Program Equivalence for Assisted Grading of Functional Programs

JOSHUA CLUNE, Carnegie Mellon University, United States of America VIJAY RAMAMURTHY, Carnegie Mellon University, United States of America RUBEN MARTINS, Carnegie Mellon University, United States of America UMUT A. ACAR, Carnegie Mellon University, United States of America

In courses that involve programming assignments, giving meaningful feedback to students is an important challenge. Human beings can give useful feedback by manually grading the programs but this is a time-consuming, labor intensive, and usually boring process. Automatic graders can be fast and scale well but they usually provide poor feedback. Although there has been research on improving automatic graders, research on scaling and improving human grading is limited.

We propose to scale human grading by augmenting the manual grading process with an equivalence algorithm that can identify the equivalences between student submissions. This enables human graders to give targeted feedback for multiple student submissions at once. Our technique is conservative in two aspects. First, it identifies equivalence between submissions that are algorithmically similar, e.g., it cannot identify the equivalence between quicksort and mergesort. Second, it uses formal methods instead of clustering algorithms from the machine learning literature. This allows us to prove a soundness result that guarantees that submissions will never be clustered together in error. Despite only reporting equivalence when there is algorithmic similarity and the ability to formally prove equivalence, we show that our technique can significantly reduce grading time for thousands of programming submissions from an introductory functional programming course.

CCS Concepts: • Theory of computation \rightarrow Automated reasoning.

Additional Key Words and Phrases: Program Equivalence, Assisted Grading, Formal Methods, Functional Programming

ACM Reference Format:

Joshua Clune, Vijay Ramamurthy, Ruben Martins, and Umut A. Acar. 2020. Program Equivalence for Assisted Grading of Functional Programs. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 171 (November 2020), 29 pages. https://doi.org/10.1145/3428239

1 INTRODUCTION

There have been many efforts to develop techniques for automated reasoning of programming assignments at scale. This has lead to the rise of automatic graders, programs that take in a set of student submissions and output grades or feedback for those submissions without requiring any human input. While recent years have yielded substantial improvements in automatic grading techniques [Gulwani et al. 2018; Kaleeswaran et al. 2016; Liu et al. 2019; Perry et al. 2019; Singh

Authors' addresses: Joshua Clune, Carnegie Mellon University, United States of America, josh.seth.clune@gmail.com; Vijay Ramamurthy, Carnegie Mellon University, United States of America, vrama628@gmail.com; Ruben Martins, Carnegie Mellon University, United States of America, rubenm@andrew.cmu.edu; Umut A. Acar, Carnegie Mellon University, United States of America, umut@cs.cmu.edu.

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2020 Copyright held by the owner/author(s).

2475-1421/2020/11-ART171

https://doi.org/10.1145/3428239

et al. 2013; Wang et al. 2018], automatic graders are still more limited in the feedback they can provide than human graders.

This creates a trade-off between scale and quality. For small courses, it makes sense to utilize human graders in order to provide the best feedback possible. For Massive Open Online Courses, human involvement in grading all submissions is often logistically impossible, so it makes sense to use automatic graders. But neither option is ideal for large, in-person, introductory functional courses. When introductory functional courses use automatic graders, it hurts the students because they receive less targeted feedback, and it can hurt the teaching staff to lose a valuable avenue for addressing uncommon misunderstandings. But when introductory functional courses use human graders, it creates a large burden on the teaching staff, and it may require capping the size of the class, hurting students by limiting their opportunity to take the class.

To provide an option that eases the cost of human grading without sacrificing feedback quality, we propose a method of enabling human graders to give targeted feedback to multiple students at once. Our approach takes a pair of expressions submitted by students and deconstructs them simultaneously to build up a formula that is valid only if the expressions are equivalent. This pairwise equivalence test is used to cluster student submissions into buckets for which all submissions can be graded and given feedback simultaneously. Our approach recognizes expressions as equivalent by finding equivalences in each expression's subexpressions. To do this, it uses a variety of inference rules to simultaneously deconstruct the expressions down to their atomic subexpressions. It then outputs formulas that are valid only if the atomic subexpressions are equivalent. Finally, our inference rules recursively use the formulas of these subexpressions as subformulas to build up a larger formula that indicates the equivalence of the overall expression. This final formula's validity can be checked by an SMT Solver to determine whether the two expressions are equivalent.

A central benefit of our approach is that when two expressions are recognized as equivalent, this fact does not merely reflect that the two expressions produce the same outputs on shared inputs. In input/output grading, the correctness of code is determined entirely by whether a student submission produces correct outputs when given a large and diverse set of inputs. But in our approach, all equivalences arise from similarities in subexpressions, so equivalences found by our technique are discoverable only due to underlying algorithmic similarities. This enables instructors to give feedback based not only on whether a problem was solved correctly, but based on the algorithmic decisions that were involved in the student's solution.

Three primary factors that impact the grading and feedback of student programs are correctness, algorithmic approach, and style. While our approach is meant to enable providing better feedback concerning algorithmic approach, as opposed to simply providing feedback concerning correctness as in input/output grading, evaluating style is outside of the scope of our technique. For that reason, we believe that our approach is best utilized in conjunction with the methods courses already use to evaluate style. For courses already doing automatic grading, this should not be an issue because if they are already doing automatic grading, they are already automatically doing style checking, and can, therefore, use that in conjunction with our approach to provide all of the same style feedback the course already provided, but additionally provide human feedback for algorithmic content.

For courses already doing fully human grading, even if it is still necessary to grade each assignment individually to address style concerns, we believe our approach can make it possible to better allocate human resources for the grading process. A grader focusing entirely on one or two large buckets can be more efficient by not being forced to figure out which common approach is being taken by every individual submission. This can help the grader more quickly move on from understanding the student's solution to addressing any style concerns, and it also helps ensure fairer grading in guaranteeing that the same grader will grade all similar submissions. A grader focusing entirely on grading submissions that were clustered with few if any other programs can

anticipate ahead of time that their grading will likely require providing more frequent and/or detailed comments. This can enable course staffs to give more submissions to graders of large buckets, easing the burden of singleton/small bucket graders.

The differences between our approach and other state-of-the-art automatic graders and clustering techniques [Gulwani et al. 2018; Perry et al. 2019; Wang et al. 2018] stem from differences in motivation. Since each bucket generated by our approach is meant to be graded by a human, it is more important for our technique to distinguish nonequivalent submissions than to ensure that all equivalent submissions are placed in the same bucket. Ensuring that all equivalent submissions are placed in the same bucket reduces time spent grading equivalent programs, enabling instructors to spend more time giving detailed feedback. This is an important goal, but it is of lower priority than preserving the accuracy of human feedback because it does not matter how detailed feedback is if it does not apply to the student to whom it is given. To secure the accuracy of human feedback while using our approach, we guarantee the correctness of our technique's recognized equivalences by proving a soundness theorem that states that if our technique recognizes two expressions as equivalent, they necessarily exhibit identical behavior.

In summary, the contributions of our paper are as follows:

- We define an effective and efficient technique for identifying equivalences between purely functional programs. The technique's design ensures that only algorithmically similar programs will be recognized as equivalent.
- We prove the soundness of this technique, showing that if our approach identifies an equivalence between two expressions, then the two expressions must exhibit identical behavior.
- We implement our approach in a tool called ZEUs and demonstrate its effectiveness in assisting the grading of more than 4,000 student submissions from a functional programming course taught at the college level in Standard ML.

2 MOTIVATING EXAMPLES

Our approach is meant to cluster expressions that are algorithmically similar, but potentially syntactically different. In this section, we show two examples of similar implementations of the same function that are successfully identified by our tool as equivalent, and describe one example in which two solutions to a task are not recognized as equivalent due to algorithmic dissimilarities.

```
fun add_opt x y =
    case (x, y) of
    (SOME m, SOME n) =>
        SOME (m + n)
    | (NONE, _) => NONE
    | (_, NONE) => NONE
    | fun add_opt x y =
        bind x (fn m =>
        bind y (fn n =>
        return (m + n)
    ))
```

Fig. 1. Two implementations of adding two optional numbers

Figure 1 contains two functions that take in two int options as input, and adds the ints in the options if possible, returning NONE otherwise. The right expression's conditional logic is modeled after Haskell-style monads, interacting with the higher order bind function to case on x

first, and then potentially y depending on the value of x, whereas the left expression cases on x and y simultaneously. Still, our approach is able to fully encode both expressions' conditional logic structures and produce a valid formula. A demonstration of how our approach specifically encodes these conditional logic structures is included in Section 5.

```
fun split [] = ([], [])
fun split [] = ([], [])
  | split [x] = ([x], [])
                                       | split (x::xs) =
  | split (x::y::L) =
                                            case xs of
                                                 [] \Rightarrow ([x], [])
        val (A, B) = split L
                                              | (y::ys) =>
                                                   let
        (x::A, y::B)
                                                     val (A, B) = split ys
      end
                                                     (x::A, y::B)
fun merge([], L) = L
                                                   end
  | merge(L, []) = L
  \mid merge(x::xs, y::ys) =
                                     fun merge (11, 12) =
      if x < y
                                        case 11 of
                                          [] => 12
      then x :: merge (xs, y::ys)
                                        | x::xs =>
      else y :: merge (x::xs, ys)
                                          case 12 of
fun msort [] = []
                                            [] => 11
  | msort [x] = [x]
                                          | y::ys =>
  | msort L =
                                            if x < y
      let
                                            then x :: merge (xs, 12)
        val (A, B) = split L
                                            else y :: merge (11, ys)
        merge(msort A, msort B)
                                      fun msort [] = []
                                        | msort [x] = [x]
      end
                                        | msort L =
                                            let
                                               val(A, B) = split L
                                              merge(msort A, msort B)
```

Fig. 2. Two implementations of mergesort

Figure 2 contains two functions that implement mergesort. The left implementation uses a style that emphasizes pattern matching on input arguments while the right implementation uses a style that emphasizes nesting binding structures. Despite their syntactic differences, both functions implement the same underlying algorithm. Therefore, our approach recognizes them as equivalent.

Our approach is not intended to cluster programs just by correctness, or final input/output behavior, but by structure. This enables our approach to distinguish between correct submissions that use different algorithms. For instance, one of the benchmarks we use in Section 7 to evaluate our tool is a task called slowDoop. The goal of this task is to take in an arbitrary list L and return a list in which all elements in L appear exactly once. Consider a similar task in which the goal is the same but has the added stipulation that the final list must be sorted. A reasonable $O(n^2)$ solution to this task would be to iterate over L, only keeping elements that do not appear later in the list, and then sort the result. But a better $O(n\log n)$ solution would be to first sort L, and then

iterate over the resulting list once to remove duplicate elements. While correct implementations of these algorithms are identical from an input/output perspective, our approach would cluster them separately, and we believe that they merit different feedback.

3 LAMBDAPIX

Our approach operates over a language which we call LambdaPix. LambdaPix is designed to be a target for transpilation from functional programming languages such as Standard ML, OCaml, or Haskell. Our techniques apply to purely functional programs only and do not allow for state (e.g., references) but are otherwise unrestricted and make no further assumptions about the programs. In this section, we present the syntax and semantics for LambdaPix.

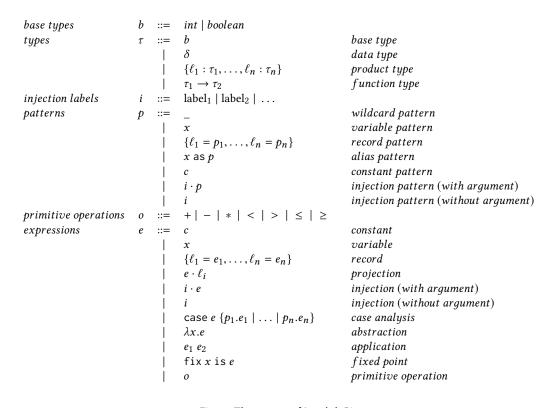


Fig. 3. The syntax of LambdaPix

We give the syntax for LambdaPix in Figure 3. Arbitrary labeled product types are supported as labeled records. For sum types and recursive types, LambdaPix is defined over an arbitrary fixed set of algebraic data types, with associated injection labels. We use meta-variables x, y, and z (and variants) to range over an unspecified set of variables.

