

# Casts and Costs: Harmonizing Safety and Performance in Gradual Typing

ANONYMOUS AUTHOR(S)

Gradual typing allows programmers to use both static and dynamic typing in a single program. However, a well-known problem with sound gradual typing is that the interactions between static and dynamic code can cause significant performance degradation. These performance pitfalls are hard to predict and resolve, and discourage users from using gradual typing features. For example, when migrating to a more statically typed program, often adding a type annotation will trigger a slowdown that can be resolved by adding more annotations elsewhere, but since it's not clear where the additional annotations must be added, the easier solution is to simply remove the annotation.

To address these problems, we develop: (1) a static cost semantics that accurately predicts the overhead of static-dynamic interactions in a gradually typed program, (2) a technique for efficiently inferring such costs for all combinations of inferrable type assignments in a program, and (3) a method for translating the results of this analysis into specific recommendations and explanations that can help programmers understand, debug, and optimize the performance of gradually typed programs. We have implemented our approach in *HERDER*, a tool for statically analyzing the performance of different typing configurations for Reticulated Python programs. An evaluation on 15 Python programs shows that *HERDER* can use this analysis to accurately and efficiently recommend type assignments that optimize the performance of these programs without sacrificing the safety guarantees provided by static typing.

## 1 INTRODUCTION

Static and dynamic typing have different strengths and weaknesses. Gradual typing [Siek and Taha 2006; Tobin-Hochstadt and Felleisen 2006] attempts to combine the strengths of both by allowing programmers to statically type some parts of their program while dynamically typing other parts. Ideally, programmers could easily migrate between more or less statically typed programs by adding or removing type annotations. Intuitively, more static programs might be expected to have better performance and reliability, while more dynamic programs are more flexible and can be executed even if they are statically ill-typed. Unfortunately, migrating gradually typed programs is difficult [Campora et al. 2018; Tobin-Hochstadt et al. 2017] and can cause reliability and performance issues [Allende et al. 2014]. In particular, Takikawa et al. [2016] observed that adding type annotations can cause a more than 100 times slowdown in Typed Racket.

### 1.1 Performance Problem of Gradual Typing

To illustrate the performance implications of migrating between gradually typed programs, consider the program in Figure 1(a) for computing spectral norms of matrices, which was adapted from the Python Benchmark Suite.<sup>1</sup> The loop labels ( $l_1$ ,  $l_2$ ,  $l_3$ , and  $l_4$ ) in the figure can be ignored for now. *Reticulated Python* [Vitousek et al. 2014] is a gradually typed variant of Python, where type annotations can be added, using the Python type hints syntax [van Rossum et al. 2014], to introduce static checking. For example, we could add type annotations to the `eval_A` function as shown below.

```
def eval_A(i:float, j:float)->float
```

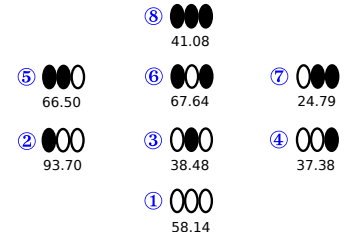
In general, we can separately decide whether or not to annotate each of the four functions in the program, yielding  $2^4 = 16$  potential typing *configurations*. Perhaps surprisingly, the choice of which functions to annotate has a significant and non-monotonic impact on the performance of the

<sup>1</sup>[https://github.com/python/performance/blob/master/performance/benchmarks/bm\\_spectral\\_norm.py](https://github.com/python/performance/blob/master/performance/benchmarks/bm_spectral_norm.py)

```

50 def bench_spectral_norm(loops):
51     range_it = xrange(loops)
52     for _ in range_it:                                      $l_3$ 
53         u = [1.0] * DEFAULT_N
54         for x in xrange(10):                                $l_4$ 
55             part_A_times_u((x,u))
56             part_At_times_u((x,u))
57 def eval_A(i, j):
58     return 1.0 / ((i + j) * (i + j + 1) // 2 + i + 1)
59 def part_A_times_u(i_u):
60     i, u = i_u
61     partial_sum = 0.0
62     for j, u_j in enumerate(u):                            $l_1$ 
63         partial_sum += eval_A(i, j) * u_j
64     return partial_sum
65 def part_At_times_u(i_u):
66     i, u = i_u
67     partial_sum = 0.0
68     for j, u_j in enumerate(u):                            $l_2$ 
69         partial_sum += eval_A(j, i) * u_j
70     return partial_sum

```



(a) A program adapted from Python Benchmark Suite. The functions `part_A_times_u` and `part_At_times_u` are identical except for the underlined parts. The loop labels  $l_1$ ,  $l_2$ ,  $l_3$ , and  $l_4$  are referenced in the cost lattice in Figure 2.

(b) Each node indicates whether the `eval_A`, `part_A_times_u`, and `part_At_times_u` functions are annotated (filled oval) or not (unfilled).

Fig. 1. A Python program (left) and its performance lattice (right)

program. In Figure 1(b), we illustrate this with a lattice that shows the execution time<sup>2</sup> in Reticulated for each of the 8 configurations where `bench_spectral_norm` is left unannotated. Each node in the lattice indicates whether the functions `eval_A`, `part_A_times_u`, and `part_At_times_u` are annotated (filled) or not (unfilled). For example, the node ① represents the configuration where no functions are annotated, while node ② represents the configuration where only `eval_A` is annotated.

Figure 1(b) shows that as we migrate the program to be more static (move up the lattice), the change in performance is unpredictable. For example, following the path ① → ② → ⑤ → ⑧, we see the performance first decreases, then increases twice. On the other hand, performance in ① → ④ → ⑦ → ⑧ first increases twice, then decreases. Additionally, the execution times at different configurations are very different, for example, the execution time at ② is about 3 times more than at ⑦. Together, these phenomena make it difficult for programmers to reason about what configurations lead to acceptable performance. Moreover, sometimes programmers are presented with a dilemma of whether static type checking is worth the performance slowdown to their application.

The unpredictable and sometimes severely negative performance impact of adding type annotations makes programmers reluctant to add them, and so the potential benefits of gradual typing go unrealized. This has lead Takikawa et al. to raise the question, “Is sound gradual typing dead?” [Takikawa et al. 2016]. Of course, adding type annotations can also improve performance.

<sup>2</sup>All times in this paper are in seconds and are measured on a laptop with 4 GB of RAM and an AMD A6-3400M quad-core processor running 64 bit Fedora 23.

In Figure 1(b), the performance at node ⑦ is much better than the original at node ①, in addition to providing increased safety from static checking. Similarly, while annotating one subset of modules led to the 100 times slowdown observed by Takikawa et al. in Typed Racket, annotating a different subset of modules improved performance [Takikawa et al. 2016]. The crux of the problem is that it's not clear in advance what annotations the programmer should or should not add in order to both increase static type safety and maximize performance.

## 1.2 Program Migration Scenarios

Programmers are reluctant to use gradual typing features due to the difficulty of predicting the performance impact of type annotations, as described in Section 1.1. The goal of our work is to remove this barrier by providing tooling that helps programmers understand and reason about the tradeoffs between the safety guarantees given by increased static type checking and performance during program migrations. In this subsection we enumerate four scenarios that such tooling should support.

(S1) *Maximizing static typing.* In this scenario, the programmer's primary goal is to maximize the amount of static type checking, while performance is a secondary concern. Maximizing static checking typically entails adding as many type annotations as possible. However, often the most static migration of a gradually typed program is not unique, and so the programmer wants to pick the most performant migration amongst the set of possibilities. The following program, adapted from [Campora et al. 2018], illustrates the non-uniqueness of most-static migrations. In this example, a type annotation may be added to the parameter `fixed` or to `widthFunc`, but not to both.

```
def width(fixed, widthFunc):
  if (fixed):
    widthFunc(fixed)
  else:
    widthFunc(5)
```

Campora et al. [2018] reported hundreds of different ways to maximize static typing in larger programs. Since each of the most-static migrations may have significantly different performance profiles, it is important that the programmer can make a rational selection among them.

