball starts movement when the space key is pressed. 6) *edge_bounce*: The ball bounces when touching the stage edge, unless touching the wall. 7) *paddle_bounce*: The ball bounces when touching the paddle. 8) *paddle_score*: If the ball touches the paddle, points increase by 1. 9) *reset_score*: If the ball touches the wall behind the paddle, the points are reset to 0. 10) *reset_ball*: If the ball touches the wall behind the paddle, the ball is reset.

We selected this dataset for three reasons: 1) Our prior work conducted formative analysis on the 42 submitted snapshots [40], and found it to represent the typical properties that are shared by many visual, interactive programs:

dynamic user interaction; object-specific program states; and properties that change over time. 2) It is a widely-used introductory programming assignment [20, 34, 46], suitable for many types of programming environments, and is therefore representative of many learning contexts and environments.

4.1 Classification Accuracy (RQ1)

We investigate RQ1 by a) comparing performance of ETF features and syntax-based features **across models**; b) comparing ETF features and syntax-based features **across rubrics** on a fixed model; and c) investigating the features produced by ETF algorithm and syntax-based approaches to understand how they caused differences in performance, described in Section 4.1.4. Our analysis of a) and b) follows the same procedure, where we started by generating ETF and syntax-based features separately (Section 4.1.1). We next performed the same feature filtering, training and evaluation procedure on the features we created (Section 4.1.2).

4.1.1 Generating Features. Generating ETF features. We used the procedure described in the Step 1 – 3 (Section 3.1–3.3) of the ETF approach to generate ETF features. First, to automatically execute student programs, we started by defining 5 *input generation rules* (explained in Section 3.1), corresponding to the requirements of rubric items: 1) *space key once*: when game starts, press space key once; corresponding to the *space_start* rubric item; 2) & 3) *intermittent up / down key*: press up / down arrow key for 2 steps in every 100 milliseconds, for a total of 15 seconds; corresponding to *key_up* and *key_down* rubric items; 4) & 5) *follow / evade ball*: when `ball.y > paddle.y`, press up key (to follow) / down key (to evade); when `ball.y < paddle.y`, press down key (to follow) / up key (to evade); *follow ball* were defined to ensure the paddle catches the ball, for the *paddle_bounce*, *paddle_score* rubric items; and *evade ball* were defined for the ball to hit back wall, for the *reset_score*, *reset_ball* rubric items. For each program snapshot, we re-executed it 5 times, each time performing the following sequence of inputs: 1,2 (meaning *space key once* followed with *intermittent up key*); 1,3; 1,4; 1,4,5 and 1,4,5. Each run of student programs generated one execution table. Overall, this resulted in 5×162 *execution trace tables*, each with an average of 1182 rows. After generating execution traces, we applied Steps 2 and 3 of ETF Algorithm (Section 3.2 and 3.3) to summarize each execution trace based on 2-step-based sliding windows, creating a *summary trace* with an average of 1181 rows per summary trace. Next, we created 1-Gram and 2-Gram-based *full feature sets* for each program. This step collected from each student an average of 518 distinct features, and created a total of 11148 distinct features in the full feature set.

Generating AST n -Gram and 1-Gram Features To examine how well ETF performs, we first compared it with a representative, syntax-based feature extraction approach that has performed well in prior evaluations by using the AST n -Gram feature extraction approach described by Akram et al. in [2]. Similar to Akram et al.’s work, we extracted all n -Grams from all ASTs, using $n = 1$ to 5 for *vertical n -Grams*, and $n = 1$ to 4 for *horizontal n -Grams* (explained in Section 2). Similar to many AST feature extraction approaches [29, 35, 47], we used a single label for all literals (`literal`). This step collected from each student an average of 153 distinct features, and created a total of 1145 distinct features as the full feature set. In addition to the structured n -Gram based approach, we also included a naive baseline, where we only generate the full feature set by extracting 1-Grams from all students’ ASTs (n -Grams with $n = 1$). This is essentially a 1-hot encoding of the AST nodes, where each node is either present (1) or absent (0). Each student’s program had an average of 20 distinct 1-Grams, creating a full feature set with 57 distinct 1-Grams.

4.1.2 Feature Filtering & Evaluation. We next applied the same feature filtering and evaluation approach to the ETF, AST n -Gram, and AST 1-Gram features, described below:

Table 2. F1 scores of AST 1-Grams, AST n-Grams, and ETF Features, over different models.