3.1 Static Semantics

We assume an arbitrary fixed set of disjoint algebraic data types with unique associated injection labels (by unique, it is meant that there are no shared injection labels between distinct data types). In particular, we assume a fixed set of judgments of the form $i: \tau \hookrightarrow \delta$ for injection labels that take in an argument of type τ to produce an expression of data type δ , and a fixed set of judgments of

the form $i:\delta$ for injection labels of data type δ that do not take in an argument. We take $i:\tau \hookrightarrow \delta$ to mean that the type δ has a label i which accepts an argument of type τ , and we take $i:\delta$ to mean that the type δ has a label i that does not accept an argument. Note that by allowing τ to contain instances of δ , this data type system affords LamdbaPix a form of inductive types.

$$\frac{- : \tau + PATTY_{1}}{\left\{l_{1} = p_{1}, \dots, l_{n} = p_{n}\right\} :: \left\{\ell_{1} : \tau_{1}, \dots, \ell_{n} : \tau_{n}\right\} + \Gamma_{1} \dots \Gamma_{n}} PATTY_{3} \qquad \frac{p : \tau + \Gamma}{x \text{ as } p :: \tau + \Gamma, x : \tau} PATTY_{4}$$

$$\frac{c :: b + PATTY_{5}}{i :: b + i} PATTY_{5} \qquad \frac{i : \delta}{i :: \delta + i} PATTY_{6} \qquad \frac{i : \tau \hookrightarrow \delta \quad p :: \tau + \Gamma}{i \cdot p :: \delta + \Gamma} PATTY_{7}$$

Fig. 4. Pattern typing in LambdaPix

Figure 4 defines an auxiliary judgment used in the typechecking of case expressions. This pattern typing judgment $p :: \tau \dashv \Gamma$ defines that expressions of type τ can be matched against the pattern p, and that doing so produces new variable bindings whose types are captured in Γ .

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \dots \quad \Gamma \vdash e_n : \tau_n}{\Gamma \vdash e_1 : \tau_1 \quad \dots \quad \Gamma \vdash e_n : \tau_n} \quad \text{Ty}_3$$

$$\frac{\Gamma \vdash e : \{\dots, \ell_i : \tau_i, \dots\}}{\Gamma \vdash e \vdash \ell_i : \tau_i} \quad \text{Ty}_4$$

$$\frac{\frac{\Gamma \vdash e : \{\dots, \ell_i : \tau_i, \dots\}}{\Gamma \vdash e \vdash \ell_i : \tau_i} \quad \text{Ty}_5$$

$$\frac{\frac{i : \sigma}{\Gamma \vdash i : \sigma} \quad \Gamma \vdash e : \tau}{\Gamma \vdash i : e : \sigma} \quad \text{Ty}_6$$

$$\frac{\Gamma \vdash e : \tau \quad p_1 :: \tau \vdash \Gamma_1 \quad \Gamma, \Gamma_1 \vdash e_1 : \tau' \quad \dots \quad p_n :: \tau \vdash \Gamma_n \quad \Gamma, \Gamma_n \vdash e_n : \tau'}{\Gamma \vdash \text{case } e \mid \{p_1 . e_1 \mid \dots \mid p_n . e_n\} : \tau'} \quad \text{Ty}_7$$

$$\frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x . e : \tau_1 \to \tau_2} \quad \text{Ty}_8$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \to \tau_2 \quad \Gamma \vdash e_2 : \tau_1}{\Gamma \vdash e_1 e_2 : \tau_2} \quad \text{Ty}_9$$

$$\frac{\Gamma, x : \tau_1 \vdash e : \tau}{\Gamma \vdash \text{fix } x \text{ is } e : \tau} \quad \text{Ty}_{10}$$

$$\frac{\Gamma, \sigma : \tau_1 \to \sigma_2 \quad \Gamma \vdash \sigma_2 : \tau_1 \to \tau_2}{\Gamma \vdash \sigma : \tau_1 \to \tau_2 \vdash \sigma : \tau_1 \to \tau_2} \quad \text{Ty}_{11}$$

Fig. 5. Expression typing in LambdaPix

Figure 5 defines typing for expressions in LambdaPix.

Definition 3.1 (Well-formed). A LambdaPix expression e is well-formed if there exists a type τ such that $\Gamma_{\text{initial}} \vdash e : \tau$, where Γ_{initial} only contains the typing judgments for primitive operations.

Not captured in the type system of LambdaPix are the following two restrictions:

- No variable may appear more than once in a pattern.
- The patterns of a case expression must be exhaustive.

3.2 Dynamic Semantics

Here we define how LambdaPix expressions evaluate. We define evaluation as a small-step dynamic semantics where the judgment $e \mapsto e'$ means that e steps to e' and the judgment e val means that e is a value and doesn't step any further. LambdaPix enjoys progress and preservation.

Definition 3.2 (Progress and Preservation). For any typing context Γ and expression e such that $\Gamma \vdash e : \tau$ it is either the case that e val or there exists an e' such that $\Gamma \vdash e' : \tau$ and $e \mapsto e'$.

$$\frac{v_{1} \not / p_{1} + B_{1} \qquad c_{1} \neq c_{2}}{v_{1} \not / c_{2} + A} \text{ MATCH}_{2} \qquad \frac{c_{1} = c_{2}}{c_{1} \not / c_{2} + A} \text{ MATCH}_{3} \qquad \frac{c_{1} \neq c_{2}}{c_{1} \not / c_{2}} \text{ MATCH}_{4}$$

$$\frac{v_{1} \not / p_{1} + B_{1} \qquad v_{n} \not / p_{n} + B_{n}}{\{\ell_{1} = v_{1}, \dots, \ell_{n} = v_{n}\} \not / \{\ell_{1} = p_{1}, \dots, \ell_{n} = p_{n}\} + B_{1} \dots B_{n}} \text{ MATCH}_{5}$$

$$\frac{v_{i} \not / p_{i}}{\{\ell_{1} = v_{1}, \dots, \ell_{n} = v_{n}\} \not / \{\ell_{1} = p_{1}, \dots, \ell_{n} = p_{n}\}} \text{ MATCH}_{6} \qquad \frac{v \not / p + B}{v \not / x \text{ as } p + B, v / x} \text{ MATCH}_{7}$$

$$\frac{v \not / p}{v \not / x \text{ as } p} \text{ MATCH}_{8} \qquad \frac{i_{1} \neq i_{2}}{i_{1} \not / i_{1}} \text{ MATCH}_{10} \qquad \frac{v \not / p + B}{i \cdot v \not / i \cdot p + B} \text{ MATCH}_{11}$$

$$\frac{i_{1} \neq i_{2}}{i_{1} \cdot v \not / i_{2} \cdot p} \text{ MATCH}_{12} \qquad \frac{v \not / p}{i \cdot v \not / i \cdot p} \text{ MATCH}_{13} \qquad \frac{i_{1} \cdot v \not / i_{2}}{i_{1} \cdot v \not / i_{2}} \text{ MATCH}_{14} \qquad \frac{i_{1} \not / i_{2} \cdot p}{i_{1} \not / i_{2} \cdot p} \text{ MATCH}_{15}$$

Fig. 6. Pattern matching in LambdaPix

LambdaPix also enjoys the finality of values: it is never the case that both $e \mapsto e'$ and e val.

To define evaluation we first define two helper judgments to deal with pattern matching (Figure 6). The judgment $v \not\parallel p \dashv B$ means the value v matches to the pattern p producing B, where B is a set of bindings of the form v'/x that indicate the value v' is bound to the variable x. The judgment $v \not\parallel p$ means the expression v does not match to the pattern p. It is assumed as a precondition to these judgements that v val, v val

$$\frac{e_1 \text{ val} \quad e_2 \text{ val} \quad \dots \quad e_{i-1} \text{ val} \quad e_i \mapsto e_i'}{\{\dots, \ell_i = e_i, \dots\} \mapsto \{\dots, \ell_i = e_i', \dots\}} \text{ DyN}_2 \qquad \frac{e_1 \text{ val} \quad \dots \quad e_n \text{ val}}{\{\ell_1 = e_1, \dots, \ell_n = e_n\} \text{ val}} \text{ DyN}_3$$

$$\frac{e \mapsto e'}{e \cdot \ell_i \mapsto e' \cdot \ell_i} \text{ DyN}_4 \qquad \frac{\{\dots, \ell_i = e_i, \dots\} \text{ val}}{\{\dots, \ell_i = e_i, \dots\} \cdot \ell_i \mapsto e_i} \text{ DyN}_5 \qquad \frac{e \mapsto e'}{i \cdot e \mapsto i \cdot e'} \text{ DyN}_6 \qquad \frac{e \text{ val}}{i \cdot e \text{ val}} \text{ DyN}_7$$

$$\frac{e \mapsto e'}{\text{case } e \text{ } \{p_1.e_1 \mid \dots \mid p_n.e_n\} \mapsto \text{case } e' \text{ } \{p_1.e_1 \mid \dots \mid p_n.e_n\}} \text{ DyN}_9$$

$$\frac{e \text{ val} \quad e \text{ } p_1 \quad \dots \quad e \text{ } p_{i-1} \quad e \text{ } p_i \text{ } + B}{\text{case } e \text{ } \{\dots, p_i.e_i \mid \dots\} \mapsto [B]e_i} \text{ DyN}_{10} \qquad \frac{e_1 \mapsto e'_1}{k^2 \cdot e_1 \cdot e_2 \mapsto e'_1 \cdot e_2} \text{ DyN}_{11}$$

$$\frac{e_1 \text{ val} \quad e_2 \mapsto e'_1}{e_1 \cdot e_2 \mapsto e'_1 \cdot e'_2} \text{ DyN}_{13} \qquad \frac{e_2 \text{ val}}{(\lambda x.e) \cdot e_2 \mapsto [e_2/x]e} \text{ DyN}_{14} \qquad \frac{e_1 \mapsto e'_1}{\text{ fix } x \text{ is } e \mapsto [\text{fix } x \text{ is } e/x]e} \text{ DyN}_{15}$$

$$\frac{e \text{ val}}{e \text{ } e_2 \mapsto e'_1} \text{ DyN}_{16} \qquad \frac{e \text{ val}}{e \mapsto e'} \text{ DyN}_{17} \qquad \frac{v \text{ val}}{v \mapsto v} \text{ BigDyN}_1 \qquad \frac{e \mapsto e'}{e \mapsto v} \text{ BigDyN}_2$$

Fig. 7. Dynamic semantics of LambdaPix

In Figure 7 we use these helper judgments to define the evaluation judgments. In DyN₁₇, e' is meant to be understood as a hard-coded value dependent on the primitive operation o. We use

these judgments to define what it means for an expression to evaluate to a value. We use $e \mapsto v$ to denote that expression e evaluates to value v. In rules $BigDyn_1$ and $BigDyn_2$, big-step dynamics are defined as the transitive closure of the small-step dynamics.

4 SOUND EQUIVALENCE INFERENCES

Our approach takes as input two LambdaPix expressions of the same type and outputs a logic formula which is valid only if the two expressions are equivalent. We construct this logic formula by constructing a proof tree of sound equivalence inferences.

4.1 Logic Formulas

$$\sigma ::= \begin{array}{ccc} \sigma_1 & \equiv t_2 & term \ equivalence \\ \sigma_1 \wedge \sigma_2 & conjunction \\ \sigma_1 \vee \sigma_2 & disjunction \\ \sigma_1 \Rightarrow \sigma_2 & implication \\ \neg \sigma & negation \\ \end{array}$$

Fig. 8. Logic Formulas

$$\frac{t \text{ Term }}{c \text{ Term }} \text{ Term}_1 \qquad \frac{t_1 \text{ Term }}{x \text{ Term }} \text{ Term}_2 \qquad \frac{t_1 \text{ Term }}{\{\ell_1 = t_1, \dots, \ell_n = t_n\} \text{ Term }}{\{\ell_1 = t_1, \dots, \ell_n = t_n\} \text{ Term }} \text{ Term}_3 \qquad \frac{t \text{ Term }}{t \cdot \ell_i \text{ Term }} \text{ Term}_4 \qquad \frac{t \text{ Term }}{t \cdot \ell_i \text{ Term }} \text{ Term}_4 \qquad \frac{t \text{ Term }}{x \text{ as } t \text{ Term }} \text{ Term}_8 \qquad \frac{t \text{ Term }}{o \text{ Term }} \text{ Term}_9 \qquad \frac{t \text{ Term }}{v \text{ Term }} \qquad \frac{t \text{ Term }}{v \text{ Term }} \text{ Term}_9 \qquad \frac{t \text{ Term }}{v \text{ Term }} \qquad$$

Fig. 9. Term Judgment

Figure 8 defines the form of the formulas generated by our approach. The leaves of these formulas are equalities between base terms t, defined in Figure 9. These base terms encode three things: LambdaPix values, patterns, and the application of a primitive operation and a value. Encoding all of these things as terms allows a term equivalence to state that either two values are the same, that a value matches with a pattern, or that a primitive operation application yields a value that is equal to another value or matches with a pattern.