(S2) *Maximizing performance.* In this scenario, the primary goal is to maximize performance, while increasing static type checking is a secondary concern. Therefore, the programmer needs support locating places to add type annotations that will either improve or at least not degrade performance. For example, starting from configuration ① in Figure 1, our tool should recommend configuration ⑦, since it has the least running time. This scenario can be extended in two ways: First, the programmer may want to maximize performance while providing type annotations in key places where increased static checking is judged to be more important or beneficial. For example, assume the function `eval_A` in Figure 1 must be annotated. Then the tool should consider configurations ②, ⑤, ⑥, and ⑧, and suggest configuration ⑧ as the most performant. Second, the programmer may want to maximize performance while restricting the number of proposed type annotations in order to manage the migration in a more incremental way. For example, starting from ①, if the programmer wants to add only a single annotation, then configurations ③ and ④ are preferable to ②.

(S3) *Increasing static type information without sacrificing performance.* In this scenario, the programmer wants to increase the static checking present in the program, but only if it does not decrease its performance. To support this scenario, the tool should be able to identify type annotations that can be added that do not incur a performance overhead. For example, starting from configuration ①, the tool might recommend migrations ③, ④, ⑦, or ⑧, all of which increase the amount of static typing without sacrificing performance. Additionally, if a previous migration hurt performance, the

148 tool should be able to identify subsequent migrations to improve it again. For example, migrating  
 149 from ① to ② significantly decreases performance, but the tool should be able to recommend ⑧ as  
 150 migration that will further increase safety guarantees from increased static typing while restoring  
 151 performance to (better than) previous levels.

152 (S4) *Explaining performance degradation.* In this scenario, the programmer has experienced a per-  
 153 formance degradation after adding type annotations and wants to understand why. Or, alternatively,  
 154 the programmer wants to understand why the tool recommends against a particular migration.  
 155 To support this, the tool should not only identify which program migrations will perform poorly,  
 156 but provide an *explanation* of why this performance degradation occurs. For example, when the  
 157 tool recommends against a migration from ① to ②, it might also explain that ② contains expensive  
 158 type casts in a deeply nested loop. Such explanations help the programmer to develop their own  
 159 mental model of how gradual typing affects the performance of their program. This enables them  
 160 to use gradual typing features more effectively, and to make more informed program migrations.

161 In each of these scenarios, we assume the programmer wants to make decisions with respect  
 162 to performance without decreasing the amount of typing. That is, we assume the tool will not  
 163 recommend the removal of type annotations. The approach described in this paper follows this  
 164 assumption, but extending it to also support the removal of type annotations poses no fundamental  
 165 difficulties.

### 166 1.3 Capabilities of a Tool to Support Program Migration

168 Each of the scenarios in Section 1.2 involve exploring  
 169 many alternative configurations of a gradually typed  
 170 program. Without tool support, this is extremely te-  
 171 dious since it requires manually adding and removing  
 172 type annotations and rerunning the program to mea-  
 173 sure its performance. Worse, this exploration cannot  
 174 hope to be complete for large programs since the search  
 175 space is simply too large, so the best configuration for  
 176 the scenario will not likely be found. Therefore, tool  
 177 support is necessary to support the scenarios and to  
 178 help programmers effectively migrate gradually typed  
 179 programs. This subsection identifies the key capabilities  
 180 needed to build such a tool, and outlines the techniques  
 181 we use to provide them.

182 We propose a methodology for systematically explor-  
 183 ing the entire space of potential program migrations needed to support each scenario and to  
 184 efficiently identify the most performant type configurations in that space. We can break this  
 185 methodology down into three fundamental capabilities: (C1) A way to enumerate and efficiently  
 186 represent all valid type configurations of a gradually typed program. (C2) A way to statically  
 187 approximate and compare the runtime performance of a single configuration of a gradually typed  
 188 program. (C3) A way to combine C1 and C2 to efficiently compute and compare the performance  
 189 approximations for all type configurations.

190 Regarding C1, in a program with  $n$  parameters, an upper limit on the number of possible type  
 191 configurations is  $2^n$  since each parameter can either be assigned a static type or remain unannotated  
 192 (and thus dynamically typed). However, in general, not every combination of parameters can be  
 193 statically annotated in a consistent way, as illustrated by the width example in Section 1.2. We can  
 194 use *variational type inference* [Chen et al. 2014] to efficiently explore all  $2^n$  possibilities by inferring  
 195 static types for all parameters in one pass while keeping track of which combinations of types  
 196

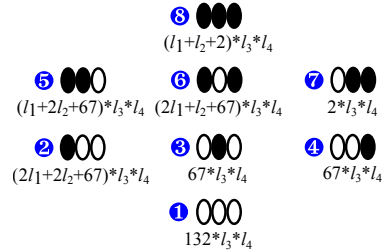


Fig. 2. Cost lattice for the program in Figure 1. We have omitted an addend 2 that is shared by all configurations. The letters  $l_1$ ,  $l_2$ ,  $l_3$ , and  $l_4$ , introduced in Figure 1(a), represent the number of iterations for the respective for loops.

are compatible with each other. The idea to use type inference to help migrate gradually typed programs is inspired by the observation that “static type systems accommodate common untyped programming idioms” by Takikawa et al. [2016], and by previous successes combining gradual typing and type inference [Campora et al. 2018; Garcia and Cimini 2015; Rastogi et al. 2012; Siek and Vachharajani 2008]. In particular, we reuse the machinery developed in Campora et al. [2018] to transform the output of variational type inference into an efficient representation of all valid type configurations of a program.

As in previous work, the success of type inference in a gradually typed setting can be expected to vary significantly across programs. In our benchmarks, we observed a broad range of outcomes, successfully inferring types for anywhere from 25% to 100% of parameters in a program (Section 6.2). Fortunately, effectively supporting the migration scenarios does not require inferring types for all parameters. We simply infer types for as many parameters as possible, then reason about this space of migrations. Subsequent migrations may require fundamental changes to the code to remove behavior that relies on dynamic typing, expanding the space we can reason about.

Capability C2 requires a way to estimate the overhead of gradual typing in each configuration, allowing us to estimate and rank their expected runtime performance. To enable C2, we develop a static *cost semantics* for gradually typed programs. The insight underlying our cost semantics is that the overhead of gradual typing is mostly caused by inserted casts [Takikawa et al. 2016] and checks [Vitousek et al. 2014, 2017]. Therefore, we statically approximate the number and complexity of cast and check operations that will be performed while executing a gradually typed program. Figure 2 shows the result of applying our cost semantics to each of the configurations of the program in Figure 1. Since we do not know statically how many times each loop body in the program will be executed, the costs are parameterized by symbolic values representing the number of iterations ( $l_1$ ,  $l_2$ ,  $l_3$ , and  $l_4$ ).

Note that our cost semantics does not approximate the absolute runtime of different configurations of the program, but only the overhead of gradual typing. This is similar to other cost analyses focused on specific aspects of program execution [Hoffmann and Hofmann 2010; Hoffmann and Shao 2015]. For example, the absence of loop labels  $l_1$  and  $l_2$  in ❶ does not suggest that these loops are not executed, but rather that no casts or checks will be performed in these loops in the corresponding configuration. The factors 2, 67, and 132 that multiply these loop labels correspond to the estimated overhead of individual casts and checks performed in the body of these loops.