	AST 1-Grams	AST n-Grams	ETF Features
Logistic Regression	0.771	0.779	0.932
AdaBoost	0.78	0.78	0.922
Random Forest	0.763	0.773	0.926
MLP	0.764	0.771	0.908
Gaussian Process	0.739	0.728	0.923
SVM	0.759	0.771	0.93

Feature Filtering. For fairness of comparison, after collecting features, we used the Step 4 (Section 3.4) from the ETF algorithm to filter features for all ETF, AST n -Gram, and AST 1-Gram features. This steps reduced the number of features from 11148, 1145, and 57 to 2876, 769, and 36 for ETF, AST n -Grams, and AST 1-Grams respectively. For each type of the three features, we started by using ETF to automatically merge duplicate features (Step 4.a), and remove features that have support smaller than a certain threshold in the training set (Step 4.b). The threshold is set as a hyperparameter (tuned as described below.)

Classification Models. We used a variety of models, including those used in prior work (e.g. [21, 35, 45]), to ensure that our comparison of feature extraction approaches was not model-specific, and to investigate whether the best-performing approach was consistent across models. We used 6 models on the feature set, to evaluate the performance of the features: Logistic Regression, AdaBoost, Random Forest, Multi-layer perceptron (MLP), Gaussian Process, and SVM. Among them, the Gaussian Process model with an RBF kernel was also employed by Akram et al., and has shown to be the best performing model in the rubric-based performance inference task that they have applied [2].

Training & Evaluation. We employed 10-fold cross-validation to evaluate how accurately these different features predict the rubric-based performance. Within each round of cross-validation, we used another 2-fold cross-validation to tune the hyperparameters (i.e. nested cross-validation [39]). For all models, we included a minimum feature support threshold hyperparameter, T , below which we exclude the feature (e.g. ETF feature or AST n -Gram feature) from the final feature set, with the minimum support threshold as a hyperparameter, tuned based on 5 values: {0, 5%, 10%, 15%, 20%}. Additionally, since different feature extraction approaches may perform best with different values of model-specific hyperparameters, we also tuned hyperparameters for each classification models, based on the following values: *Logistic regression*: with penalty in {L1, L2}; *Random Forest*: with $n_estimators$ (i.e., number of trees in the forest) in {100, 200, 300, 400, 500}; *AdaBoost*: with learning rate in {0.01, 0.1, 1}; *MLP*: with learning rate in {0.001, 0.01, 0.1}; *SVM*: We used a linear kernel, with the regularization parameter (C) in {0.01, 0.1, 1, 10, 100}; Gaussian Process models optimize kernel hyperparameters during model fitting, we therefore did not tune hyperparameters for the Gaussian Process classifier. The values of the minimum feature support threshold and the model-specific hyperparameters were determined by their F_1 scores in the nested 2-fold cross-validation, based on a grid search on $5 \times \#(\text{model-specific hyperparameter values})$ possible types of hyperparameter combinations, during each round of the 10-fold cross-validation. Since many of our target labels are imbalanced, the accuracy score offers less information on how well our model performs in predicting target labels. We therefore use F_1 scores to tune hyperparameters. To ensure that data from a given student was not contained in *both* training and testing sets, all cross-validation splits were done on the 42 **students** (instead of on the 162 snapshots).

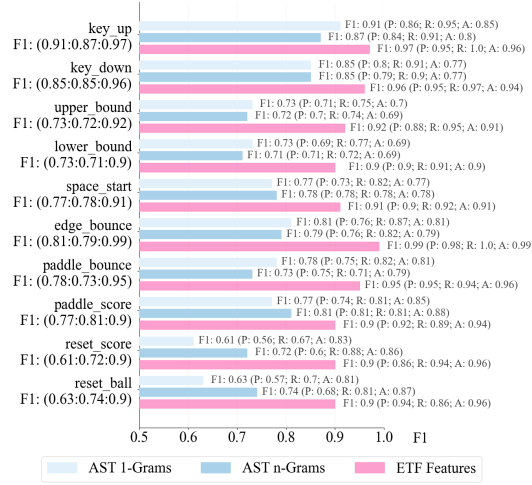


Fig. 4. The F_1 score (F1), Precision (P), Recall (R) and Accuracy (A) of ETF, n -Gram, and 1 -Gram features on each rubric, using an SVM model. x-axis starts from $F_1 = 0.5$.

4.1.3 Results. Comparison Across Models. For each of the 6 models described earlier, we used the model to predict students' rubric-based performance, and calculated its F_1 score among the 162 students' data using 10-fold cross validation, creating one F_1 score for each rubric. We next averaged F_1 scores for each classifier across all rubrics. Table 2 reports the average F_1 of each classifier in its prediction across all rubric items. First, an overall 0.73+ F_1 score of the AST 1 -Gram features shows that each node in students' code AST offered relatively useful information to indicate whether the student have successfully completed a certain rubric. Second, we saw that the AST n -Gram approach performed relatively higher than the 1 -Gram approach (5 in 6 cases), but only by 0.002–0.02. In addition, we also saw that ETF features generated F_1 scores that were consistently 0.14 to 0.2 higher than the AST n -Gram features, showing that all classifiers benefited from the execution-trace-based information extracted by the ETF features, with overall F_1 scores between 0.9 and 0.93. This result shows potential for us to make use of ETF features to correctly analyze students' current progress, which should enable automated, formative feedback in future work, to help a student who is stuck in the middle of programming.