The inclusion of primitive operation applications as base terms is somewhat strange since they are not values in the actual dynamics of LambdaPix, but this inclusion enables the resulting formula to include all of the information pertaining to the theory from which the primitive operation originates. For instance, since the theory of quantifier-free linear integer arithmetic knows that addition is commutative, this inclusion makes it possible for the expressions $\lambda x.\lambda y.(x+y)$ and $\lambda x.\lambda y.(y+x)$ to be recognized as equivalent.

Except when a variable, primitive operation, as pattern, or wildcard pattern is included in one of the terms, term equivalence is identical to syntactic equality. When a primitive operation is included in a term, the specific primitive operation is used to determine how to understand the term equivalence (e.g. $1 + 2 \equiv 3$ is a valid term equivalence using the primitive operation "+").

When an as pattern is included in a term equivalence: x as $e_1 \equiv e_2$, the term equivalence is the same as $x \equiv e_1 \land e_1 \equiv e_2$. When a wildcard pattern is included in a term equivalence: $\underline{} \equiv e$, the term equivalence can simply be interpreted as "true".

When one or more free variables are included in a formula, they must be resolved to determine the formula's truth. Throughout our approach, contexts are used to keep track of the types of all of a formula's free variables. Expressions can be substituted for variables of the same type in a formula to resolve it (e.g. $[3/x](x \equiv 1 \land x \equiv 2)$ yields $3 \equiv 1 \land 3 \equiv 2$). A formula is valid if it is true under all possible substitutions of its variables. To denote this, we define a new form of judgment:

Definition 4.1 $(\forall_{\Gamma}.j)$. If $\Gamma = \vec{x} : \vec{\tau}$, then the judgement $\forall_{\Gamma}.j$ holds if for all \vec{v} where $v_i : \tau_i$ and v_i val for all $v_i \in \vec{v}$, it is the case that $[\vec{v}/\vec{x}]j$ holds. Implicitly, although the types of primitive operations are included in Γ_{initial} , and therefore Γ , we omit typings of the form $o : \tau_1 \to \tau_2$ from $\vec{x} : \vec{\tau}$ so that we do not range over all possible meanings for LambdaPix's primitive operations. Then if Γ is a typing context with a mapping for every free variable in a formula σ , the validity of σ is denoted $\forall_{\Gamma}.\sigma$.

The validity of formulas will be what determines whether our approach recognizes two LambdaPix expressions as equivalent. Our approach takes as input two LambdaPix expressions and uses them to output a logic formula. In Section 6, we show that if the output formula is valid by Definition 4.1, then the two expressions are necessarily equivalent. To define our approach's method of constructing the logic formula from the original LambdaPix expressions in Section 4.4, we begin by first defining a few helper judgments pertaining to weak head reduction and freshening.

4.2 Weak Head Reduction $e \downarrow e'$

We do not have the option of fully evaluating the expressions during execution, as expressions may contain free variables in redex positions. For this reason we use weak head reduction at each step; this eliminates head-position redexes until free variables get in the way. The result is a weak head normal form expression.

$$\frac{e \leadsto e' \quad e' \downarrow e''}{e \downarrow e''} \text{ BigWhnf}_1 \qquad \frac{e \not \leadsto}{e \downarrow e} \text{ BigWhnf}_2 \qquad \frac{e_1 \leadsto e'_1}{e_1 \ e_2 \leadsto e'_1 \ e_2} \text{ Whnf}_1$$

$$\frac{(\lambda x. e_1) e_2 \leadsto [e_2/x] e_1}{e_1 \bowtie e_2 \leadsto e'_1 \bowtie e'} \text{ Whnf}_3 \qquad \frac{e \leadsto e'}{e \cdot \ell_i \leadsto e' \cdot \ell_i} \text{ Whnf}_3 \qquad \frac{e_1 \leadsto e'_1}{e_1 \bowtie e_2 \leadsto e'_1 \bowtie e} \text{ Whnf}_4$$

Fig. 10. Weak Head Reduction

4.3 Freshening

It is sometimes useful to generate fresh variables (globally unique variables) to avoid variable capture. As single variables are not the only form of binding sites in LamdbaPix, we generalize this notion to patterns. When freshen $p.e \hookrightarrow p'.e'$, p'.e' is the same as p.e except all variables bound by p are alpha-varied to fresh variables. The definition of the freshen judgment is given in Figure 11.

In addition to creating fresh variables to avoid variable capture, our approach sometimes generates fresh variables in order to couple the binding sites between two expressions being considered. For instance, if our approach knows that the same expression e is being matched to variable x in one expression and variable y in another expression, it is useful to equate these bindings so that

$$\frac{y \text{ fresh}}{\text{freshen } \underline{.e} \hookrightarrow \underline{.e}} \text{ Freshen}_1 \qquad \frac{y \text{ fresh}}{\text{freshen } x.e \hookrightarrow y.[y/x]e} \text{ Freshen}_2$$

$$\frac{\text{freshen } p_1.e \hookrightarrow p'_1.e_1 \text{ freshen } p_2.e_1 \hookrightarrow p'_2.e_2 \dots \text{ freshen } p_n.e_{n-1} \hookrightarrow p'_n.e_n}{\text{freshen } \{\ell_1 = p_1, \dots, \ell_n = p_n\}.e \hookrightarrow \{\ell_1 = p'_1, \dots, \ell_n = p'_n\}.e_n} \text{ Freshen}_3$$

$$\frac{y \text{ fresh } \text{ freshen } p.e \hookrightarrow p'.e'}{\text{freshen } x \text{ as } p.e \hookrightarrow y \text{ as } p'.[y/x]e'} \text{ Freshen}_4 \qquad \frac{\text{freshen } p.e \hookrightarrow p'.e'}{\text{freshen } i.e \hookrightarrow i.e} \text{ Freshen}_5$$

Fig. 11. Freshening

$$\frac{y \text{ fresh}}{\mathsf{EB}(_.e_1,_.e_2)} \hookrightarrow (_.e_1,_.e_2) \to (_.e_1,_.e_2) \to (B_1 \\ \frac{y \text{ fresh}}{\mathsf{EB}(_.e_1,x.e_2)} \hookrightarrow (y.e_1,y.[y/x]e_2) \to (B_2 \\ \frac{y \text{ fresh}}{\mathsf{EB}(x.e_1,_.e_2)} \hookrightarrow (y.[y/x]e_1,y.e_2) \to (B_3 \\ \frac{y \text{ fresh}}{\mathsf{EB}(x.e_1,x'.e_2)} \hookrightarrow (y.[y/x]e_1,y.[y/x']e_2) \to (B_4 \\ \frac{y \text{ fresh}}{\mathsf{EB}(x \text{ as } p_1.e_1,p_2.e_2)} \to (y \text{ as } p'.[y/x]e'_1,p'.e'_2) \to (B_5 \\ \frac{y \text{ fresh}}{\mathsf{EB}(p_1.e_1,p_2.e_2)} \to (y \text{ as } p'.[y/x]e'_1,y \text{ as } p'.e'_2) \to (B_5 \\ \frac{y \text{ fresh}}{\mathsf{EB}(p_1.e_1,x \text{ as } p_2.e_2)} \to (y \text{ as } p'.e'_1,p'.e'_2) \to (B_6 \\ \frac{\mathsf{EB}(p_1.e_1,p'_1.e_2)}{\mathsf{EB}(p_1.e_1,p'_1.e_2)} \to (p''_1.e'_1,p''_1.e'_2) \to (B_6 p_1.e'_1,p'_1.e'_2) \to (B_6 p_1.e'_1.e'_1.e'_1.e'_1.e'_2) \to (B_6 p_1.e'_1.e'_1.e'_1.e'_2) \to (B_6 p_1.e'_1.e'_1.e'_1.e'_1.e'_2) \to (B_6 p_1.e$$

Fig. 12. Equate Bindings Judgment

as our approach proceeds, it is able to know that x in the first expression is the same as y in the second expression. The judgment $\mathsf{EB}(p_1.e_1,p_2.e_2) \hookrightarrow (p.e_1',p.e_2')$ defined in Figure 12 does exactly that, taking in two bindings and returning freshened versions of those bindings that use the same variables so long as the two bindings $p_1.e_1$ and $p_2.e_2$ can be alpha-varied to use a shared pattern p.

The benefit of the equate bindings judgment specifically comes into play when comparing case expressions. If two case expressions are casing on the same e, and they have identical or near identical binding structures, then it is sometimes useful to freshen the case expressions together, so that as our approach proceeds to consider all of the possible outcomes of the case expressions, it is able to know that the same e was bound in the same way in both expressions. The judgment $\mathsf{FT}(\{p_1.e_1 \mid \ldots \mid p_n.e_n\}, \{p_1'.e_1' \mid \ldots \mid p_m'.e_m'\}) \overset{s}{\hookrightarrow} (\{p_1''.e_1'' \mid \ldots \mid p_n''.e_n''\}, \{p_1'''.e_1''' \mid \ldots \mid p_n'''.e_n'''\})$

$$\frac{\mathsf{EB}(p_{1}.e_{1},p_{2}.e_{2})\hookrightarrow(p.e_{1}',p.e_{2}')}{\mathsf{FT}(\{p_{1}.e_{1}\mid\cdot\},\{p_{2}.e_{2}\mid\cdot\})\overset{1}{\hookrightarrow}(\{p.e_{1}'\mid\cdot\},\{p.e_{2}'\mid\cdot\})}\;\mathsf{FT}_{1}}{\frac{\mathsf{EB}(p_{1}.e_{1},p_{2}.e_{2})\hookrightarrow(p.e_{1}',p.e_{2}')\quad\mathsf{FT}(rest_{1},rest_{2})\overset{n}{\hookrightarrow}(rest_{1}',rest_{2}')}{\mathsf{FT}(\{p_{1}.e_{1}\mid rest_{1}\},\{p_{2}.e_{2}\mid rest_{2}\})\overset{n+1}{\longleftrightarrow}(\{p.e_{1}'\mid rest_{1}'\},\{p.e_{2}'\mid rest_{2}'\})}}\;\mathsf{FT}_{2}}{\frac{\mathsf{V}_{i\in[n]}(\mathsf{freshen}\;p_{i}.e_{i}\hookrightarrow p_{i}''.e_{i}'')\quad\mathsf{V}_{i\in[m]}(\mathsf{freshen}\;p_{i}'.e_{i}'\hookrightarrow p_{i}''.e_{i}''')}{\mathsf{FT}(\{p_{1}.e_{1}\mid\ldots\mid p_{n}.e_{n}\},\{p_{1}'.e_{1}'\mid\ldots\mid p_{m}'.e_{m}'\})\overset{0}{\hookrightarrow}(\{p_{1}''.e_{1}''\mid\ldots\mid p_{n}''.e_{n}''\},\{p_{1}''.e_{1}'''\mid\ldots\mid p_{m}''.e_{m}'''\})}}\;\mathsf{FT}_{3}$$

Fig. 13. Freshen Together Judgment

defined in Figure 13 takes in two lists of bindings from case expressions, and equates the first s bindings, independently freshening the rest. The judgment is defined so that once a pair of bindings cannot be equated, all subsequent bindings are freshened independently. This is done to ensure that no bindings are unsoundly equated. The rules listed in Figure 13 are listed in order of precedence (i.e. if it is possible to apply FT_2 or FT_3 , it will apply FT_2).

4.4 Formula Generation $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$

The judgment that connects the validity of logic formulas with the equivalence of LambdaPix expressions is $\Gamma \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$. The judgment that defines how our approach generates said logic formulas is $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$.