Since gradual typing overhead is only a (smaller or larger) fraction of the running time of a program, we can see that the ratios of estimated costs for two configurations do not correspond precisely to the ratios of their running time. For example, the cost of ❶ in Figure 2 is 66 times that of ❷, while the running time of ❶ in Figure 1(b) is roughly 2.5 times that of ❷. However, the relative ordering of the costs in Figure 2 do correspond to the relative ordering of running times in Figure 2. For example, the cost at ❷ is  $2 * l_3 * l_4$  is the lowest estimated cost and corresponds to the lowest running time at ❷. Similarly, the cost decreases along the path ❶  $\rightarrow$  ❹  $\rightarrow$  ❸, just as the running time decreases along the path ❶  $\rightarrow$  ❹  $\rightarrow$  ❸. This illustrates that the estimated costs are useful as a tool for predicting the relative performance of different typing configurations. This makes sense since the different configurations differ only in their type assignments, and so any difference in their running time should be explained by the overhead of gradual typing.

One may ask how can we compare, for example, the costs of ❶ and ❸ since one mentions  $l_1$  and  $l_2$  and the other does not? The answer is: we can not. This makes sense because the relation of the running times between these two variants is not fixed. In this case, it depends on the value of DEFAULT\_N (see Figure 1(a)), which determines how many times the loop bodies of  $l_1$  and  $l_2$  are executed. In general, different loop bodies can be expected to execute different numbers of times based on different inputs and environment settings. The loop labels make this uncertainty explicit.

246 While costs involving different loop labels cannot be directly compared, they can still be used to  
247 produce explanations to help a programmer make an informed decision when deciding between  
248 different migrations. For example, the tool might explain that the configuration at ⑧ induces casts  
249 in the loops  $l_1$  and  $l_2$ , but reduces the cost of casts in the loops  $l_3$  and  $l_4$ , relative to configuration ①.

250 With capabilities C1 and C2, we can statically enumerate all type configurations of a program  
251 and statically compare the performance of different configurations using our cost semantics.  
252 Hypothetically, we could apply these together in order to identify the best configurations to support  
253 our program migration scenarios. The problem is that, since the number of type configurations  
254 produced by C1 scales exponentially with the number of parameters, the search space quickly  
255 grows overwhelming to perform the search by applying C2 directly. This problem is solved by C3,  
256 which enables efficiently computing and comparing costs for all valid type configurations.

257 The key observation to support C3 is that a substantial amount of work can be reused between  
258 the cost analyses of different configurations. For example, calculating the cost of configurations ①  
259 and ③ can share the computations of costs for `eval_A` and `part_At_times_u`. Similarly, calculating  
260 the costs of ③ and ④ can share the computations of `eval_A`. In large programs, the majority of the  
261 computations can be shared when computing the costs of two similar configurations.

262 Unfortunately, sharing cannot be achieved by simply analyzing the costs of different functions  
263 separately then composing the results since interactions between functions do affect the analyses.  
264 However, by locally capturing differences between sets of configurations and preserving these  
265 differences throughout the analysis, we can effectively reuse results wherever possible. Specifically,  
266 we apply the ideas of *variational programming* [Chen et al. 2016; Erwig and Walkingshaw 2013]  
267 and *variational typing* [Chen et al. 2014] to systematically reuse computations during the cost  
268 analysis. Instead of enumerating all configurations and computing the cost of each separately, we  
269 perform a *variational cost analysis* that analyzes the program once to compute a *variational cost*  
270 that compactly represents the cost of all valid type configurations.

271 With these three capabilities, we can support all of the scenarios outlined in Section 1.3 by  
272 performing a variational cost analysis, then querying the variational cost to identify the desired  
273 configuration. For example, to support S1 (maximizing static typing), we would select the lowest  
274 cost among the configurations that include as many annotations as possible.

#### 275 1.4 Relation with Previous Work and Contributions of this Work

276 There have been several lines of research addressing the performance problem of gradual typing  
277 since Takikawa et al.'s [2016] report on the prohibitive overhead of sound gradual typing. Previous  
278 work has addressed the problem through new languages with more efficient gradually typed  
279 semantics [Muehlboeck and Tate 2017; Vitousek et al. 2017] and through new implementation  
280 techniques for existing languages [Bauman et al. 2017; Richards et al. 2017]. There are two main  
281 differences between our work and previous efforts: First, our approach does not require changes to  
282 existing gradually typed languages or implementations; it works with the prevailing implementation  
283 technique of translating a typed variant of a language into an underlying untyped language. Second,  
284 in addition to addressing the common goal of reducing or avoiding performance problems, our  
285 work also provides a way to understand and debug performance problems when they occur. We  
286 discuss the relation with existing work in more detail in Section 7.2.

287 By drawing insights from gradual typing and type inference [Campora et al. 2018; Garcia  
288 and Cimini 2015; Siek and Vachharajani 2008], cost analysis [Danner et al. 2015; Hoffmann and  
289 Hofmann 2010], and variational programming [Chen et al. 2012, 2014], we develop a methodology  
290 for understanding, debugging, and optimizing the performance of gradual programs based on a  
291 deep understanding of how types affect performance. To test the feasibility of this methodology,  
292 we have implemented our variational cost analysis as HERDER, a tool that efficiently and accurately  
293

294

analyzes the costs of many type configurations of a Reticulated Python program. Overall, this paper makes the following contributions:

- (1) We develop a cost semantics for casts in gradually typed programs with a guarded semantics in Section 4. The cost semantics is simple and enables automating cost analysis, yet still allows authentically comparing the relative run times of different configurations for the same program.
- (2) We combine variations and cost analysis, yielding a variational cost analysis in Section 5. Instead of computing costs for all configurations separately, variational cost analysis systematically reuses computations to compute a variational cost that encodes the cost of all the configurations that can be inferred.
- (3) We have implemented our approach as HERDER on top of Reticulated and evaluate it in Section 6. We evaluate the accuracy of the cost semantics by taking the configuration HERDER reports as having the lowest cost and testing whether it has the fastest runtime amongst the measured configurations. Our evaluation demonstrates that HERDER can efficiently find configurations yielding good performance. In most benchmarks, the recommended configuration is one of the top 3 in terms of execution time. Moreover, our approach is scalable, taking exponentially less time than a brute-force approach as the number of configurations becomes large, and 2–4 times as long as the cost analysis of a single configuration.

The rest of the paper is organized as follows. Section 2 provides necessary background on gradual typing and variational typing. Section 3 informally introduces our static cost semantics, while Section 4 gives the formal definition. Section 5 gives the formal definition of the full variational cost semantics. The implementation of HERDER is described in Section 6, together with an evaluation of its accuracy and scalability. Section 7 describes related work and Section 8 concludes.

## 2 TYPING, GRADUALLY AND VARIATIONALLY

This section provides background needed to understand the rest of the paper. In Section 2.1, we describe why and where casts are inserted into gradually typed programs. In Section 2.2, we review variational typing [Chen et al. 2012] as a way to efficiently analyze many variants of a program.

### 2.1 Gradual Typing

Gradual typing allows mixing dynamically typed and statically typed code within a single program. In a gradually typed program, statically typed values can flow into dynamically typed code and vice versa. For example, consider the following dynamically typed function `double`.

```
def double(x):
    return x * 2
```

Suppose the multiplication operation `*` is statically typed as  $\text{Int} \rightarrow \text{Int} \rightarrow \text{Int}$ . Then the dynamically typed argument `x` flows into this statically typed operation. Similarly, if we invoke this function as `double(3)`, where `3` has static type `Int`, then a statically typed value flows into the dynamically typed function.