Performance Across Rubrics. We next investigate the performance of the three feature extraction approaches across rubrics. Since all model show similar trends, we select SVM and present performance on all rubrics in Figure 4. We found 1) The naive AST 1 -Gram features had relatively lower F_1 scores on rubrics that had less prevalence in data (e.g., *reset_score*, *reset_ball*); 2) Comparing to AST 1 -Gram, the AST n -Gram features produced higher F_1 scores for *paddle_score*, *reset_score*, *reset_ball*, showing that AST n -Gram extracted more useful feature for these three rubric items. However, the ETF features performed relatively well across all rubrics, with its F_1 scores ranging from 0.9 to 0.99, showing that ETF features have strong potential to enable formative feedback on a variety of fine-grained, specific rubrics.

4.1.4 Case Study. Why does ETF perform better than the features extracted by the syntax-based n -Gram approach? We explore this question in two dimensions: 1) What are the possible ETF / AST n -Gram features that were useful / misleading? 2) What are possible reasons that some programs were misclassified? To answer (1), We looked for features

that are strongly associated with a given rubric score using a support-based feature selection method from prior work⁴ [16]. To do that, for each feature in the *full feature set*, we calculated its support (i.e., frequency of appearance) in the positive instances as S_{pos} , and its support in the negative instances as S_{neg} . We used the *absolute difference of support* between the two ($|Diff| = |S_{neg} - S_{pos}|$) to sort features based on their absolute difference, and inspect features that have the highest $|Diff|$ s. Next, to answer (2), for each one of the ETF / AST n -Gram features, we inspect all programs that were wrongly classified to identify possible reasons of misclassification. We used this approach to inspect all rubric items, and present our findings on the *upper_bound* rubric as an example.

Upper_bound requires that the paddle stops moving upwards when it hits the top of the stage (which implicitly requires completion of *key_up*). Our analysis of the top 10 ETF features with the highest $|Diff|$ for *upper_bound* identified three representative features⁵:

1) $\{Paddle.Move(\uparrow), KeysDown(\uparrow)\}$; with $S_{pos} = 1$, $S_{neg} = 0.51$. This feature describes scenarios when the paddle moves up and the up key is pressed simultaneously. All programs that satisfy *key_up* would include this feature in their execution trace, which is a prerequisite for *upper_bound*⁶. For simplicity, we refer to this feature as $F(key_move)$ below.

2) $\{Paddle.Move(\uparrow), KeysDown(\uparrow), \{KeysDown(\uparrow)\}\}$; with $S_{pos} = 0.58$, $S_{neg} = 0.01$. This feature describes scenarios in the trace where in the first timestamp, the paddle moves upward, with an up arrow key; and in the next timestamp, the paddle stops moving, despite the up arrow key still being pressed. This feature may appear if a player keeps pressing the up arrow key, but the paddle stops when it reaches the top of the stage. We refer to this feature as $F(key_stop)$.

3) $\{Paddle.OffStageemptyempty(empty), Paddle.TouchingEdgeemptyempty(empty)\}, \{Paddle.OffStageemptyempty(empty), Paddle.Moveemptyempty(empty), KeysDownemptyempty(empty), Paddle.TouchingEdgeemptyempty(empty)\}$; with $S_{pos} = 0$, $S_{neg} = 0.44$. This feature only appeared in the negative samples, where the paddle moves while OffStage when a key is pressed. This might occur if the player did not implement a proper bound, allowing the paddle to move despite being offstage. Note this feature is unparameterized, so it does not specify *which* key is pressed. We refer to this feature as $F(off_stage)emptyempty(empty)$.

These three features, if logged correctly, would capture almost all possible scenarios that could happen when a piece of code satisfied or violated the *upper_bound* rubric. A correct implementation *should* only have \vec{x} vector (described in Section 3.4) $[1, 1, 0]$, which means that $F(key_move)$ and $F(key_stop)$ are present (paddle moves up and stops at the edge) and $F(off_stage)emptyempty(empty)$ is absent (paddle does not move when offstage) “OffStage”. On the other hand, an incorrect implementation can have two possible types of \vec{x} vector: 1) $[0, 0, 0]$, when the paddle does not move with key at all; or 2) $[1, 0, 1]$, caused by the paddle moving up with the arrow key, but not stopping at the edge (causing an absent $F(key_stop)$ and a present $F(off_stage)emptyempty(empty)$). Our manual grading on the rubrics and their prevalence (Table 1) also shows that S_{neg} of $F(key_move)$ and $F(off_stage)emptyempty(empty)$ only differed slightly from students’ actual implementation: among the 75 programs that did not satisfy *upper_bound* rubric, 40 (53%) were caused by failures of implementing *key_up*; 35 (47%) caused by moving out of the stage upper bound; which only differed by 0.03 with the S_{neg} we found in $F(off_stage)emptyempty(empty)$.