When $\Gamma \vdash e_1 \stackrel{\sigma}{\iff} e_2 : \tau \dashv \Gamma'$ or $\Gamma \vdash e_1 \stackrel{\sigma}{\iff} e_2 : \tau \dashv \Gamma'$, the only free variables appearing in e_1 and e_2 are in Γ , so $\Gamma \vdash e_1 : \tau$ and $\Gamma \vdash e_2 : \tau$. However, σ can contain more free variables than just those in Γ . The purpose of Γ' is to describe the rest of the variables in σ . Γ and Γ' are disjoint and between them account for all variables which may appear in σ .

$$\frac{e_1 \downarrow e_1' \quad e_2 \downarrow e_2' \quad \Gamma \vdash e_1' \overset{\sigma}{\longleftrightarrow} e_2' : \tau \dashv \Gamma';}{\Gamma \vdash e_1 \overset{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'} \text{ IsoExp}$$

Fig. 14. IsoExp Rule

The judgment $\Gamma \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ is defined by Figure 14 and is mutually recursive with $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$. We use it to define what it means for two expressions to be isomorphic.

Definition 4.2 (Isomorphic). We call two expressions e_1 and e_2 where $\Gamma_{\text{initial}} \vdash e_1 : \tau$ and $\Gamma_{\text{initial}} \vdash e_2 : \tau$ isomorphic if $\Gamma_{\text{initial}} \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ and $\stackrel{\text{val}}{\forall}_{\Gamma'} . \sigma$.

The purpose of the distinction between the two judgments is to allow our approach to perform weak head reduction exactly when needed. The judgment $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$ assumes as a precondition that e_1 and e_2 are in weak head normal form, and is defined by Figures 15, 16, and 17.

Each rule in Figure 15 is written to address a particular syntactic form that e_1 and e_2 might take. Since each rule targets a particular syntactic form, the premises of each rule are motivated by the

$$\frac{\Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau \quad e_1 \text{ Term} \quad e_2 \text{ Term}}{\Gamma \vdash e_1 \overset{e_1 \equiv e_2}{\longleftrightarrow} e_2 : \tau \dashv \cdot} \text{ Iso}_{\text{atomic}}$$

$$\frac{\Gamma \vdash e_1 \overset{\sigma_1}{\longleftrightarrow} e_1' : \tau_1 \dashv \Gamma_1' \quad \dots \quad \Gamma \vdash e_n \overset{\sigma_n}{\longleftrightarrow} e_n' : \tau_n \dashv \Gamma_n'}{\Gamma \vdash \{\ell_1 = e_1, \dots, \ell_n = e_n\}} \overset{\sigma_1 \land \dots \land \sigma_n}{\longleftrightarrow} \{\ell_1 = e_1', \dots, \ell_n = e_n'\} : \{\ell_1 : \tau_1, \dots, \ell_n : \tau_n\} \dashv \Gamma_1', \dots, \Gamma_n'} \text{ Iso}_{\text{record}}$$

$$\frac{\Gamma \vdash e_1 \overset{\sigma}{\longleftrightarrow} e_2 : \{\dots, \ell_i : \tau_i, \dots\} \dashv \Gamma'}{\Gamma \vdash e_1 \vdash \ell_i \overset{\sigma}{\longleftrightarrow} e_2 : \ell_i : \tau_i \dashv \Gamma'} \text{ Iso}_{\text{projection}}$$

$$\frac{i : \tau \hookrightarrow \delta \quad \Gamma \vdash e_1 \overset{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'}{\Gamma \vdash i \cdot e_1 \overset{\sigma}{\longleftrightarrow} i \cdot e_2 : \delta \dashv \Gamma'} \text{ Iso}_{\text{injection}}$$

$$\frac{x \text{ fresh} \quad \Gamma, x : \tau \vdash [x/x_1]e_1 \overset{\sigma}{\longleftrightarrow} [x/x_2]e_2 : \tau' \dashv \Gamma'}{\Gamma \vdash \lambda x_1 . e_1 \overset{\sigma}{\longleftrightarrow} \lambda x_2 . e_2 : \tau \to \tau' \dashv x : \tau, \Gamma'} \text{ Iso}_{\text{lambda}}$$

$$\frac{x \text{ fresh} \quad \Gamma, x : \tau \vdash [x/x_1]e_1 \overset{\sigma}{\longleftrightarrow} [x/x_2]e_2 : \tau \dashv \Gamma'}{\Gamma \vdash \text{fix } x_1 \text{ is } e_1 \overset{\sigma}{\longleftrightarrow} \text{fix } x_2 \text{ is } e_2 : \tau \dashv x : \tau, \Gamma'}} \text{ Iso}_{\text{fix}}$$

Fig. 15. Formula Generation Rules

semantics of that form. For example, $\operatorname{Iso}_{\operatorname{lambda}}$ has the premises x fresh and $\Gamma, x : \tau \vdash [x/x_1]e_1 \stackrel{\sigma}{\Leftrightarrow} [x/x_2]e_2 : \tau' \dashv \Gamma'$. The former premise simply declares x as a previously unused variable, and the latter premise states that if any value x of type τ (the input type to both expressions) is substituted for x_1 in the left expression and x_2 in the right expression, then the two expressions will be equivalent if σ is valid. This reflects the fact that two functions are equivalent if and only if their outputs are equivalent for all valid inputs.

Although the soundness of these rules is guaranteed, their completeness is not. For instance, Iso_{projection} has the premise $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \{\dots, \ell_i : \tau_i, \dots\} \dashv \Gamma'$. If this premise holds, then the conclusion that $\Gamma \vdash e_1 \cdot \ell_i \stackrel{\sigma}{\longleftrightarrow} e_2 \cdot \ell_i : \tau_i \dashv \Gamma'$ necessarily follows, as if two records are equivalent, then each of the records' respective entries must also be equivalent. But it is not the case that in order for two projections to be equivalent, they must project from equivalent records.

Each rule in Figure 16 addresses the case in which at least one of the expressions being compared is an application. When the two expressions being compared are both applications of equivalent arguments onto equivalent functions, $\operatorname{IsO}_{\operatorname{application1}}$ can be used to infer equivalence of the resulting applications. For situations in which an application is being compared to another syntactic form, or two applications that cannot be recognized as equivalent via $\operatorname{IsO}_{\operatorname{application1}}$ are being compared, the remaining rules take an application and replace it with a shared fresh variable in both expressions. For example, if the expressions f(x) and f(x+0) are being compared, $\operatorname{IsO}_{\operatorname{application1}}$ is sufficient to find equivalence because f can be found equivalent to f and f(x) and f(x) are being compared, $\operatorname{IsO}_{\operatorname{application1}}$ alone would be insufficient, as the outermost function of the first expression is f and the outermost function of the second expression is f. For this situation, $\operatorname{IsO}_{\operatorname{application2}}$ is needed to replace f(x) with the fresh variable f0, yielding the expressions f1 and f2, which can be immediately found equivalent via $\operatorname{IsO}_{\operatorname{atomic}}$ 2.

The current formula generation application rules have multiple limitations. First, the rules only allow applications to be replaced with shared fresh variables when the application being replaced is at the outermost level of one of the expressions. This has the consequence that although

$$\frac{\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \to \tau' \dashv \Gamma' \quad \Gamma \vdash e_1' \stackrel{\sigma'}{\longleftrightarrow} e_2' : \tau \dashv \Gamma''}{\Gamma \vdash e_1 e_1' \stackrel{\sigma \land \sigma'}{\longleftrightarrow} e_2 e_2' : \tau' \dashv \Gamma', \Gamma''} \quad \text{Iso}_{application1}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash y \stackrel{\sigma}{\longleftrightarrow} [y/(x e_1)] e_2 : \tau \dashv \Gamma'}{\Gamma \vdash x e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'} \quad \text{Iso}_{application2}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash [y/(x e_2)] e_1 \stackrel{\sigma}{\longleftrightarrow} y : \tau \dashv \Gamma'}{\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} x e_2 : \tau \dashv \Gamma'} \quad \text{Iso}_{application3}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash y \stackrel{\sigma}{\longleftrightarrow} [y/(o e_1)] e_2 : \tau \dashv \Gamma'}{\Gamma \vdash o e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'} \quad \text{Iso}_{application4}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash [y/(o e_2)] e_1 \stackrel{\sigma}{\longleftrightarrow} y : \tau \dashv \Gamma'}{\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} o e_2 : \tau \dashv \Gamma'} \quad \text{Iso}_{application5}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash y \stackrel{\sigma}{\longleftrightarrow} [y/((\text{fix } x_1 \text{ is } e_1) e_2)] e : \tau \dashv \Gamma'}{\Gamma \vdash (\text{fix } x_1 \text{ is } e_1) e_2 \stackrel{\sigma}{\longleftrightarrow} e : \tau \dashv \Gamma'} \quad \text{Iso}_{application6}$$

$$\frac{y \text{ fresh} \quad \Gamma, y : \tau \vdash [y/((\text{fix } x_1 \text{ is } e_1) e_2)] e \stackrel{\sigma}{\longleftrightarrow} y : \tau \dashv \Gamma'}{\Gamma \vdash (\text{fix } x_1 \text{ is } e_1) e_2 : \tau \dashv \Gamma'} \quad \text{Iso}_{application6}$$

Fig. 16. Formula Generation Application Rules

f(x) + f(x) and 2 * f(x) are obviously equivalent, and the substitution of f(x) for a shared fresh variable y would enable Iso_{atomic} to prove that fact, our current rules do not support this inference. Second, $Iso_{application6}$ and $Iso_{application7}$ require substituting an entire fixed point application in an expression, so unless if the two expressions being compared have essentially identical fixed points included, these rules will be ineffective. Still, despite these limitations, the current formula generation application rules are sufficient for their most common purpose of working with Iso_{fix} to ensure that recursive function calls are recognized as equivalent when given equivalent arguments.

Each rule in Figure 17 addresses the situation in which at least one of the expressions being compared is a case analysis. These rules can be grouped into two broad approaches. For situations in which only one of the expressions being compared is a case analysis, or both expressions are case analyses but the expressions being cased on are not equivalent, $\operatorname{Iso}_{\operatorname{casel}}$ and $\operatorname{Iso}_{\operatorname{case2}}$ are used to unpack one case analysis at a time. If case $e\{p_1.e_1 \mid \ldots \mid p_n.e_n\}$ is being compared to e', then the formula generated by these rules states that if e can be pattern matched with p_i and no prior patterns, e_i needs to be equivalent to e' in order for the two overall expressions to be equivalent.

For situations in which the two expressions being compared are case analyses that are casing on equivalent expressions, Iso_{case3} , Iso_{case4} , and Iso_{case5} are used to deconstruct both case expressions simultaneously. To do this, Iso_{case3} is always used first to ensure that the expressions being cased on are equivalent. If the expressions being cased on are not equivalent, then σ in the formula generated by Iso_{case3} will not be valid, and so the output formula $\sigma \wedge \sigma'$ will not be valid

$$\frac{e \ \mathsf{Term} \quad \forall_{i \in [n]} \left(\mathsf{freshen} \ p_i.e_i \hookrightarrow p_i'.e_i' \quad p_i' :: \tau' \dashv \Gamma_i \quad \Gamma, \Gamma_i \vdash e_i' \overset{\sigma_i}{\Longleftrightarrow} e' : \tau \dashv \Gamma_i' \right)}{\Gamma \vdash \mathsf{case} \ e \ \{p_1.e_1 \mid \ldots \mid p_n.e_n\}} \overset{\wedge_{i \in [n]}((\wedge_{j \in [i-1]}(e \not\equiv p_j')) \land e \equiv p_i') \Rightarrow \sigma_i}{\longleftrightarrow} e' : \tau \dashv \forall_{i \in [n]} \Gamma_i, \Gamma_i' \end{aligned} \\ \frac{e \ \mathsf{Term} \quad \forall_{i \in [n]} \left(\mathsf{freshen} \ p_i.e_i \hookrightarrow p_i'.e_i' \quad p_i' :: \tau' \dashv \Gamma_i \quad \Gamma, \Gamma_i \vdash e_i' \overset{\sigma_i}{\Longleftrightarrow} e' : \tau \dashv \Gamma_i' \right)}{\mathsf{Iso}_{\mathsf{case2}}} \\ \frac{e \ \mathsf{Term} \quad \forall_{i \in [n]} \left(\mathsf{freshen} \ p_i.e_i \hookrightarrow p_i'.e_i' \quad p_i' :: \tau' \dashv \Gamma_i \quad \Gamma, \Gamma_i \vdash e_i' \overset{\sigma_i}{\Longleftrightarrow} e' : \tau \dashv \Gamma_i' \right)}{\mathsf{Iso}_{\mathsf{case2}}} \\ \frac{e \ \mathsf{Term} \quad \forall_{i \in [n]} \left(\mathsf{freshen} \ p_i.e_i \hookrightarrow p_i'.e_i' \quad p_i' :: \tau' \dashv \Gamma_i \quad \Gamma, \Gamma_i \vdash e_i' \overset{\sigma_i}{\Longleftrightarrow} e' : \tau \dashv \Gamma_i' \right)}{\mathsf{Iso}_{\mathsf{case2}}} \\ \frac{e \ \mathsf{Term} \quad \forall_{i \in [n]} \left(\mathsf{freshen} \ p_i.e_i \hookrightarrow p_i'.e_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i' \mid \Gamma, \Gamma_i' \mid P_i' :: \tau' \dashv \Gamma_i$$

Fig. 17. Formula Generation Case Rules

as a result. If the expressions being cased on are equivalent, then Iso_{case4} and Iso_{case5} can be used to generate σ' . This approach is needed in addition to Iso_{case1} and Iso_{case2} because Iso_{case1} and Iso_{case2} require that the expression being cased on is a base term.