During static type checking, the interface between static and dynamic code is defined by a *consistency* relation [Garcia et al. 2016; Siek and Taha 2006]. Consistency (denoted by  $\sim$ ) weakens type equality by making every type consistent with the type `Dyn`, representing the type of dynamically typed code. So the expression `x * 2` in the `double` function is statically type correct since `x` has type `Dyn` and  $\text{Dyn} \sim \text{Int}$ . Of course, in order to preserve the dynamic type safety of gradually typed programs, additional type checking may have to be performed at runtime. For example, the `double` function must be translated to a version with an explicit *type cast*  $[\text{Int} \Leftarrow \text{Dyn}]$  as shown below.

```
def double(x):
    return [Int ← Dyn]x * 2
```

344 During program translation, such casts are inserted into the program wherever dynamically and  
 345 statically typed code interact. Note that we do not need to cast the value 2 to `Int` since its type  
 346 is statically known. Similarly, if we annotated the type of `double` to be `Int → Int`, we would not  
 347 need to cast `x` to `Int`, and the function call `double(3)` would not involve any casts since all types  
 348 are statically known.

349

## 350 2.2 Variational Typing

351 Variational typing was developed by Chen et al. [2014] to provide types for *variational programs*. A  
 352 variational program represents several related program variants using *choices* [Erwig and Walking-  
 353 shaw 2011] to denote where the variants differ, as illustrated below.

```
354 result = B⟨odd, double⟩(3) (vprog)
```

355 The variational program `vprog` contains a choice named `B` with *alternatives* `odd` and `double`. Two  
 356 distinct programs can be generated by *selecting* the first or second alternative of choice `B`. For  
 357 example, selecting the first alternative, denoted  $\lfloor \text{vprog} \rfloor_{B.1}$ , yields the program `result = odd(3)`,  
 358 while selecting  $\lfloor \text{vprog} \rfloor_{B.2}$  yields `result = double(3)`.

359 Choices with the same name in a variational program are synchronized, while choices with  
 360 different names are independent. That is,  $\lfloor e \rfloor_{d.i}$  selects the  $i$ th alternative of all choices named  
 361  $d$  in  $e$ . We call  $d.i$  a *selector* and range over selectors with  $s$ . Obtaining a plain (non-variational)  
 362 program from a variational program may require several selections. We call a set of selectors a  
 363 *decision*, ranged over by  $\delta$ , and generalize the notation of selection to decisions as  $\lfloor e \rfloor_\delta$ .

364 Variation in expressions naturally gives rise to variation in types. For example, the expression  
 365  $B\langle \text{odd}, \text{double} \rangle$  can be assigned the type  $B\langle \text{Int} \rightarrow \text{Bool}, \text{Int} \rightarrow \text{Int} \rangle$ . Variational type systems extend  
 366 traditional type systems to accommodate choices at the expression and type level. Typing function  
 367 applications is complicated by the fact that both the function and argument may be variational. To  
 368 accommodate this, variational type systems are equipped with a type equivalence relation. Two  
 369 types  $T_1$  and  $T_2$  are equivalent, denoted  $T_1 \equiv T_2$  if  $\lfloor T_1 \rfloor_\delta = \lfloor T_2 \rfloor_\delta$  for every decision  $\delta$ . So, for example,  
 370  $B\langle \text{Int} \rightarrow \text{Bool}, \text{Int} \rightarrow \text{Int} \rangle \equiv \text{Int} \rightarrow B\langle \text{Bool}, \text{Int} \rangle$  since both sides of the equivalence yield  $\text{Int} \rightarrow \text{Bool}$   
 371 when selecting `B.1` and both sides yield  $\text{Int} \rightarrow \text{Int}$  when selecting `B.2`. Taking the right-hand side  
 372 of this equivalence, it is easy to see that the function application  $B\langle \text{odd}, \text{double} \rangle(3)$  is well typed  
 373 and yields a result of type  $B\langle \text{Bool}, \text{Int} \rangle$ .

374 A crucial property of variational type systems is that *selection preserves typing*. That is, if  $e : T$   
 375 then  $\forall s. \lfloor e \rfloor_s : \lfloor T \rfloor_s$ . Previous work showed how variational types enable efficiently reasoning  
 376 about all possible assignments of dynamic or static types to function parameters in gradually  
 377 typed programs [Campora et al. 2018]. This enables efficiently migrating between gradually typed  
 378 programs with different type assignments. In this work, we tackle the problem of estimating the  
 379 performance overhead associated with different migrations for a gradual program.

380

## 381 3 THE WORKFLOW OF HERDER

382 HERDER works by generating and reasoning about variational programs that represent all possible  
 383 type configurations at once. We use the following program to illustrate how this works, including  
 384 how variational casts are inserted and how variational costs are computed. We assume the `+` operator  
 385 has the static type `Int → Int → Int`.

```
387 def add(x, y):  
388     return x + y
```

389 Casts are only inserted when passing dynamically typed values to statically typed code, which can  
 390 occur when some parts of the program contain type annotations, when primitive operations are  
 391 assigned static types by the language implementation, or when typed external code is called. We

392



focus on how programmer-added type annotations change the behavior of cast insertion. Therefore, before we can reason about the costs of inserted casts arising via different annotations, we need to infer what combinations of type annotations can be added to the program.

Migrational typing by Campora et al. [2018] can efficiently infer types for all configurations of a program in a gradually typed language with type inference. It is only necessary for type inference to be sound (not complete), so migrational typing can be applied even to languages with features, such as subtyping, that prevent complete type inference. For add, we can use migrational typing to infer that each parameter can independently have type `Dyn` or `Int` (note that an unannotated parameter is equivalent to one annotated by `Dyn`). This yields four potential configurations of `add`, which we can represent in a single variational program, `addV`, with two independent choices.

```
def addV(x: B⟨Dyn, Int⟩, y: D⟨Dyn, Int⟩):
  return x + y
```

In this paper, we focus on the problem of reasoning about and comparing costs for different configurations and reuse necessary machinery from Campora et al. [2018] to get type information.

Instead of generating all configurations of `addV` and separately reasoning about the casts in each, we instead add *variational casts* to `addV` and reason about `addV` directly. In this case, we need to insert variational casts to ensure that the arguments to `+` have the type `Int` in each configuration. Specifically, we insert the cast  $[Int \Leftarrow B\langle Dyn, Int \rangle] x$ , and similarly for `y`. Conceptually, this represents the cast  $[Int \Leftarrow Dyn]x$  in `B.1`, where `x` has static type `Dyn`, and it represents the cast  $[Int \Leftarrow Int]x$  in `B.2`, where `x` has static type `Int`. Since  $[Int \Leftarrow Int]$  is a no-op, we can transform the variational cast applied to `x` to  $B\langle [Int \Leftarrow Dyn], \epsilon \rangle x$ , making it clear that no cast ( $\epsilon$ ) is performed in the `B.2` case where passing `x` to `+` can be statically type checked.

After cast insertion and simplification, we obtain `addVC`:

```
def addVC(x: B⟨Dyn, Int⟩, y: D⟨Dyn, Int⟩):
  return B⟨[Int \Leftarrow Dyn], \epsilon⟩ x + D⟨[Int \Leftarrow Dyn], \epsilon⟩ y
```

Our next goal is to find a way to measure the overhead of the inserted casts, and more importantly to compare the overhead of different configurations. For this simple program, we can simply count the number of inserted casts. Therefore, we assign the *variational cost*  $B\langle 1, 0 \rangle$  to the cast  $B\langle [Int \Leftarrow Dyn], \epsilon \rangle x$ , and the cost  $D\langle 1, 0 \rangle$  for casting `y`. The cost for `addV` is then  $B\langle 1, 0 \rangle + D\langle 1, 0 \rangle$ . Applying standard variational programming techniques [Erwig and Walkingshaw 2013], we can reduce this cost to  $B\langle D\langle 2, 1 \rangle, D\langle 1, 0 \rangle \rangle$ , which captures the costs of all four configurations of the program.