⁴Because ETF and AST n -Gram both include a large number of correlated features, we cannot infer the most useful features from feature importance scores [11].

⁵Others were highly correlated with one of the three.

⁶Note that having this feature does not mean students implemented *key_up* correctly— Move and KeysDown may co-occur without causality, as explained in Section 3.3.

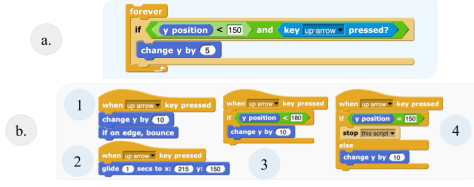


Fig. 5. One common correct implementation (a) and 4 types of less common correct implementations (b) for *upper_bound*.

However, the ETF features still caused 4 false negative and 11 false positive predictions. If the patterns have accurately defined what should happen in the trace, how did the wrong predictions occur? We analyzed the falsely-predicted programs and their feature sets, and found two reasons:

1) **Intermittent up key input caused $F(\text{key_stop})\text{emptyempty}(\text{empty})$ to not always occur in correct programs.** We would expect $F(\text{key_stop})\text{emptyempty}(\text{empty})$ to show up in *all* correct programs, but actually S_{pos} is only 0.58. As described in Section 4.1.1, we defined the *intermittent up key* input generation rule to programmatically press an up arrow key for 2 continuous steps every 100 milliseconds (a very short keypress). Consequently, each of these 2-step-long key-press instances may be interleaved by rows of other *property sets* in the *summary trace* (explained in Section 3.2). Thus, up key may not have been pressed long enough to capture a consecutive movement event (before touching edge) and stopping event (touching edge), and the $\{\uparrow\}\text{Paddle.Move}, \{\uparrow\}\text{KeysDown}\}$ property would not immediately follow $\{\uparrow\}\text{Paddle.Move}, \{\uparrow\}\text{KeysDown}\}$. As 2-Grams are only extracted for successive properties it may be missed even for correct programs. Furthermore, as $F(\text{key_stop})\text{emptyempty}(\text{empty})$ only has an S_{pos} of 0.58, other less relevant features (e.g., $F(\text{key_move})\text{emptyempty}(\text{empty})$) weighted higher, causing both *false negative* and *false positive* predictions. We therefore hypothesize that by increasing the time of pressing the key, we would increase S_{pos} of the $F(\text{key_stop})\text{emptyempty}(\text{empty})$ feature and improve the classification accuracy.

2) **Slow paddle movement caused the paddle to never achieve upper bound.** We discussed earlier that 47% of the students actually made the paddle go offstage, so we would expect the S_{neg} of $F(\text{off-stage})\text{emptyempty}(\text{empty})$ to be 0.47, not 0.44. We found that the two programs that had an un-detected $F(\text{off-stage})\text{emptyempty}(\text{empty})$ moved the paddle too slowly, such that the upper bound was not reached within our 15-second input sequence. Because $F(\text{off-stage})\text{emptyempty}(\text{empty})$ is an indicator of *not* fulfilling *upper_bound*, missing $F(\text{off-stage})\text{emptyempty}(\text{empty})$ also caused *false positive* predictions. This is an example where *one* single set of *input generation rules* does not always generalize to *all* possible students' code, since 15-second input can be too *long* for faster paddles that reach top in 1 seconds; and too *short* slower paddles (e.g., the two programs in this example) to enable the detection of *upper_bound*.

We next inspect AST n -Gram features. We found none of the 10 features with the highest $|Diff|$ s to connect strongly to what may be useful to predict *upper_bound*. The most relevant feature we found is one that signifies a $<$ comparison operator inside an if condition, with $S_{pos} = 0.69$, $S_{neg} = 0.27$. This feature would be present in one of the common correct solutions we identified from students' code (Figure 5 (a)) for determining if `paddle.y` is out of bound, but it does not appear in the 4 less common, but also correct solutions, shown in Figure 5 (b). Our analysis of misclassified programs showed a similar trend: Among the 23 false negatives predictions, 13 were implemented by **less common, but correct code** (Figure 5 (b)). This may have caused these false-negative predictions since their **code** does not have the most common AST n -Gram features; but these programs' **execution traces** led to the same distinctive features as the common programs, causing them to be predicted accurately by the ETF features. On the other hand, among the 28 false positive predictions, 19 used code that was only **one or two blocks different** from a correct solution,