All rules in Figure 15 are deterministic in the sense that for all possible expressions, at most one rule is applicable. However, the rules in Figures 16 and 17 are non-deterministic. If two case expressions or two applications are being compared, there may be multiple applicable rules. For instance, if case 1 {1.2|_.3} is being compared to case 2 {_.2}, then $\rm Iso_{case1}$, $\rm Iso_{case2}$, and $\rm Iso_{case3}$ are all applicable. Our approach handles this by considering all formulas that can be generated by applying any applicable rule and outputs the disjunction of all generated formulas. We will later show that applying any applicable rule in such a situation is sound and that therefore, taking the disjunction of all generated formulas is also sound. The only exception to this is that $\rm Iso_{case3}$ cannot be applied multiple times in a row because it is never useful to do so and allowing this would cause an infinite loop.

In instances where there is no applicable rule, such as if i_1 e_1 is compared with i_2 e_2 where $i_1 \neq i_2$ and either e_1 or e_2 cannot be encoded into a term, our approach simply outputs the formula σ =False, which is always sound.

4.5 Limitations and Further Extensions

The current set of rules is comprehensive and covers a wide range of operators that are often found in many functional programming assignments. However, there are limitations to the current set of rules, some of which have already been noted. The current main limitations include:

- Our current handling of projections in Iso_{projection} requires that in order for two projections to be recognized as equivalent, they must project from equivalent records.
- Our current handling of recursive function calls occurs entirely through the interplay between
 Iso_{fix} and the formula generation application rules. Because of how these rules are currently
 defined, recursive functions can only be recognized as equivalent if in all situations they
 recurse on equivalent arguments or do not recurse at all.
- The approach taken by the formula generation application rules is limited in that applications can only be replaced with shared fresh variables when the application being replaced is at the outermost level of one of the expressions being compared.
- Iso_{application6} and Iso_{application7} both require substituting a variable for an entire fixed point application, which will only be useful if the two expressions being compared have essentially identical fixed points included.
- Since Iso_{casel} and Iso_{case2} require that the expression being cased on is a base term, the current set of rules cannot identify equivalence between a case analysis in which the expression being cased on isn't a base term and any other syntactic form.
- The current definition of LambdaPix does not allow for state, and so our approach cannot identify the equivalence of any programs that use state.

Compared to other potential extensions that could be implemented to address an aforementioned limitation, extending LambdaPix to support state would likely require a significant number of changes to our approach. However, this could be potentially achieved by handling sequential state-altering declaration similar to how we handle local declaration. Currently, we handle local declaration by encoding the declaration into the SMT formula in the same way that we would encode a single pattern case expression (i.e. let val $x = e_1$ in e_2 end becomes case e_1 { $x.e_2$ } at the transpilation to LambdaPix stage). It would not be possible to do the same procedure for sequential declaration since the scoping would have to be global. However, we believe that it may be feasible to treat reference declaration/assignment similar to variable declaration/initialization and reference update similar to variable shadowing with modified scoping.

One advantage of the structure of our approach is that extending our system to address some of the previously listed limitations is straightforward. As soon as a new rule that addresses one of the system's current limitations is found to be sound, it can be simply tacked on to the current system without needing to modify any preexisting rules. This also applies to extensions of the underlying language LambdaPix itself. Adding new base types to LambdaPix such as strings or reals requires no modification of the current rules whatsoever, and adding additional syntactic expression forms requires only the addition of rules for comparing the new form against itself and arbitrary expressions. Even though it is easy to extend the LambdaPix language and add additional rules, the current version is already rich enough to capture common behavior in programming assignments of introductory courses.

5 OPERATION

To provide a better understanding of our approach, we step through our approach's operation on a pair of simple Standard ML expressions provided above. As we step through this example, we will refer to the inference rules from the previous section to illustrate how they are applied.

```
fun bind a f =
fun add_opt x y =
  case (x, y) of
                                               case a of
    (SOME m, SOME n) =>
                                                SOME b => f b
      SOME (m + n)
                                               | NONE => NONE
  | (NONE, _) => NONE
  | (_, NONE) => NONE
                                            val return = SOME
                                             fun add_opt x y =
                                               bind x (fn m \Rightarrow
                                               bind y (fn n \Rightarrow
                                                 return (m + n)
                                               ))
```

First, we transpile both expressions to LambdaPix. This is shown above. Since much of the proof derivation which drives our approach is free of branching, through most of this section we will view our approach as transforming the above expressions through the application of rules, rather than building up a proof tree.

```
λχ.λγ.
                                                               \lambda x . \lambda y .
   case (x,y) of
                                                                  (\lambda a.\lambda f.
   { (SOME \cdot m, SOME \cdot n) . SOME \cdot (m+n)
                                                                     case a of
   | (NONE,_).NONE
                                                                     { SOME·b.f b
   | (_,NONE).NONE }
                                                                     | NONE.NONE }
                                                                  ) x (\lambda m)
                                                                  (\lambda a.\lambda f.
                                                                     case a of
                                                                     { SOME·b.f b
                                                                     | NONE.NONE }
                                                                  ) y (\lambda n.
                                                                     (\lambda e. SOME \cdot e) (m+n)
```

The entry point to our approach is the $\Gamma \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ judgement, defined by the rule IsoExp. By this rule, we reduce both expressions to weak head normal form then apply the $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$ judgement to them. However, since the expressions in consideration are abstractions, the expressions are already in weak head normal form, so no transformation is necessary to apply this rule.

```
case (x,y) of
                                                              (\lambda a.\lambda f.
{ (SOME \cdot m, SOME \cdot n). SOME \cdot (m+n)
                                                                 case a of
| (NONE,_).NONE
                                                                 { SOME · b . f b
| (_,NONE).NONE }
                                                                | NONE.NONE }
                                                              ) \times (\lambda m.
                                                               (\lambda a.\lambda f.
                                                                  case a of
                                                                  { SOME·b.f b
                                                                  | NONE.NONE }
                                                               ) y (λn.
                                                                  (\lambda e.SOME \cdot e) (m+n)
                                                              ))
```

Next, since both expressions are lambda expressions with two curried arguments, we proceed with two applications of the rule Iso_{lambda} . This requires us to create two new fresh variables and substitute them for the first two function arguments in both expressions. For simplicity, we will simply call the first fresh variable x and the second fresh variable y even though these names conflict with the original variable names. The key difference between before and after this process is that before this process, the two functions had the same variable names x and y by coincidence, whereas after this process, the two functions use the same fresh variables x and y by design. These applications of Iso_{lambda} yield the above expressions.

```
case (x,y) of (SOME \cdot m, SOME \cdot n) \cdot SOME \cdot (m+n) { SOME \cdot b \cdot b \cdot (\lambda m \cdot n) \cdot NONE } (\lambda a \cdot \lambda f \cdot b \cdot n) \cdot (\lambda a \cdot \lambda f \cdot b \cdot n) case a of { SOME \cdot b \cdot f \cdot b \cdot n | NONE \cdot NONE \cdot b \cdot n | NONE \cdot NONE \cdot n) | NONE \cdot NONE \cdot n | NON
```

Since the premise of Iso_{lambda} invokes the \Leftrightarrow judgement, IsoExp requires that we reduce both expressions to weak head normal form. The left expression is already in weak head normal form, so no transformation is necessary, but the right expression must undergo two beta reductions before it is in weak head normal form. The result of these beta reductions is above.

Since both expressions are case expressions, our approach has multiple options for how to proceed. Formally, our approach pursues all of these options, generating separate formulas for each option, and finally outputting a disjunction of all of the generated formulas. This ensures that if any option can generate a valid formula, then the final result will be the disjunction of the valid formula with several other formulas, which altogether is valid. In this case, attempting to proceed with $\rm Iso_{case3}$ will not yield a valid formula because the two expressions are casing on different things, but applying either $\rm Iso_{case1}$ or $\rm Iso_{case2}$ can yield a valid formula. For this demonstration, we step through the derivation that results from applying $\rm Iso_{case1}$ and call the formulas generated by applying $\rm Iso_{case2}$ or $\rm Iso_{case3}$ $\sigma_{\rm Iso_{case3}}$ and $\sigma_{\rm Iso_{case3}}$ respectively.

As there are three branches in the left case expression, our approach's proof tree now splits into three branches. For this demonstration, we just step through the first of these branches, as the other two branches work similarly. We call the formulas generated by the other two branches of the proof tree $\sigma_{\text{branch 2}}$ and $\sigma_{\text{branch 3}}$.

```
SOME·(m1+n1) case x of  \{ \text{ SOME·b.} \\ (\lambda m. \\ (\lambda a. \lambda f. \\ \text{ case a of } \\ \{ \text{ SOME·b.f b} \\ | \text{ NONE. NONE } \}  ) y (\lambda n. (\lambda e. \text{SOME·e}) (m+n)) ) b | NONE. NONE }
```

We "freshen" the branch selected to avoid variable capture. In this situation we will freshen the first branch of the left case expression by replacing m and n with m1 and n1, respectively. From this

branch we will generate a formula of the form

```
((x,y) \equiv (SOME \cdot m1, SOME \cdot n1)) \Rightarrow ...
```

where the ellipses is what we are going to fill in as we complete this branch of the proof tree.

Since the left expression has been simplified to a base term, our approach proceeds to work on the right expression. Our approach applies Iso_{case2} twice (using beta reduction to reduce the expression to weak head normal form as appropriate), and finishes each branch of the proof tree by using Iso_{atomic} to compare base terms.

Putting everything together, the final formula is:

```
(\sigma_{branch\ 1} \land \sigma_{branch\ 2} \land \sigma_{branch\ 3}) \lor \sigma_{Iso_{case2}} \lor \sigma_{Iso_{case3}} where \sigma_{branch\ 1} is ((x,y) \equiv (SOME \cdot m1, SOME \cdot n1)) \Rightarrow \\ ((x \equiv SOME \cdot b1) \Rightarrow \\ (y \equiv SOME \cdot b2) \Rightarrow (SOME \cdot (m1 + n1) \equiv SOME \cdot (b1 + b2)) \land \\ (y \not\equiv SOME \cdot b2 \land y \equiv NONE) \Rightarrow (SOME \cdot (m1 + n1) \equiv NONE) \\ ) \land \\ ((x \not\equiv SOME \cdot b1 \land x \equiv NONE) \Rightarrow \\ (y \equiv SOME \cdot b2) \Rightarrow (SOME \cdot (m1 + n1) \equiv NONE) \land \\ (y \not\equiv SOME \cdot b2 \land y \equiv NONE) \Rightarrow (SOME \cdot (m1 + n1) \equiv NONE) \\ )
```

and $\sigma_{\text{branch 2}}$ and $\sigma_{\text{branch 3}}$ are similar.

Since the two original expressions were equivalent, this formula is valid. The validity of this formula can be verified either by hand or by an SMT Solver.

6 SOUNDNESS

We prove the soundness of our approach: if our approach takes in two expressions and outputs a valid formula, then the two expressions must be equivalent.