We can obtain the cost of each configuration by selecting from the variational cost with the corresponding decision. For example, selecting with  $\delta = \{B.1, D.2\}$  yields 1, corresponding to the single cast of `x` in the configuration of `addV` produced by selecting with the same decision  $\delta$ . From the variational cost, we can see that the configuration of `add` that annotates both parameters with `Int` leads to the lowest cost.

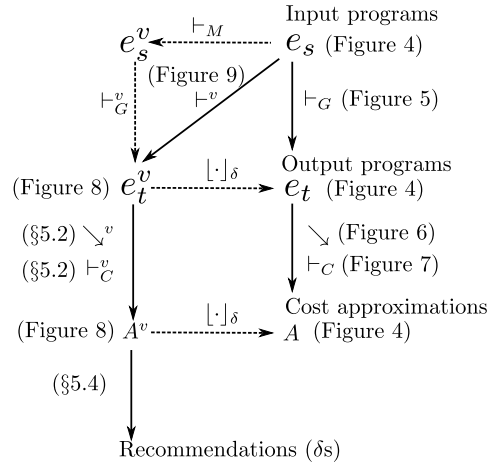


Fig. 3. Overview of computing costs for all valid type configurations. The operations and transformations attached to dashed arrows are not presented in this paper, they are either from previous work (such as  $\vdash_M$  from Campora et al. [2018] and  $\lfloor \cdot \rfloor_\delta$  from Chen et al. [2014]) or for aiding conceptual understanding ( $\vdash_C^v$ ).

Sections 4 and 5 present a formal treatment of this process. To separate concerns, we first present the cost semantics for a single program in Section 4, then a method that computes the costs for all valid configurations for the given program in Section 5.

To give a high-level view of the formalization, we present the connections between various relations and syntaxes in Figure 3. Moving down the right side of the figure illustrates the process of assigning costs to individual plain programs, while moving down the left side of the figure illustrates the variational cost analysis process. Plain cost analysis works as follows (starting from the upper right corner and moving down): the syntax  $e_s$  represents input gradual programs, which may contain type annotations; the type-directed transformation  $\vdash_G$  inserts casts into  $e_s$ , erases type annotations, and generates  $e_t$ , which can be directly executed on the underlying interpreter; we then use  $\searrow$  to further transform programs in  $e_t$  into programs in the same syntax but which are more amenable to cost analysis, and compute costs  $A$  with the  $\vdash_C$  relation.

Conceptually, variational cost analysis works as follows (starting from the upper right corner, moving left, then down): starting from  $e_s$ , we apply migrational typing,  $\vdash_M$ , [Campora et al. 2018] to compute all valid type configurations, which we encode as  $e_s^v$ ; we apply a variational transformation,  $\vdash_G^v$ , to insert variational casts yielding a variational version of the target language  $e_t^v$ ; a variational transformation  $\searrow^v$  makes the program more amenable to cost analysis; finally, we compute variational costs  $A^v$  with  $\vdash_C^v$ . In practice, however, we move directly from  $e_s$  to  $e_t^v$  using  $\vdash^v$  since this simplifies the formalization. Since  $A^v$  encodes costs for all configurations, we can use this result to make recommendations to satisfy the current program migration scenario. The arrows annotated by  $[\cdot]_\delta$  help establish the correctness of variational cost analysis by relating variational results with their corresponding plain results through selection (see Theorem 5.2).

## 4 COSTS FOR A SINGLE CONFIGURATION

Casts are the major source of performance degradation in gradually typed programs [Takikawa et al. 2016]. Therefore, we need a way to accurately estimate the costs associated with the inserted casts. In this section, we address this need by developing a static cost semantics in the style of Danner et al. [2015]. Our cost semantics produces an expression in a *cost language* that evaluates to a cost estimate for the corresponding gradually typed program. Specifically, in Section 4.1, we define a function for estimating the cost of an individual cast. In Section 4.2, we introduce the cost language and cost semantics informally through an example, then give a formal treatment in Sections 4.3 through 4.5. In Section 4.6, we discuss important properties of the cost semantics.

### 4.1 The Cost of a Basic Cast

First, we consider the cost of an individual cast  $\lceil G_1 \Leftarrow G_2 \rceil$ , which checks that a value with gradual type  $G_2$  can be converted into a value with gradual type  $G_1$ , and performs the conversion if necessary. The easiest cast to reason about is casts to  $\text{Dyn}$ . When  $G_1 = \text{Dyn}$  and  $G_2$  is a base type (such as  $\text{Int}$  or  $\text{Bool}$ ) the cost of the cast is  $b$ , representing the cost of boxing a value.

The simplest cast that does some work is a cast where  $G_2 = \text{Dyn}$  and  $G_1$  is a base type. In this case, a dynamic type check must be performed. We use  $c$  to represent the constant cost of such basic casts. For casts involving list types in both the source and target, we use  $s$  to represent the overhead of checking the elements of a list, and we recursively compute the cost for the types inside the respective list constructors. Our cost function is below:

$$\text{cost}(\lceil G_1 \Leftarrow G_2 \rceil) = \begin{cases} b & G_1 = \text{Dyn} \\ s \cdot \text{cost}(\lceil G'_1 \Leftarrow G'_2 \rceil) & G_i = [G'_i] \\ c & \text{otherwise} \end{cases}$$

The *cost* function handles basic casts, which are simply performed at the locations they are encountered in the program. In contrast, casts involving function types or reference types are trickier since they are not applied immediately, but rather when the corresponding function or reference is used [Siek et al. 2009]. We illustrate how to account for the cost of such casts and give a precise cost calculation in Section 4.4.

## 4.2 Estimating Cast Costs in Programs

Armed with a cost function for individual casts, we can estimate the cost of a sequence of statements by simply summing up the costs of each. However, more interesting programs present two challenges: (1) How do we represent the costs of functions, whose bodies may contain casts whose individual costs depend on the types of the arguments they are applied to? (2) How do we estimate the cost of loops, whose bodies may contain casts that are executed a statically indeterminate number of times? In this subsection, we introduce a cost language and cost semantics for addressing these challenges. We use the following function `mult` as a running example.

```
def mult(md, mr):
    sum = 0
    for i in range(md):
        sum = add(sum, mr)
    return sum
```

This function multiplies the multiplicand `md` with the multiplier `mr` by iteratively adding `mr` to a local variable `sum`. The built-in Python function `range : Int → [Int]` returns a list  $[1, 2, \dots, n]$  for the given  $n$ , and the helper function `add` returns the sum of its arguments.

First, we consider how to represent the cast cost of a function. A function by itself incurs no cost, but rather represents a *potential* cost when invoked. Moreover, the cast cost of the body of a function will vary depending on the arguments that are passed in. The potential cost of a function can be represented by a corresponding function in the cost language, and the cost of a function application can be approximated by executing a corresponding application in the cost language.

More concretely, following Danner et al. [2015], we represent the cost *approximation* of a term by a pair  $(C, P)$ , where  $C$  represents the *immediate cost* of the term and  $P$  represents its *potential cost*. We use  $A$  to range over approximations and use *cost* and *potential* to refer to immediate costs and potential costs, respectively. The potential of a function abstraction  $\lambda x.e$  in the source language is captured by a *potential abstraction* in the cost language of the form  $\Lambda x.A$ . The cost of a function application  $e_1 e_2$  can then be computed by applying the corresponding potentials of  $e_1$  and  $e_2$ .

Returning to our example, we can sketch a template for the approximation of `mult` as  $(0, \Lambda md.(0, \Lambda mr.A_m))$ , where the function itself (and its partial application) has no immediate cost, and the body of the function is approximated by  $A_m$ , which may refer to the potentials of its parameters, `md` and `mr`.