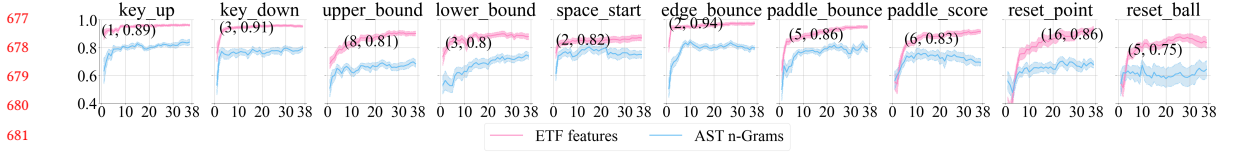


Fig. 6. Cold start curves for the ETF features (pink) and n -Gram features (blue), x-axis is the training data size on the number of students, y-axis is the average F_1 score by the model on the 10% hold-out test set (**Note**: the y-axis starts at 0.4).

e.g., reversing or incorrectly using the $<$ comparison inside the `if` condition, shown in the common correct solution (Figure 5 (a)). This caused the syntax-based approach to be unable to detect the small difference between correct and incorrect code, while ETF was able to extract the difference in output regardless.

4.2 Cold Start Analysis (RQ2)

One potential limitation of data-driven code classification approaches is that they require a lot of data labeling that the instructor may not otherwise do [7, 42]. Therefore, our RQ2 asks how the amount of labeled data available affects the performance of ETF-based classification. We adapted a cold start procedure outlined by Price et al. to measure the relationship between data quantity and the resulting performance of the features [30]. We used 20 random trials, each started by shuffling the 42 students in the dataset. With each trial, we used the last 10% of students (i.e., 4 students with 8–16 snapshots) as the testing set, and started adding the remaining students (38) one by one into the training data set. Note that both during shuffling and increasing training data, we use one **student** as one unit of analysis, instead of one **snapshot**, to ensure that a given student’s program is not contained in both the training and testing set. At each iteration with a training set of size $|T|$, this represents a scenario in which an instructor or researcher chooses $|T|$ random students to label, and we evaluate how well the both the ETF features and the n -Gram features would perform with this amount of data. Inside the training set, we used an SVM model with a linear kernel, and used the same approach described in Section 4.1 to grid-search and tune hyperparameters when applicable⁷.

Figure 6 shows the average F_1 score of the SVM model prediction across 20 trials, at the increase of training data size, with shades inside $+$ and $-$ standard error. We found that the ETF features consistently performed better than the AST n -Gram features, starting to show higher prediction F_1 scores at less than 5 students’ code snapshots ($n < 20$). We also noticed the more unstable, higher-error cold-start curves in both ETF and AST n -Gram features for the last rubric item `reset_ball`, which may have been caused by both the randomness in the shuffling process and the lack of predictive, useful features. In addition, we marked on Figure 6 the first time ETF features achieves higher than 90% of its final F_1 score, and found it required a range of 1 – 54 snapshots from 1 to 16 students. 6 rubric items took less than 4 students for training, but others (e.g., `upper_bound`, `reset_score`) needed more training data, which may have been caused by misleading features we discussed in Section 4.1.4. This shows that ETF features are able to produce high predictive accuracy without requiring a large amount of training data, and that an instructor only needs to grade at most half of the students’ snapshots to enable relatively accurate (with an $F_1 > 90\%$ of its final F_1) analysis of student progress.

⁷In a small training sample, there may not be enough data for train-validation separation. For those data, we used the default hyperparameters (i.e., regularization parameter $C = 1$, no minimum support threshold).

5 DISCUSSION & LIMITATIONS

We found that ETF creates features that lead to highly accurate rubric-based classification of students' in-progress and submitted code (average $F_1 = 0.93$), and it only requires relatively little labeled data to achieve high performance on most rubrics. Therefore, this work shows potential to use the ETF approach to generate formative feedback to students, indicating their progress on each rubric item as they work. This would only require an instructor to manually grade a small amount of students' in-progress and submitted code to create the model. We envision that the ETF approach can be used to build a tool where, whenever a student runs their code, it logs an execution trace, extracts a set of meaningful features from that execution trace, and use the presence and absence of features to accurately predict students' rubric-based performance. This could be used to offer encouraging feedback as students progress that can improve performance and affective outcomes (e.g. [8, 23]). It could also monitor students' progress over time, and prompt them to seek help when they have not progressed for a while (e.g. [22, 33]), or to inform a teacher dashboard summarizing students' progress (e.g. [43]). While our model is not perfect, its accuracy is high enough to ensure the overall quality of this feedback, and comparable with expert-authored formative assessment models, used in existing systems [23, 41], which are also fallible. However, despite the strong potential of applying execution trace-based feature engineering to visual, interactive programs, our work does include several limitations, discussed below:

Limitation: ETF features are sensitive to inputs. All trace-based program analysis takes in inputs to execute the programs. However, defining inputs for visual, interactive programs, as we have shown in Section 4.1.4, is particularly difficult. We have shown that defining one set of input generation rules may not always produce high-coverage executions for all programs. This caused ETF to fail to generate features that are indicative of a certain behavior, as shown in Section 4.1.4. Similar to the case study in Section 4.1.4, most of the erroneous predictions produced by the ETF features were caused by insufficient input sequences, showing that the input sequence had a strong influence on the classification results. This difficulty about producing high-coverage input sequences is shared by all program analysis tools for visual, interactive programs [12, 37, 38], and requires future work such as fully automated input generation to handle this problem. However, we also imagine that the real-world scenario that we envision ETF to work best in does not always require writing *input generation rules* for automated input generation. We may simply log execution traces when a student or instructor manually executes a program (e.g., to test or grade it), and use these execution traces to train the model (from instructor-generated traces) and to offer feedback (for student-generated traces). These execution traces would almost certainly include meaningful inputs that students / instructors use to test and grade programs, which may be even more effective than our input generation rules.

Limitation: Selection of candidate properties require human expertise. For the Pong assignment, we designed ETF to include 9 candidate properties, such as detecting movement and variable changes in the execution trace. While we expect the candidate properties we included to be generalizable, as they are found in many other novice programming tasks, we did not evaluate their generalizability here. In addition, there could potentially be other properties that are important for different assignments (e.g., sound effects, mouse movement, etc), which would require adding new candidate properties. We make no claim that the 9 properties we designed for Pong are an exhaustive list of useful properties, but suggest them as a core set that should be generally useful. In future work, we will explore expanding this set of candidate properties to cover additional assignments and to explore the generalizability of the ETF approach.

6 CONCLUSION

In this work, we explore feature engineering for classifying visual, interactive programs. We present a novel, effective approach that extracted useful features from execution traces (ETF), leading to high predictive accuracy. Our results show strong potential for using ETF to monitor student progress and offer automated, formative feedback.