6.1 Extensional Equivalence

To prove the soundness of our approach, we must first define what it means for two expressions to be equivalent. For this, we introduce extensional equivalence, a widely accepted notion of equivalence. Extensional equivalence is the same as contextual equivalence, and so two extensionally equivalent expressions are indistinguishable in terms of behavior. This implies that extensional equivalence is closed under evaluation. Extensional equivalence is also an equivalence relation, so we may assume that it is reflexive, symmetric, and transitive. LambdaPix enjoys referential transparency, meaning that extensional equivalence of LambdaPix expressions is closed under replacement of subexpressions with extensionally equivalent subexpressions.

We use $e_1 \cong e_2 : \tau$ to denote that expressions e_1 and e_2 are extensionally equivalent and both have the type τ .

Definition 6.1 (Extensional Equivalence). We define that $e_1 \cong e_2 : \tau$ if $\Gamma_{initial} \vdash e_1 : \tau$, $\Gamma_{initial} \vdash e_2 : \tau$, $e_1 \mapsto v_1$, $e_2 \mapsto v_2$, and

- (1) Rule EQ₁: In the case that $\tau = \tau_1 \to \tau_2$, for all expressions v such that $\Gamma_{\text{initial}} \vdash v : \tau_1$, $v_1 \ v \cong v_2 \ v : \tau_2$.
- (2) Rule EQ₂: In the case that τ is not an arrow type, for all patterns p such that p :: τ , either $v_1 /\!\!/ p \dashv B$ and $v_2 /\!\!/ p \dashv B$ or $v_1 /\!\!/ p$ and $v_2 /\!\!/ p$.

Unlike our approach, extensional equivalence inducts over the types of the expressions rather than their syntax, and is defined only over closed expressions. As we are only concerned with proving our approach sound over valuable expressions, we leave extensional equivalence undefined for divergent expressions.

This is an atypical formalization of extensional equivalence; it is typically defined in terms of the elimination forms of each type connective. However, since pattern matching in LambdaPix subsumes the elimination of all connectives other than arrows, we simply define equivalence at all non-arrow types in terms of pattern matching.

The soundness theorem for our approach connects our technique's definition of isomorphic with this definition of extensional equivalence. It is as follows:

Theorem 6.2 (Soundness). For any expressions e_1 and e_2 , if $\Gamma_{\text{initial}} \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ and $\stackrel{\text{val}}{\forall}_{\Gamma'}.\sigma$, then $e_1 \cong e_2 : \tau$.

6.2 Proof Sketch

As the $\Gamma \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ judgement is defined simultaneously with the $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$ judgement, we prove the theorem by simultaneous induction on both of these judgements. We also use the $\overset{\text{val}}{\forall}_{\Gamma}.j$ judgement to strengthen the inductive hypotheses to account for variables. Recall that if $\Gamma = \vec{x} : \vec{\tau}$, then the judgement $\overset{\text{val}}{\forall}_{\Gamma}.j$ holds if for all \vec{v} where $v_i : \tau_i$ and v_i val for all $v_i \in \vec{v}$, it is the case that $[\vec{v}/\vec{x}]j$ holds (implicitly, we omit any primitive operations from the context $\Gamma = \vec{x} : \vec{\tau}$ as to not range over all possible meanings for LambdaPix's primitive operations). The theorem we wish to show by induction is then:

• If
$$\Gamma \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$$
 then $\stackrel{\mathsf{val}}{\forall}_{\Gamma}$. $\left(\text{if} \left(\stackrel{\mathsf{val}}{\forall}_{\Gamma'} . \sigma \right) \text{ then } e_1 \cong e_2 : \tau \right)$.

• If $\Gamma \vdash e_1 \stackrel{\sigma}{\longleftrightarrow} e_2 : \tau \dashv \Gamma'$ then $\stackrel{\mathsf{val}}{\forall}_{\Gamma}$. $\left(\text{if} \left(\stackrel{\mathsf{val}}{\forall}_{\Gamma'} . \sigma \right) \text{ then } e_1 \cong e_2 : \tau \right)$.

We first verify that the above statements imply the soundness theorem. Indeed, when $\Gamma_{\text{initial}} \vdash \Gamma_{\text{initial}} = \Gamma_$

We first verify that the above statements imply the soundness theorem. Indeed, when $\Gamma_{\text{initial}} \vdash e_1 \stackrel{\sigma}{\Longleftrightarrow} e_2 : \tau \dashv \Gamma'$ we have $\forall \Gamma_{\text{Initial}} \cdot \left(\text{if} \left(\begin{matrix} \forall \mathsf{q}' \\ \forall \mathsf{q}' \cdot \sigma \end{matrix} \right) \right)$ then $e_1 \cong e_2 : \tau \right)$. Since Γ_{initial} contains only primitive operations, which are omitted from the $\nabla_{\Gamma} \cdot J$ judgment, the outer quantifier quantifies over no variables, so we have that if $\left(\begin{matrix} \mathsf{q} \\ \forall \mathsf{q}' \cdot \sigma \end{matrix} \right)$ then $e_1 \cong e_2 : \tau$. This together with the assumption that $\nabla_{\Gamma} \cdot \sigma \in \mathcal{T}$ allows us to conclude that $e_1 \cong e_2 : \tau$.

The full proof of each rule's soundness has 18 cases and uses 14 lemmas and can be found in the extended version of this paper [Clune et al. 2020]. Two cases are included below as examples:

Iso_{record}: Let $\Gamma = \vec{x} : \vec{\tau}$ and let \vec{v} be arbitrary where $v_i : \tau_i$ and v_i val for all $v_i \in \vec{v}$. Assume $[\vec{v}/\vec{x}] \begin{pmatrix} \mathsf{val} \\ \forall_{\Gamma_1', \dots, \Gamma_n'} . \sigma_1 \wedge \dots \sigma_n \end{pmatrix}$. It must be shown that $[\vec{v}/\vec{x}] (\{\ell_1 = e_1, \dots, \ell_n = e_n\}) \cong \{\ell_1 = e_1', \dots, \ell_n = e_n'\}$).

Lemma 6.3. If $e_1 \cong e_2 : \tau$, $e_1 \mapsto v_1$, and $e_2 \mapsto v_2$, then for all patterns p where $p :: \tau \dashv \Gamma$, it is the case that either $v_1 \not\parallel p \dashv B$ and $v_2 \not\parallel p \dashv B$ or $v_1 \not\parallel p$ and $v_2 \not\parallel p$.

Proof: by induction on $e_1 \cong e_2 : \tau$. If $\tau = \tau_1 \to \tau_2$ then by inversion of $p :: \tau \dashv \Gamma$, p must either be a wildcard or a variable. Then by MATCH₁ and MATCH₂, we have that $v_1 \not\parallel p \dashv B$ and $v_2 \not\parallel p \dashv B$. If τ isn't an arrow type, then we conclude by EQ_2 .

By conjunction and that all the Γ'_i are disjoint, we have that for all $i \in [n]$, $\forall \Gamma'_{i} \cdot \sigma_i$. Then by the inductive hypotheses, we have that $[\vec{v}/\vec{x}](e_i \cong e'_i : \tau_i)$. Since we are only concerned with proving our approach sound over valuable expressions, without loss of generality, we can assume that $[\vec{v}/\vec{x}]e_i \mapsto v_i$ and $[\vec{v}/\vec{x}]e'_i \mapsto v'_i$ for some values v_i and v'_i . By Lemma 6.3, we have that for all p_i where $p_i :: \tau_i \dashv \Gamma_i$, either $v_i \not| p_i \dashv B_i$ and $v'_i \not| p_i \dashv B_i$ or $v_i \not| p_i$ and $v'_i \not| p_i$.

To appeal to EQ₂, let p be an arbitrary pattern such that $p:\{\ell_1:\tau_1,\ldots,\ell_n:\tau_n\}\dashv\Gamma'$. We proceed by cases:

- In the case that for all $i \in [n]$ $v_i /\!\!/ p_i + B_i$ and $v_i' /\!\!/ p_i + B_i$, by MATCH₅ we have $\{\ell_1 = v_1, \ldots, \ell_n = v_n\} /\!\!/ p + B_1 \ldots B_n$ and $\{\ell_1 = v_1', \ldots, \ell_n = v_n'\} /\!\!/ p + B_1 \ldots B_n$.
- In the case that there is some $i \in [n]$ where $v_i \not | p_i$ and $v_i' \not | p_i$, by MATCH₆ we have $\{\ell_1 = v_1, \ldots, \ell_n = v_n\} \not | p$ and $\{\ell_1 = v_1', \ldots, \ell_n = v_n'\} \not | p$.

Since in all cases either $\{\ell_1 = v_1, \dots, \ell_n = v_n\} \ /\!\!/ p + B$ and $\{\ell_1 = v_1', \dots, \ell_n = v_n'\} \ /\!\!/ p + B$ or $\{\ell_1 = v_1, \dots, \ell_n = v_n\} \ /\!\!/ p$ and $\{\ell_1 = v_1', \dots, \ell_n = v_n'\} \ /\!\!/ p$, by EQ₂, we may conclude

$$[\vec{v}/\vec{x}](\{\ell_1 = e_1, \dots, \ell_n = e_n\}) \cong \{\ell_1 = e'_1, \dots, \ell_n = e'_n\})$$

Iso_{application2}: Let $\Gamma = \vec{z} : \vec{\tau}$ and let \vec{v} be arbitrary where $v_i : \tau_i$ and v_i val for all $v_i \in \vec{v}$. Assume $[\vec{v}/\vec{z}] \begin{pmatrix} v_1 \\ \forall r \end{pmatrix}$. It must be shown that $[\vec{v}/\vec{z}](x e_1 \cong e_2)$.

By the inductive hypothesis we have

$$\overset{\text{val}}{\forall}_{\Gamma,y:\tau}.\left(\text{if }\left(\overset{\text{val}}{\forall}_{\Gamma'}.\sigma\right)\text{ then }y\cong[y/(x\ e_1)]e_2:\tau\right)$$

Since we are only concerned with proving our approach sound over valuable expressions, without loss of generality, we can assume that $x e_1 \mapsto w$ for some value w such that $\Gamma \vdash w : \tau$ and w val. Since y is fresh, the inductive hypothesis written above implies

if
$$[w/y][\vec{v}/\vec{z}] \begin{pmatrix} val \\ \forall \Gamma'.\sigma \end{pmatrix}$$
 then $[w/y][\vec{v}/\vec{z}](y \cong [y/(x e_1)]e_2 : \tau)$

By assumption, we already have $[w/y][\vec{v}/\vec{z}] \begin{pmatrix} val \\ \forall \Gamma'.\sigma \end{pmatrix}$. Therefore we have

$$[w/y][\vec{v}/\vec{z}](y \cong [y/(x e_1)]e_2 : \tau)$$

which is equivalent to

$$[\vec{v}/\vec{z}](w \cong [w/(x e_1)]e_2 : \tau)$$

Since $x e_1 \mapsto w$, the two are extensionally equivalent. By the referential transparency of LambdaPix, the above expression is equivalent to

$$[\vec{v}/\vec{z}](x e_1 \cong [(x e_1)/(x e_1)]e_2 : \tau)$$

which is simply

$$[\vec{v}/\vec{z}](x e_1 \cong e_2 : \tau)$$

Proc. ACM Program. Lang., Vol. 4, No. OOPSLA, Article 171. Publication date: November 2020.

7 EXPERIMENTAL RESULTS

We implemented our approach in a tool called zeus to serve as a grading assistant by clustering equivalent programs into equivalence classes. The goal of our evaluation is to answer the following:

- Q1. Can zeus automatically identify equivalent programs in programming assignments for introductory functional programming courses?
- Q2. How many equivalence classes are found by zeus?
- Q3. What is the runtime performance of ZEUS?

7.1 Implementation

ZEUS is implemented in Standard ML and is publicly available as open-source at https://github.com/CMU-TOP/zeus. zeus takes as input a set of homework assignments from an introductory functional programming course at the college level taught in Standard ML. Each submission is transpiled from Standard ML into LambdaPix, and then zeus is run pairwise on the transpiled expressions and outputs a logical formula. If this formula is valid, then both expressions are algorithmically similar and guaranteed to be equivalent, so they are placed into the same equivalence class. As an optimization, since extensional equivalence is transitive, if zeus verifies that two programs p_1 and p_2 are (not) equivalent, and that p_1 is also (not) equivalent to p_3 , then zeus does not check that p_2 is equivalent to p_3 . This optimization significantly reduces the number of comparisons that otherwise would be quadratic in the number of assignments.