While we produce costs, we maintain a *cost environment* that maps each source language variable to its potential. Let us assume that `sum` has type `ref Dyn`. The first statement of the body is `sum = 0`, which has approximation  $(0, 0)$ . The cost is 0 since we generate no casts and 0 is a constant, and its potential is 0 since the statement does not involve functions or loops. In reality, `sum` is a reference and has associated potential. To handle this, we define a transformation in Section 4.4, but for simplicity while illustrating this example we will assume that reference creation and assignment contributes no overhead relating to proxies. Therefore, the statement contributes no cost to  $A_m$  and the environment is extended with `sum`  $\mapsto 0$ .

We now turn to producing a cost for the loop in `mult`. For the subexpression `range(md)`, suppose the cost environment maps the built-in function `range` to the potential  $\Lambda u.(0, u)$ , indicating that it incurs no cast costs except the potential costs of its argument. Since `md` has type `Dyn`, a cast  $[Int \leftarrow$

| Term variables | $x, y, z$ | Base types                                      | $\gamma$ | Cost variables | $x, y, z$    | Loop labels  | $l$ |
|----------------|-----------|---|----------|----------------|--------------|--|-----|
| Source         | $e_s ::=$ | $x \mid \lambda x: G. e_s \mid e_s e_s$         |          | Type constr.   | $T ::=$      | $[] \mid \text{ref}$                                   |     |
| expr.          |           | $\text{for } x \text{ in } e_s \text{ do } e_s$ |          | Gradual types  | $G ::=$      | $\gamma \mid T G \mid G \rightarrow G \mid \text{Dyn}$ |     |
|                |           | $\text{let } x = e_s \text{ in } e_s$           |          | Approximations | $A ::=$      | $(C, P) \mid C \oplus P$                               |     |
|                |           | $\text{letrec } x = e_s \text{ in } e_s$        |          | Costs          | $C ::=$      | $n \mid l \cdot C \mid C + C$                          |     |
|                |           | $\text{ref } e_s \mid !e_s \mid e_s := e_s$     |          | Potentials     | $P ::=$      | $n \mid x \mid l \cdot P \mid P + P$                   |     |
| Target         | $e_t ::=$ | $x \mid \lambda x. e_t \mid e_t e_t$            |          |                |              | $\mid \Lambda x. A \mid P P$                           |     |
| expr.          |           | $\text{for } x \text{ in } e_t \text{ do } e_t$ |          | Type env.      | $\Gamma ::=$ | $\emptyset \mid \Gamma, x \mapsto G$                   |     |
|                |           | $\text{let } x = e_t \text{ in } e_t$           |          | Cost env.      | $\Phi ::=$   | $\emptyset \mid \Phi, x \mapsto P$                     |     |
|                |           | $\text{ref } e_t \mid !e_t \mid e_t := e_t$     |          |                |              |  |     |
|                |           | $[G \leftarrow G]e_t$                           |          |                |              |  |     |

Fig. 4. Syntax for our gradually typed functional language and its corresponding cost language.

$\text{Dyn}]$  must be performed before passing  $\text{md}$  into  $\text{range}$ . This cast has cost  $c$ . The overall approximation of  $\text{range}(\text{md})$  is this immediate cost added to the approximation returned by applying the potential of  $\text{range}$  to the potential of the argument  $\text{md}$ , which is  $(\Lambda u.(0, u)) \text{ md} = [\text{md}/u](0, u) = (0, \text{md})$ . This yields a final approximation of  $(c, \text{md})$  for the subexpression  $\text{range}(\text{md})$ .

The cost of a loop is the cost of its body times the number of iterations. In general, the number of iterations is unknown statically. Therefore, in our approximations we introduce a unique symbolic value to stand for the number of iterations each loop executes, which we call a *loop label*. The cost approximation of a loop is the cost and potential of its body, each multiplied by the loop label.

In  $\text{mult}$ , the body of the loop is  $\text{sum} = \text{add}(\text{sum}, \text{mr})$ . For simplicity, let us assume for now that the assignment does not generate a reference cast and that the use of  $\text{sum}$  in the body (a dereference) does not generate a cast from a proxy. Now we can focus on the application of  $\text{add}$  to  $\text{sum}$  and  $\text{mr}$ . Suppose  $\text{add}$  is annotated as  $\text{Int} \rightarrow \text{Int} \rightarrow \text{Int}$ , then we must cast both arguments for an immediate cost of  $2c$ . Also suppose  $\text{add}$  has potential  $\Lambda x.(0, \Lambda y.(c, x + y))$ . As described above, the cost environment maps  $\text{sum} \mapsto 0$  and  $\text{mr} \mapsto \text{mr}$ , so the resulting approximation of the application is  $(c, 0 + \text{mr}) = (c, \text{mr})$ . Adding the immediate costs of the casts to this approximation yields an overall approximation of  $(3c, \text{mr})$  for the loop body. To approximate the overall cost of the loop, we multiply the cost of the body by a new loop label  $l$ , yielding  $(3c \cdot l, \text{mr} \cdot l)$ , then add the approximation of  $\text{range}(\text{md})$ , yielding  $(3c \cdot l + c, \text{mr} \cdot l + \text{md})$ .

Finally, the function returns a dereference of  $\text{sum}$ . This has cost  $(0, 0)$  since the return type of  $\text{mult}$  is  $\text{Dyn}$ . Since the loop is the only source of cast costs in  $\text{mult}$ , we can set  $A_m$  in our template above to the loop approximation to yield our final approximation for the whole function:  $(0, \Lambda \text{md}.(0, \Lambda \text{mr}.(3c \cdot l + c, \text{mr} \cdot l + \text{md})))$ .

### 4.3 Cast Insertion Rules

As the running example in Section 4.2 shows, our idea of computing costs is to measure the number and complexity of casts in programs. This subsection presents rules for inserting casts into gradually typed programs. In Figure 4, we define a simple calculus that captures the essential features of gradually typed programs with respect to cost analysis. The syntax of expressions is lambda calculus extended by references and a loop construct. The syntax of gradual types consists of base types ( $\gamma$ ), list types, function types, and the dynamic type. List types can be introduced through user type annotations or the initial type environment. The syntax of approximations, costs, and potentials

$$\begin{array}{c}
589 \\
590 \\
591 \\
592 \\
593 \\
594 \\
595 \\
596 \\
597 \\
598 \\
599 \\
600 \\
601 \\
602 \\
603 \\
604 \\
605 \\
606 \\
607 \\
608 \\
609 \\
610 \\
611 \\
612 \\
613 \\
614 \\
615 \\
616 \\
617 \\
618 \\
619 \\
620 \\
621 \\
622 \\
623 \\
624 \\
625 \\
626 \\
627 \\
628 \\
629 \\
630 \\
631 \\
632 \\
633 \\
634 \\
635 \\
636 \\
637
\end{array}$$