REFERENCES

- [1] B. Akram et al. Assessment of students' computer science focal knowledge, skills, and abilities in game-based learning environments. 2019.
- [2] B. Akram, W. Min, E. Wiebe, A. Navied, B. Mott, K. E. Boyer, J. Lester, et al. Automated assessment of computer science competencies from student programs with gaussian process regression. In *Proceedings of the 2020 Educational Data Mining Conference*, 2020.
- [3] M. Allamanis, E. T. Barr, P. Devanbu, and C. Sutton. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)*, 51(4):1–37, 2018.
- [4] U. Alon, M. Zilberstein, O. Levy, and E. Yahav. A general path-based representation for predicting program properties. *PLDI'18*, 2018.
- [5] U. Alon, M. Zilberstein, O. Levy, and E. Yahav. code2vec: Learning distributed representations of code. *POPL'19*, 2019.
- [6] C. R. Association et al. Generation cs: Computer science undergraduate enrollments surge since 2006. Retrieved March, 20:2017, 2017.
- [7] D. Azcona, P. Arora, I.-H. Hsiao, and A. Smeaton. user2code2vec: Embeddings for profiling students based on distributional representations of source code. In *LAK'19*, 2019.
- [8] A. T. Corbett and J. R. Anderson. Locus of feedback control in computer-based tutoring: Impact on learning rate, achievement and attitudes. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 245–252, 2001.
- [9] N. Diana, M. Eagle, J. Stamper, S. Grover, M. Bienkowski, and S. Basu. Data-driven generation of rubric criteria from an educational programming environment. In *Proceedings of the 8th International Conference on Learning Analytics and Knowledge*, pages 16–20, 2018.
- [10] S. H. Edwards and K. P. Murali. Codeworkout: short programming exercises with built-in data collection. In *Proceedings of the 2017 ACM Conference on Innovation and Technology in Computer Science Education*, pages 188–193, 2017.
- [11] D. E. Farrar and R. R. Glauber. Multicollinearity in regression analysis: the problem revisited. *The Review of Economic and Statistics*, pages 92–107, 1967.
- [12] B. for review. Blinded for review. Blinded for reiew.
- [13] S. Grover. Toward a framework for formative assessment of conceptual learning in k-12 computer science classrooms. 2021.
- [14] D. Hovemeyer and J. Spacco. Cloudcoder: a web-based programming exercise system. *Journal of Computing Sciences in Colleges*, 28(3):30–30, 2013.
- [15] S. Iyer, I. Konstas, A. Cheung, and L. Zettlemoyer. Summarizing source code using a neural attention model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2073–2083, 2016.
- [16] J. S. Kinnebrew, K. M. Loretz, and G. Biswas. A contextualized, differential sequence mining method to derive students' learning behavior patterns. *JEDM/ Journal of Educational Data Mining*, 5(1):190–219, 2013.
- [17] A. J. Ko, B. A. Myers, and H. H. Aung. Six learning barriers in end-user programming systems. In *2004 IEEE Symposium on Visual Languages-Human Centric Computing*, pages 199–206. IEEE, 2004.
- [18] M. Kölling. The greenfoot programming environment. *ACM Transactions on Computing Education (TOCE)*, 10(4):1–21, 2010.
- [19] A. N. Kumar. Explanation of step-by-step execution as feedback for problems on program analysis, and its generation in model-based problem-solving tutors. *Technology, Instruction, Cognition and Learning (TICL) Journal*, 4(1), 2006.
- [20] N. Lytle, A. Milliken, V. Cateté, and T. Barnes. Investigating different assignment designs to promote collaboration in block-based environments. In *Proceedings of the 51st ACM Technical Symposium on Computer Science Education, SIGCSE '20*, page 832–838, New York, NY, USA, 2020. Association for Computing Machinery.
- [21] Y. Mao, R. Zhi, F. Khoshnevisan, T. W. Price, T. Barnes, and M. Chi. One minute is enough: Early prediction of student success and event-level difficulty during a novice programming task. *International Educational Data Mining Society*, 2019.
- [22] S. Marwan, A. Dombe, and T. W. Price. Unproductive help-seeking in programming: What it is and how to address it. In *Proceedings of the 2020 ACM Conference on Innovation and Technology in Computer Science Education*, pages 54–60, 2020.
- [23] S. Marwan, G. Gao, S. Fisk, T. W. Price, and T. Barnes. Adaptive immediate feedback can improve novice programming engagement and intention to persist in computer science. In *Proceedings of the 2020 ACM Conference on International Computing Education Research*, pages 194–203, 2020.
- [24] J. Moenig and B. Harvey. Byob build your own blocks (a/k/a snap!). URL: <http://byob.berkeley.edu/>, accessed Aug, 2012.
- [25] L. Mou, G. Li, L. Zhang, T. Wang, and Z. Jin. Convolutional neural networks over tree structures for programming language processing. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [26] B. Paaßen, J. Jensen, and B. Hammer. Execution traces as a powerful data representation for intelligent tutoring systems for programming. *International Educational Data Mining Society*, 2016.
- [27] A. Patil. Automatic grading of programming assignments. 2010.
- [28] C. Piech, J. Huang, A. Nguyen, M. Phulsuksombati, M. Sahami, and L. Guibas. Learning program embeddings to propagate feedback on student code. In *International Conference on Machine Learning*, pages 1093–1102, 2015.