In the definition of LambdaPix, we assumed an arbitrary fixed set of disjoint algebraic datatypes with unique associated injection labels. This is unrealistic for an implementation since Standard ML includes datatype declarations. Our transpilation from Standard ML to LambdaPix instead scrapes all datatype declarations from the original Standard ML submission and uses those datatypes and their constructors as LambdaPix's set of datatypes and injection labels.

To determine the validity of the formulas generated by ZEUS, we use the SMT solver Z3 [de Moura and Bjørner 2008] using the theory of quantifier-free linear integer arithmetic and the theory of datatypes. From the theory of quantifier-free linear integer arithmetic, we use the built-in functions "+", "-", "*", "≤", "<", and ">", corresponding to the primitive operations of LambdaPix. We use the theory of datatypes to represent base terms of all types aside from *ints* and *booleans*.

Although we only use these two theories in ZEUS, nothing restricts a different implementation from using additional theories. For instance, another implementation could leverage the theory of *strings* by adding strings as a base type in LambdaPix and adding the SMT solver's built-in string functions to LambdaPix's set of primitive operations.

7.2 Benchmarks

To evaluate zeus, we used more than 4,000 student submissions from an introductory functional programming course. The number of submissions varies between 318 and 351 per assignment. Table 1 describes the twelve assignments that were used in our evaluation. These assignments show a large diversity of programs that includes different datatypes and the use of pattern matching and are a good test suite to test the applicability of zeus as a grading assistant. Figure 18 shows some of the datatype declarations that are assumed by the homework assignments presented in Table 1.

7.3 Clustering of equivalent programs

Tables 2 and 3 analyze the equivalent classes detected by ZEUS. In particular, for each task, Table 2 shows the number of submissions (#), the number of equivalent classes (ECs), the number of equivalence classes that contain 90% and 75% of the submissions (90th and 75th Percentile ECs, respectively), the number of equivalent classes containing more than 1 submission (Non-singleton

Function	Signature	Description
concat	int list list → int list	concat takes a list of int lists and returns their con-
		catenation without using the built-in "@" function
prefixSum	int list → int list	prefixSum replaces each i-th element in an int list
		with the sum of the list's first $i + 1$ elements
countNonZero	int tree → int	countNonZero takes an int tree T and returns the
		number of nonzero nodes in T
quicksort	$('a * 'a \rightarrow order) * 'a list \rightarrow 'a list$	quicksort implements the quicksort algorithm
slowDoop	$(a * a \rightarrow order) * a list$	slowDoop takes a comparison function and uses
	→ 'a list	it to remove all duplicates in a list
differentiate	$(int \to real) \to (int \to $	differentiate differentiates a polynomial that is
	real)	represented with the type int \rightarrow real
integrate	$(int \to real) \to real \to$	integrate takes a polynomial p and a real c and
	$(int \rightarrow real)$	returns the antiderivative of p with constant of
		integration c
treefoldr	$(a * b \rightarrow b) \rightarrow b \rightarrow a$	(treefoldr g init T) returns (foldr g init L) where L
	tree → 'b	is the inorder traversal of T
treeReduce	$(a * a \rightarrow a) \rightarrow a \rightarrow a$	treeReduce is the same as treefoldr except that it
	tree → 'a	must have $O(\log n)$ span assuming g is associative and init is an identity for g
findN	$(a \rightarrow bool) \rightarrow (a * a \rightarrow$	(findN p eq T n s k) returns s $[x1,, xn]$ where
	bool) \rightarrow 'a shrub \rightarrow int	[x1,, xn] are the leftmost values for T such that
	\rightarrow ('a list \rightarrow 'b) \rightarrow (unit	for all i from 1 to n , p xi returns true and the xi's
	\rightarrow 'b) \rightarrow 'b	are eq-distinct. (findN p eq T n s k) returns k() if
	(/	no such [x1,, xn] exist
sat	$prop \rightarrow ((string * bool)$	sat takes in a proposition, a success function s
	$\underset{,}{\text{list}} \rightarrow \text{'a)} \rightarrow (\text{unit} \rightarrow \text{'a})$	from a list assigning booleans to free variables to
	d → 'a	'a, and a failure function from unit to 'a. If the
		proposition is satisfiable by an assignment of free
		variables A, then sat returns s(A). Otherwise, it returns k()
findPartition	'a list \rightarrow ('a list \rightarrow bool)	(findParition A pL pR) returns true if there exist
	\rightarrow ('a list \rightarrow bool) \rightarrow	an L and R such that (L, R) is a partition of A where
	bool	pL accepts L and pR accepts R. (findPartition A pL
		pR) returns false otherwise

Table 1. Description of homework assignments used in our evaluation

Fig. 18. Datatype declarations assumed by homework assignments

ECs), and the percentage of submissions found equivalent to at least one other submissions (% in Non-singleton ECs). Table 3 shows the number of correct and incorrect student submissions, the

	#	ECs	90th Percentile ECs	75th Percentile ECs	Non- singleton ECs	% in Non- singleton ECs
treefoldr	332	22	3	1	9	96
integrate	323	34	4	1	5	91
slowDoop	347	30	6	2	12	95
countNonZero	351	29	8	4	13	95
concat	351	40	8	2	10	91
treeReduce	332	57	24	8	20	89
prefixSum	351	68	33	7	23	87
differentiate	316	65	34	3	6	81
quicksort	347	73	39	9	18	84
findN	330	73	40	7	18	83
findPartition	331	83	50	12	22	82
sat	318	104	73	25	14	72

Table 2. Analysis of the number of equivalent classes (ECs)

Table 3. Analysis of correctness of student submissions

	Correct Submissions	Correct ECs	Non- singleton Correct ECs	Incorrect Submissions	Incorrect ECs	Non- singleton Incorrect ECs
treefoldr	302	12	4	30	10	5
integrate	307	23	3	16	11	2
slowDoop	346	29	12	1	1	0
countNonZero	346	26	12	5	3	1
concat	336	30	9	15	10	1
treeReduce	188	18	8	144	39	12
prefixSum	347	64	23	4	4	0
differentiate	308	59	5	8	6	1
quicksort	328	56	16	19	17	2
findN	296	48	14	34	25	4
findPartition	291	59	13	40	24	9
sat	273	72	11	45	32	3

number of equivalence classes containing only correct or incorrect submissions, and the number of equivalence classes containing multiple correct or incorrect submissions. There were no equivalence classes that contained both correct and incorrect submissions.

A common trend among all tasks was that a significant majority of student submissions were placed into a relatively small number of large equivalence classes, with the remaining submissions widely dispersed among many small equivalence classes, frequently of size 1. For instance, for the task concat, zeus detected 40 equivalent classes. However, only 10 of those classes contain more than one submission, and 8 equivalence classes contain more than 90% of the submissions. In almost all tasks, the largest equivalence classes consisted of various distinct but correct solutions to the problem. The one exception to this trend was that in the task treeReduce, a significant number

	#	Total Time (s)	Number of Comparisons	Average Time (s)
treefoldr	332	33.701	694	0.049
integrate	323	45.833	991	0.046
slowDoop	347	57.564	1,201	0.048
countNonZero	351	72.268	1,578	0.046
concat	351	76.110	1,614	0.047
treeReduce	332	146.830	3,089	0.048
prefixSum	351	205.695	4,112	0.050
differentiate	316	140.797	3,025	0.047
quicksort	347	210.554	4,303	0.049
findN	330	202.577	3,660	0.055
findPartition	331	284.502	4,807	0.059
sat	318	486.218	6,532	0.074

Table 4. Runtime analysis

of students mistook associativity for commutativity or otherwise assumed that the function passed into treeReduce was necessarily commutative. This common misunderstanding resulted in a large number of incorrect submissions for treeReduce, but because the misunderstanding was common, zeus was still able to place the majority of incorrect submissions into a small number of large equivalence classes. In all tasks, at least 72% of submissions were identified as equivalent to at least one other submission. These results support the hypothesis that zeus can be used as a grading assistant to reduce the workload of instructors in reviewing equivalent code, thus freeing their time to provide more detailed feedback.

7.4 Runtime performance

Table 4 shows the time needed by ZEUS to cluster all assignments for a given task when running on a common Mac laptop with a 1.6GHz processor and 4 GB of RAM. Specifically, for each task, it shows the number of submissions, the total time to cluster submissions in seconds, the number of pairwise comparisons performed during clustering, and the average time for a single comparison in seconds. The average time to compare two individual submissions is small and it ranges from 0.046 seconds to 0.074 seconds. When performing the clustering of a given assignment, we can observe that the number of comparisons is much less than quadratic and that the total time varies between 1 and 8 minutes. This shows that ZEUS is efficient in practice and can be used in real-time to help instructors grade assignments.

7.5 Discussion

We manually inspected the cases where zeus did not put two programs in the same equivalence class. The most common reasons for this were the following:

• The two programs are not equivalent: since these programs correspond to actual student submissions, not all of the programs are correct. When an incorrect implementation produces the wrong output on any number of inputs, our algorithm appropriately puts it in a different equivalence class from the correct submissions. Additionally, for the sat task, the correct behavior of this function when an input proposition is satisfiable by multiple assignments is not fully defined. If multiple assignments A satisfy the proposition, there are no rules about

which A to use when returning s(A). So for this task, two correct submissions could produce different outputs.

- The two programs use different recursive helper functions: we found cases where equivalence classes were distinguished by the structure of the helper functions students created. Since our current inference rules do not consider these cases, zeus fails to recognize that two programs are equivalent if they use recursive helper functions with different input structures.
- The two recursive programs use different base cases: our algorithm's treatment of fixed points causes it to never peer into a recursive call. Our algorithm's treatment of case expressions causes it to only recognize two expressions as equivalent if they handle all inputs in basically the same way. Together, these have the implication that when one expression treats a certain input as a base case while the other expression treats it as a recursive case, then the algorithm will be unable to recognize the expressions as equivalent.
- One of the programs uses built-in Standard ML functions: seven out of the twelve tasks involve list manipulation operations. For instance, the top five tasks with the largest number of equivalence classes (sat, findPartition, findN, quicksort, and prefixSum) correspond to tasks that involve list manipulation. Many of the submissions for these tasks use built-in Standard ML functions for list reversal or list concatenation. We did not use a theory of list structures in our SMT Solver, so we were only able to recognize two expressions as equivalent if they used these built-in functions on the same input inputs and order or if they did not use these built-in functions at all.

We note that even with the current limitations, zeus already shows that it can efficiently cluster the majority of the submissions into a few equivalence classes. Also, zeus could be extended by adding additional inference rules or support for additional SMT theories that would allow the identification of equivalent programs that are currently missed by zeus.

8 RELATED WORK

Proving that two problems are equivalent is a well-studied topic and has many applications ranging from hardware equivalence [Berman and Trevillyan 1989], compiler optimizations [Zuck et al. 2002], to program equivalence [Godlin and Strichman 2009]. However, the use of program equivalence for grading programming assignments is scarce [Kaleeswaran et al. 2016]. In this section, we cover related work from program equivalence and automatic grading that is closer to our approach.

8.1 Program Equivalence

Program Verification. The problem of program equivalence can be reduced to a verification problem by showing that both programs satisfy the same specification. For instance, model-checking techniques [Clarke et al. 2004, 2001] can be used to show that two C programs satisfy the same specification. This specification can be written to ensure that for the same input, the programs are equivalent if they always produce the same output. Fedyukovich et al. [Fedyukovich et al. 2016] present techniques for proving that two similar programs have the same property rather than being equivalent. Their approach requires formally verifying one of the programs and using this proof to check the validity of the property in the other program by establishing a coupling between the two programs. A similar approach can also be done for functional programs. For instance, one could write a formal specification of the functionality of a program in Why3ML [Bobot et al. 2015]. We tried this approach by writing a formal specification for programming assignments for the function concat, however, the Why3 framework [Bobot et al. 2015] was not able to prove that the program satisfied the specification. In general, proving the program equivalence concerning a specification is a more challenging task than the one we address in this paper since we can take advantage of program structure to prove that they are equivalent.