$$\begin{array}{c}
\text{VAR} \frac{x : G \in \Gamma}{\Gamma \vdash_G x \rightsquigarrow x : G} \quad \text{ABS} \frac{\Gamma, x \mapsto G \vdash_G e_s \rightsquigarrow e_t : G_1}{\Gamma \vdash_G \lambda x : G.e_s \rightsquigarrow \lambda x.e_t : G \rightarrow G_1} \\
\text{APP} \frac{\Gamma \vdash_G e_{1s} \rightsquigarrow e_{1t} : G \quad \Gamma \vdash_G e_{2s} \rightsquigarrow e_{2t} : G' \quad \text{dom}(G) \sim G'}{\Gamma \vdash_G e_{1s} e_{2s} \rightsquigarrow (\llbracket \text{dom}(G) \rightarrow \text{cod}(G) \Leftarrow G \rrbracket e_{1t} \llbracket \text{dom}(G) \Leftarrow G' \rrbracket e_{2t}) : \text{cod}(G)} \\
\text{FOR} \frac{\Gamma \vdash_G e_{1s} \rightsquigarrow e_{1t} : G_1 \quad \text{ext}_L(G_1) = G \quad \Gamma, x \mapsto G \vdash_G e_{2s} \rightsquigarrow e_{2t} : G_2 \quad l \text{ fresh}}{\Gamma \vdash_G \text{for } x \text{ in } e_{1s} \text{ do } e_{2s} \rightsquigarrow \text{for } x \text{ in } \llbracket [G] \Leftarrow G_1 \rrbracket e_{1t} \text{ do } e_{2t} : G_2} \\
\text{LET} \frac{\Gamma \vdash_G e_{1s} \rightsquigarrow e_{1t} : G_1 \quad \Gamma, x \mapsto G_1 \vdash_G e_{2s} \rightsquigarrow e_{2t} : G}{\Gamma \vdash_G \text{let } x = e_{1s} \text{ in } e_{2s} \rightsquigarrow \text{let } x = e_{1t} \text{ in } e_{2t} : G} \\
\text{LETREC} \frac{\Gamma, x \mapsto G_1 \vdash_G e_{1s} \rightsquigarrow e_{1t} : G_1 \quad \Gamma, x \mapsto G_1 \vdash_G e_{2s} \rightsquigarrow e_{2t} : G}{\Gamma \vdash_G \text{letrec } x = e_{1s} \text{ in } e_{2s} \rightsquigarrow \text{letrec } x = e_{1t} \text{ in } e_{2t} : G} \\
\text{REF} \frac{\Gamma \vdash_G e_s \rightsquigarrow e_t : G}{\Gamma \vdash_G \text{ref } e_s \rightsquigarrow \text{ref } e_t : \text{ref } G} \quad \text{DEREF} \frac{\Gamma \vdash_G e_s \rightsquigarrow e_t : G_1 \quad G = \text{ext}_R(G_1)}{\Gamma \vdash_G !e_s \rightsquigarrow !\llbracket \text{ref } G \Leftarrow G_1 \rrbracket e_t : G} \\
\text{ASSIGN} \frac{\Gamma \vdash_G e_{1s} \rightsquigarrow e_{1t} : G_1 \quad G = \text{ext}_R(G_1) \quad \Gamma \vdash_G e_{2s} \rightsquigarrow e_{2t} : G_2 \quad G \sim G_2}{\Gamma \vdash_G e_{1s} := e_{2s} \rightsquigarrow \llbracket \text{ref } G \Leftarrow G_1 \rrbracket e_{1t} := \llbracket G \Leftarrow G_2 \rrbracket e_{2t} : \text{ref } G} \\
\llbracket G \Leftarrow G \rrbracket e_t = e_t \quad \llbracket G_1 \Leftarrow G \rrbracket e_t = [G_1 \Leftarrow G] e_t \\
\text{dom}(G_1 \rightarrow G_2) = G_1 \quad \text{dom}(\text{Dyn}) = \text{Dyn} \quad \text{cod}(G_1 \rightarrow G_2) = G_2 \quad \text{cod}(\text{Dyn}) = \text{Dyn} \\
\text{ext}_L([G]) = G \quad \text{ext}_L(\text{Dyn}) = \text{Dyn} \quad \text{ext}_R(\text{ref } G) = G \quad \text{ext}_R(\text{Dyn}) = \text{Dyn}
\end{array}$$

Fig. 5. Cast insertion rules.

are as described in Section 4.2. The  $n$  in the syntax refers to  $b$ ,  $c$ , and  $s$  in Section 4.1. We discuss the syntax  $C \oplus P$  in Section 4.5.

In Figure 5, we define the cast insertion procedure for this calculus as a part of the typing process. The judgment  $\Gamma \vdash_G e_s \rightsquigarrow e_t : G$  can be read as: given a type environment  $\Gamma$  and a source expression  $e_s$ ,  $e_s$  has type  $G$  and is translated to a target expression  $e_t$ , which contains inserted casts. The formalization is fairly standard, except for the addition of rule FOR.

We use the syntactic form  $\llbracket G_2 \Leftarrow G_1 \rrbracket$  to denote casts that can potentially be inserted, that is, they are inserted only when  $G_1 \neq G_2$ . The rule definitions use several helper functions, given at the bottom of Figure 5. These functions extract certain parts from types when they have desired structures or Dyn when they are Dyn. Otherwise, these functions are undefined. For example, the function  $\text{ext}_L$  extracts the element type from a list type and Dyn from Dyn. The helper functions allow us to create a single rule for each source language construct, regardless of types.

The VAR rule for variable references is standard. No casts are inserted in the translation process. The ABS rule is also standard, except that it removes the parameter's annotation, since the target language is untyped. The APP rule uses the consistency relation ( $\sim$ ) to make sure that the type of the argument is consistent with the domain of the function. The definition of  $\sim$  is standard [Siek and Taha 2006], and we omit it here. If the two types are consistent, the rule translates both expressions

and inserts casts on each. The function  $(e_{1t})$  is cast to have a function type with the original type's domain and codomain, using the *dom* and *cod* helper functions. The argument is also cast so that its type matches the domain of the new function type.

In the FOR rule, there is a variable  $x$  associated with the loop that ranges over the list produced by  $e_{1s}$ . Typing ensures that  $e_{1s}$  has a type that can be treated as a list by using the *ext<sub>L</sub>* function. Consequently, the translation procedure inserts a cast on  $e_{1t}$  to a list type containing the extracted type. The LET and LETREC for non-recursive and recursive let-bindings are standard.

For references, the REF rule is trivial since no casts are inserted. For dereferences, in the DEREF rule, the source expression must be checked to ensure that its type is a reference via the *ext<sub>R</sub>* function. Thus, the translation process inserts a cast to a reference of the type extracted by *ext<sub>R</sub>* before dereferencing it. Finally, assignment in the ASSIGN rule also applies *ext<sub>R</sub>* to the reference being assigned to, and this similarly generates a cast. Additionally, the type of the expression being stored is cast to to the underlying type of the reference.

#### 4.4 Cost of Wrapped Casts

The cast insertion rules in Figure 5 can insert casts involving function types. Since whether the cast will be successful or not cannot be checked at definition time, the usual *guarded* approach handles functions by dynamically creating function *proxies* that wrap underlying functions and cast their inputs and outputs when they are called [Siek and Taha 2006]. The costs of such casts thus depend on how the are used, which our cost model so far does not consider. To illustrate, consider the following expression.

$$\mathbf{let} \ f = [\text{Int} \rightarrow \text{Int} \leftarrow \text{Dyn} \rightarrow \text{Dyn}] \lambda x. x \ \mathbf{in} \ f \ 1 + f \ 2$$

Following the ideas in Section 4.2, we will assign a potential cost  $\Lambda x.(0, x)$  to  $f$  (since there are no casts in the function body), and so the calls at  $f \ 1$  and  $f \ 2$  each generate a cost of  $\Lambda x.(0, x) \ 0$ , which reduces to 0. This, however, does not match the cost of guarded semantics where each call induces two casts, one from Int to Dyn and the other from Dyn to Int.