- [29] T. W. Price, Y. Dong, and T. Barnes. Generating data-driven hints for open-ended programming. *International Educational Data Mining Society*, 2016.
- [30] T. W. Price, R. Zhi, Y. Dong, N. Lytle, and T. Barnes. The impact of data quantity and source on the quality of data-driven hints for programming. In *International conference on artificial intelligence in education*, pages 476–490. Springer, 2018.
- [31] C. Reas and B. Fry. Processing: programming for the media arts. *AI & SOCIETY*, 20(4):526–538, 2006.
- [32] M. Resnick, J. Maloney, A. Monroy-Hernández, N. Rusk, E. Eastmond, K. Brennan, A. Millner, E. Rosenbaum, J. Silver, B. Silverman, et al. Scratch: programming for all. *Communications of the ACM*, 52(11):60–67, 2009.
- [33] I. Roll, V. Aleven, B. M. McLaren, E. Ryu, R. S. d Baker, and K. R. Koedinger. The help tutor: does metacognitive feedback improve students’ help-seeking actions, skills and learning? In *International Conference on Intelligent Tutoring Systems*, pages 360–369. Springer, 2006.
- [34] K. Roy. App inventor for android: Report from a summer camp. In *Proceedings of the 43rd ACM Technical Symposium on Computer Science Education*, SIGCSE ’12, page 283–288, New York, NY, USA, 2012. Association for Computing Machinery.
- [35] Y. Shi, K. Shah, W. Wang, S. Marwan, P. Penmetsa, and T. Price. Toward semi-automatic misconception discovery using code embeddings. In *The 11th International Conference on Learning Analytics Knowledge (LAK 21)*, 2021.
- [36] V. J. Shute. Focus on Formative Feedback. *Review of Educational Research*, 78(1):153–189, 2008.
- [37] A. Stahlbauer, C. Frädriich, and G. Fraser. Verified from scratch: Program analysis for learners’ programs. In *35th IEEE/ACM International Conference on Automated Software Engineering, ASE 2020, Melbourne, Australia, September 21–25, 2020*, pages 150–162. IEEE, 2020.
- [38] A. Stahlbauer, M. Kreis, and G. Fraser. Testing scratch programs automatically. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 165–175, 2019.
- [39] M. Stone. Cross-validators choice and assessment of statistical predictions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(2):111–133, 1974.
- [40] W. Wang, C. Zhang, A. Stahlbauer, G. Fraser, and T. Price. Snapcheck: Automated testing for snap programs. In *Proceedings of the 26th ACM Conference on Innovation and Technology in Computer Science Education V. 1, ITiCSE ’21*, pages 227–233, New York, NY, USA, 2021. Association for Computing Machinery.
- [41] W. Wang, R. Zhi, A. Milliken, N. Lytle, and T. W. Price. Crescendo : Engaging Students to Self-Paced Programming Practices. In *Proceedings of the ACM Technical Symposium on Computer Science Education*, 2020.
- [42] M. Wu, M. Mosse, N. Goodman, and C. Piech. Zero shot learning for code education: Rubric sampling with deep learning inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 782–790, 2019.
- [43] F. Xhakaj, V. Aleven, and B. M. McLaren. Effects of a teacher dashboard for an intelligent tutoring system on teacher knowledge, lesson planning, lessons and student learning. In *European conference on technology enhanced learning*, pages 315–329. Springer, 2017.
- [44] L. Yan, N. McKeown, and C. Piech. The pyramidsnapshot challenge: Understanding student process from visual output of programs. In *Proceedings of the 50th ACM Technical Symposium on Computer Science Education*, pages 119–125, 2019.
- [45] T. P. T. B. Ye Mao, Samiha Marwan and M. Chi. What time is it? student modeling needs to know. pages 171–182. *Proceedings of The 13th International Conference on Educational Data Mining (EDM 2020)*, 07 2020.
- [46] R. Zhi, M. Chi, T. Barnes, and T. W. Price. Evaluating the effectiveness of parsons problems for block-based programming. In *Proceedings of the 2019 ACM Conference on International Computing Education Research*, pages 51–59. ACM, 2019.
- [47] R. Zhi, S. Marwan, Y. Dong, N. Lytle, T. W. Price, and T. Barnes. Toward Data-Driven Example Feedback for Novice Programming. *Proceedings of the International Conference on Educational Data Mining*, pages 218–227, 2019.
- [48] R. Zhi, T. W. Price, N. Lytle, Y. Dong, and T. Barnes. Reducing the state space of programming problems through data-driven feature detection. In *Educational Data Mining in Computer Science Education (CSEDM) Workshop@ EDM*, 2018.

- [47] R. Zhi, S. Marwan, Y. Dong, N. Lytle, T. W. Price, and T. Barnes. Toward Data-Driven Example Feedback for Novice Programming. *Proceedings of the International Conference on Educational Data Mining*, pages 218–227, 2019.
- [48] R. Zhi, T. W. Price, N. Lytle, Y. Dong, and T. Barnes. Reducing the state space of programming problems through data-driven feature detection. In *Educational Data Mining in Computer Science Education (CSEDM) Workshop@ EDM*, 2018.

- [45] L. Yan, N. McKeown, and C. Piech. The pyramidsnapshot challenge: Understanding student process from visual output of programs. In *Proceedings of the 50th ACM Technical Symposium on Computer Science Education*, pages 119–125, 2019.
- [46] T. P. T. B. Ye Mao, Samiha Marwan and M. Chi. What time is it? student modeling needs to know. pages 171–182. Proceedings of The 13th International Conference on Educational Data Mining (EDM 2020), 07 2020.
- [47] R. Zhi, M. Chi, T. Barnes, and T. W. Price. Evaluating the effectiveness of parsons problems for block-based programming. In *Proceedings of the 2019 ACM Conference on International Computing Education Research*, pages 51–59. ACM, 2019.
- [48] R. Zhi, S. Marwan, Y. Dong, N. Lytle, T. W. Price, and T. Barnes. Toward Data-Driven Example Feedback for Novice Programming. *Proceedings of the International Conference on Educational Data Mining*, pages 218–227, 2019.
- [49] R. Zhi, T. W. Price, N. Lytle, Y. Dong, and T. Barnes. Reducing the state space of programming problems through data-driven feature detection. In *Educational Data Mining in Computer Science Education (CSEDM) Workshop@ EDM*, 2018.