Regression Verification. In regression verification [Felsing et al. 2014; Godlin and Strichman 2009], the goal is to prove that two versions of a program are equivalent. One approach is to transform loops in programs to recursive procedures and to match the recursive calls in both programs and abstract them via uninterpreted functions [Godlin and Strichman 2009]. Other approaches use invariant inference techniques to prove the equivalence of programs with loops [Felsing et al. 2014]. By using these techniques, one can encode the two versions of the program into Horn clauses and use constraint solvers to automatically find certain kinds of invariants. Alternatively, one can also use symbolic execution and static analysis to generate summaries of program behaviors that capture the modifications between the programs. These summaries can be encoded into logical formulas and their equivalence can be checked using SMT solvers [Backes et al. 2013]. Our approaches also consider that student submissions are similar but they are not different versions of the same program. Even though we do not use any invariant generation techniques, this is orthogonal to our approach and could increase the number of equivalent classes detected for recursive programs.

Contextual Equivalence. There is a broad set of work that targets contextual equivalence for functional programs. Approaches based on step-indexed logical relations [Ahmed et al. 2009; Ahmed 2006; Dreyer et al. 2009] or on bisimulations [Hur et al. 2012; Koutavas and Wand 2006; Sumii and Pierce 2005] have been used to prove context equivalence of functional programs with different fragments of ML that often include finite datatypes and integer references. While these approaches are more theoretical and focus on functional programs with state, we do not support state but can handle pattern matching which is crucial for a practical tool to cluster programming assignments of introductory functional courses. The closest approach to ours is the one recently presented by Jaber [Jaber 2020]. Jaber presents techniques for checking the equivalence of OCaml programs with state. His approach focuses in particular on contextual equivalence and developing a framework in which references can be properly accounted for. Our approach neglects references, as we require programs to be purely functional, but includes a more comprehensive treatment of datatypes. We attempted to compare our Zeus's performance against Jaber's SyTeC1 prototype, but unfortunately, all of our benchmarks included datatypes that were not supported by the available prototype.

8.2 Automatic grading

Clustering similar assignments. To help instructors to grade programming assignments, several automatic techniques have been proposed to cluster similar assignments into buckets with the purpose of giving automatic feedback [Gulwani et al. 2018; Kaleeswaran et al. 2016; Pu et al. 2016; Wang et al. 2018]. Our approach differs from these since our goal is not to replace the instructor or to fully automate the grading but rather to use ZEUS as a grading assistant with formal guarantees.

CLARA [Gulwani et al. 2018] cluster correct programs and selects a canonical program from each cluster to be considered as the reference solution. In this approach, a pair of programs p_1 and p_2 are said to be dynamic equivalent if they have the same control-flow and if related variables in p_1 and p_2 always have the same values, in the same order, during the program execution on the same inputs. In contrast, our approach has a stronger notion of equivalence since we do not depend on dynamic program analysis. CodeAssist [Kaleeswaran et al. 2016] clusters submissions for dynamic programming assignments by their solution strategy. They consider a small set of features and if two programs share these features then they are put in the same cluster. Other clustering approaches are based on deep learning techniques [Pu et al. 2016] and also provide no formal guarantees about the quality of the clustering. SemCluster [Perry et al. 2019] improves upon other clustering techniques by considering semantic program features. They use control flow features and data flow features to represent each program and merge this information to create a program feature vector. K-means clustering is used to cluster all programs based on the program

feature vectors. Even though there are no formal guarantees for the equivalence of programs in each cluster, experimental results [Perry et al. 2019] show that the number of clusters found by SemCluster is much smaller than competitive approaches.

Automatic repair. AutoGrader [Singh et al. 2013] takes as input a reference solution and an error model that consists of potential corrections and uses constraint solving techniques to find a minimum number of corrections that can be used to repair the incorrect student solution. Sarfgen [Wang et al. 2018] uses a three-stage algorithm based on search, align, and repair. It starts by searching for a small number of correct programs that can be used to repair the incorrect submission and have the same control-flow structure. Next, they compute a syntactic distance between those programs using an embedding of ASTs into numerical vectors. These programs are then aligned and the differences between aligned statements can suggest corrections that can be repaired automatically.

Automatic repair is better suited for Massive Open Online Courses where a fully automated method is needed, while our approach is better suited for large, in-person courses, where the feedback of instructors can be more beneficial. The feedback returned by automatic repair tools is limited to changes in the code, while our approach is meant to assist instructors to provide more detailed feedback for students. Each equivalent class will have specific comments that are more helpful to the student than a repaired version of their submission. Moreover, while our approach can be used for both correct and incorrect submissions, automatic repair is only useful to fix incorrect submissions and cannot give any feedback for different implementations of correct submissions.

Formal guarantees. Liu et al. [Liu et al. 2019] proposes to automatically determine the correctness of an assignment against a reference solution. Instead of using test cases, they use symbolic execution to search for semantically different execution paths between a student's submission and the reference solution. If such paths exist, then the submission is considered incorrect and feedback can be provided by using counterexamples based on path deviations. Our approach is not based on symbolic execution but instead uses inference rules to derive a formula for which both student submissions are equivalent if and only if they have the same structure and the observable behavior.

CodeAssist [Kaleeswaran et al. 2016] checks equivalence of a candidate submission from a cluster with a correct solution of that cluster that has been previously validated by an instructor. They exploit the fact of just handling dynamic programming assignments to establish a correspondence between variables and control locations of the two programs. Using this correspondence, they can encode the problem into SMT and prove program equivalence. Our approach is more general since our inference rules simulate relationships between expressions of the two programs and can be applied to several problem domains and not just dynamic programming assignments.

9 CONCLUSION

We present techniques for checking for equivalence between purely functional programs. Guided by inference rules that inform needed equivalences between two programs' subexpressions, our approach simultaneously deconstructs the expressions being compared to build up a formula that is valid only if the expressions are equivalent. We prove the soundness of our approach: if our approach takes in two expressions of the same type and outputs a valid formula, then the two expressions are equivalent. We implement our approach and show that it can assist grading by clustering over 4,000 real student code submissions from an introductory functional programming class taught at the undergraduate level.

ACKNOWLEDGMENTS

This work was partially funded by National Science Foundation (Grants CCF-1901381, CCF-1762363, and CCF-1629444).

REFERENCES

- Amal Ahmed, Derek Dreyer, and Andreas Rossberg. 2009. State-dependent representation independence. In *Proc. Symposium on Principles of Programming Languages*. ACM, 340–353. https://doi.org/10.1145/1480881.1480925
- Amal J. Ahmed. 2006. Step-Indexed Syntactic Logical Relations for Recursive and Quantified Types. In *Proc. European Symposium on Programming*, Springer, 69–83. https://doi.org/10.1007/11693024_6
- John D. Backes, Suzette Person, Neha Rungta, and Oksana Tkachuk. 2013. Regression Verification Using Impact Summaries. In Proc. International Symposium Model Checking Software. Springer, 99–116. https://doi.org/10.1007/978-3-642-39176-7_7
- C Leonard Berman and Louise H Trevillyan. 1989. Functional comparison of logic designs for VLSI circuits. In *Proc. International Conference on Computer-Aided Design*. IEEE, 456–459. https://doi.org/10.1109/ICCAD.1989.76990
- François Bobot, Jean-Christophe Filliâtre, Claude Marché, and Andrei Paskevich. 2015. Let's verify this with Why3. Int. J. Softw. Tools Technol. Transf. 17, 6 (2015), 709–727. https://doi.org/10.1007/s10009-014-0314-5
- Edmund Clarke, Daniel Kroening, and Flavio Lerda. 2004. A tool for checking ANSI-C programs. In *Proc. International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 168–176. https://doi.org/10.1007/978-3-540-24730-2 15
- Edmund M. Clarke, Armin Biere, Richard Raimi, and Yunshan Zhu. 2001. Bounded Model Checking Using Satisfiability Solving. Formal Methods Syst. Des. 19, 1 (2001), 7–34. https://doi.org/10.1023/A:1011276507260
- Joshua Clune, Vijay Ramamurthy, Ruben Martins, and Umut A. Acar. 2020. Program Equivalence for Assisted Grading of Functional Programs (Extended Version). *CoRR* (2020).
- Leonardo Mendonça de Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In Proc. International Conference on Tools and Algorithms for the Construction and Analysis of Systems. Springer, 337–340. https://doi.org/10.1007/978-3-540-78800-3 24
- Derek Dreyer, Amal Ahmed, and Lars Birkedal. 2009. Logical Step-Indexed Logical Relations. In *Proc. Annual Symposium on Logic in Computer Science*. IEEE Computer Society, 71–80. https://doi.org/10.1109/LICS.2009.34
- Grigory Fedyukovich, Arie Gurfinkel, and Natasha Sharygina. 2016. Property Directed Equivalence via Abstract Simulation. In *Proc. International Conference Computer-Aided Verification*. Springer, 433–453. https://doi.org/10.1007/978-3-319-41540-6_24
- Dennis Felsing, Sarah Grebing, Vladimir Klebanov, Philipp Rümmer, and Mattias Ulbrich. 2014. Automating regression verification. In *Proc. International Conference on Automated Software Engineering*. ACM, 349–360. https://doi.org/10.1145/2642937.2642987
- Benny Godlin and Ofer Strichman. 2009. Regression verification. In *Proc. Design Automation Conference*. ACM, 466–471. https://doi.org/10.1145/1629911.1630034
- Sumit Gulwani, Ivan Radicek, and Florian Zuleger. 2018. Automated clustering and program repair for introductory programming assignments. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 465–480. https://doi.org/10.1145/3192366.3192387
- Chung-Kil Hur, Derek Dreyer, Georg Neis, and Viktor Vafeiadis. 2012. The marriage of bisimulations and Kripke logical relations. In *Proc. Symposium on Principles of Programming Languages*. ACM, 59–72. https://doi.org/10.1145/2103656.
- Guilhem Jaber. 2020. SyTeCi: automating contextual equivalence for higher-order programs with references. *PACMPL* 4, POPL (2020), 59:1–59:28. https://doi.org/10.1145/3371127
- Shalini Kaleeswaran, Anirudh Santhiar, Aditya Kanade, and Sumit Gulwani. 2016. Semi-supervised verified feedback generation. In *Proc. International Symposium on Foundations of Software Engineering*. ACM, 739–750. https://doi.org/10. 1145/2950290.2950363
- Vasileios Koutavas and Mitchell Wand. 2006. Small bisimulations for reasoning about higher-order imperative programs. In *Proc. Symposium on Principles of Programming Languages*. ACM, 141–152. https://doi.org/10.1145/1111037.1111050
- Xiao Liu, Shuai Wang, Pei Wang, and Dinghao Wu. 2019. Automatic grading of programming assignments: an approach based on formal semantics. In *Proc. International Conference on Software Engineering: Software Engineering Education and Training.* IEEE / ACM, 126–137. https://doi.org/10.1109/ICSE-SEET.2019.00022
- David Mitchel Perry, Dohyeong Kim, Roopsha Samanta, and Xiangyu Zhang. 2019. SemCluster: clustering of imperative programming assignments based on quantitative semantic features. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation*. 860–873. https://doi.org/10.1145/3314221.3314629
- Yewen Pu, Karthik Narasimhan, Armando Solar-Lezama, and Regina Barzilay. 2016. sk_p: a neural program corrector for MOOCs. In *Proc. International Conference on Systems, Programming, Languages and Applications: Software for Humanity*. ACM, 39–40. https://doi.org/10.1145/2984043.2989222
- Rishabh Singh, Sumit Gulwani, and Armando Solar-Lezama. 2013. Automated feedback generation for introductory programming assignments. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 15–26. https://doi.org/10.1145/2491956.2462195

- Eijiro Sumii and Benjamin C. Pierce. 2005. A bisimulation for type abstraction and recursion. In *Proc. Symposium on Principles of Programming Languages*. ACM, 63–74. https://doi.org/10.1145/1040305.1040311
- Ke Wang, Rishabh Singh, and Zhendong Su. 2018. Search, align, and repair: data-driven feedback generation for introductory programming exercises. In Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, 481–495. https://doi.org/10.1145/3192366.3192384
- Lenore Zuck, Amir Pnueli, Yi Fang, and Benjamin Goldberg. 2002. VOC: A translation validator for optimizing compilers. Electronic notes in theoretical computer science 65, 2 (2002), 2–18. https://doi.org/10.1016/S1571-0661(04)80393-1