We want to adapt the potential assigned to  $f$  so that it includes the costs of casts in the generated proxies. However, the challenge is that to properly create these potentials, the cost analysis, which is static, need to know about the proxies, which are created dynamically. The trick is that we transform the program before analysis to syntactically include the proxies that will be generated at runtime. Specifically, we transform each function cast into a lambda expression that contains casts within its body, thereby creating a potential whose cost is not 0. The example above is transformed into the following expression.

$$\mathbf{let} \ f = \lambda y. [\text{Int} \leftarrow \text{Dyn}] (\lambda x. x \ [\text{Dyn} \leftarrow \text{Int}] y) \ \mathbf{in} \ f \ 1 + f \ 2$$

Now  $f$  has a potential  $\Lambda y.(b + c, y)$ , and each application  $f \ 1$  and  $f \ 2$  will be assigned cost  $b + c$ , yielding a total cost of  $2(b + c)$ , which matches the expected behavior of the guarded semantics for higher-order casts.

Similarly, for casts involving reference types, the guarded semantics will create a proxy that induces casts on all future dereferences and assignments. For such casts, we reuse the trick described above of embedding lambdas representing the proxies into the program before analysis. Correspondingly, we transform dereferences and assignments into lambda applications. Interestingly, this idea can treat proxied references (those that require casts) and raw references (those that do not) uniformly. To illustrate, consider the following expression, where  $y$  has type Dyn.

$$\mathbf{let} \ g = \lambda x. !x * !(\mathbf{ref} \ 1) \ \mathbf{in} \ f \ [\mathbf{ref} \ \text{Int} \leftarrow \mathbf{ref} \ \text{Dyn}] (\mathbf{ref} \ y)$$

$$\searrow ([G_1 \rightarrow G_2 \Leftarrow \text{Dyn}]e_t) = \searrow ([G_1 \rightarrow G_2 \Leftarrow \text{Dyn} \rightarrow \text{Dyn}](\searrow e_t)) \quad (1)$$

$$\searrow ([\text{Dyn} \Leftarrow G_1 \rightarrow G_2]e_t) = \searrow ([\text{Dyn} \rightarrow \text{Dyn} \Leftarrow G_1 \rightarrow G_2](\searrow e_t)) \quad (2)$$

$$\searrow ([G_3 \rightarrow G_4 \Leftarrow G_1 \rightarrow G_2]e_t) = \lambda x. (\searrow ([G_4 \Leftarrow G_2](\searrow e_t))) \searrow ([G_2 \Leftarrow G_3]x) \quad (3)$$

$$\searrow (\mathbf{ref} e_t) = \mathbf{let} y = \searrow (e_t) \mathbf{in} \lambda x. y \quad (4)$$

$$\searrow (!e_t) = (\searrow e_t) () \quad (5)$$

$$\searrow (e_{1t} := e_{2t}) = \searrow (e_{1t}) \searrow (e_{2t}) \quad (6)$$

$$\searrow ([\mathbf{ref} G_2 \Leftarrow \text{Dyn}]e_t) = \searrow ([\mathbf{ref} G_2 \Leftarrow \mathbf{ref} \text{Dyn}](\searrow e_t)) \quad (7)$$

$$\searrow ([\text{Dyn} \Leftarrow \mathbf{ref} G_2]e_t) = \searrow ([\mathbf{ref} \text{Dyn} \Leftarrow \mathbf{ref} G_2](\searrow e_t)) \quad (8)$$

$$\searrow ([\mathbf{ref} G_2 \Leftarrow \mathbf{ref} G_1]e_t) = \lambda x. (\searrow ([G_2 \Leftarrow G_1](\searrow e_t))) \searrow ([G_1 \Leftarrow G_2]x) \quad (9)$$

$$\textit{otherwise} \searrow ([G_2 \Leftarrow G_1]e_t) = [G_2 \Leftarrow G_1](\searrow e_t) \quad (10)$$

Fig. 6. Transformation rules after cast insertion. Proxies usually inserted at runtime after checking certain values are expanded syntactically, where possible. The *otherwise* in case (10) means that this rule applies when all others fail.

Note that this expression contains both a proxied reference  $[\mathbf{ref} \text{Int} \Leftarrow \mathbf{ref} \text{Dyn}](\mathbf{ref} y)$  and a raw reference  $\mathbf{ref} 1$ , which are transformed into  $\lambda a. [\text{Int} \Leftarrow \text{Dyn}]y$  and  $\lambda d. 1$ , respectively (in reality, a let binding is used to evaluate the expression in the  $\mathbf{ref}$  body before wrapping it in a lambda, but we just directly wrap the expressions here for simplicity). Each dereference is transformed into an application by applying it to a unit value,  $()$ . Overall, the transformation yields the following expression.

$$\mathbf{let} g = \lambda x. (x ()) * (\lambda d. 1) () \mathbf{in} f (\lambda a. [\text{Int} \Leftarrow \text{Dyn}]y)$$

Now let us analyze the cast costs. For the proxied reference, the potential is  $\Lambda a. (c, y)$ , and for the raw reference, the potential is  $\Lambda d. (0, 0)$ . When they are dereferenced, the corresponding applications lead to the costs  $(\Lambda a. (c, y)) 0$  and  $(\Lambda d. (0, 0)) 0$ , which are  $c$  and  $0$ , respectively, matching the expected costs of guarded semantics on dereferences.

Overall, by transforming function casts and dereferences into lambda abstractions, we can reuse the idea of potentials to precisely estimate the costs of these casts. In Figure 6, we present rules for transforming expressions as described above, where  $\searrow (e_t)$  applies the transformations to  $e_t$ . In the figure, we present rules that are relevant to casts only and ignore rules for other constructs of  $e_t$ , which recursively apply the transformation to their subterms, if applicable. The first three rules transform casts with function types. The next three transform reference expressions into lambdas and applications. The next three handle casts with references. The final rule terminates the recursive transformation in rules 3 and 9 when they arrive at casts between base types.

#### 4.5 Cost Computing Rules

This subsection defines a cost semantics that computes a cost approximation for any expression transformed by  $\searrow$ . The cost semantics is presented in Figure 7. The rules have the general form,  $\Phi \vdash_c e_t \mid A$ , meaning that expression  $e_t$  has approximation  $A$  in the context of cost environment  $\Phi$ .

In rule  $\text{VAR}$ , a variable reference  $x$  has no immediate cost since the language is call-by-value, but may have a potential cost that is retrieved from the cost environment. For example, a variable  $f$  can reference a function, which would have an abstraction for its potential. The cost of abstractions in rule  $\text{ABS}$  is  $0$  since they cause no evaluation, but their potential is an abstraction of the form  $\Lambda x. A$ .

$$\begin{array}{c}
736 \\
737 \\
738 \\
739 \\
740 \\
741 \\
742 \\
743 \\
744 \\
745 \\
746 \\
747 \\
748 \\
749 \\
750 \\
751 \\
752 \\
753 \\
754 \\
755 \\
756 \\
757 \\
758 \\
759 \\
760 \\
761 \\
762 \\
763 \\
764 \\
765 \\
766 \\
767 \\
768 \\
769 \\
770 \\
771 \\
772 \\
773 \\
774 \\
775 \\
776 \\
777 \\
778 \\
779 \\
780 \\
781 \\
782 \\
783 \\
784
\end{array}$$

$$\begin{array}{c}
\text{VAR} \frac{x \mapsto P \in \Phi}{\Phi \vdash_c x \mid (0, P)} \quad \text{ABS} \frac{\Phi, x \mapsto x \vdash_c e_t \mid A}{\Phi \vdash_c \lambda x. e_t \mid (0, \Lambda x. A)} \quad \text{APP} \frac{\Phi \vdash_c e_{1t} \mid (C_1, P_1) \quad \Phi \vdash_c e_{2t} \mid (C_2, P_2)}{\Phi \vdash_c e_{1t} e_{2t} \mid C_1 + C_2 \oplus (P_1 P_2)} \\
\text{FOR} \frac{\Phi \vdash_c e_{1t} \mid (C_1, P_1) \quad \Phi, x \mapsto P_1 \vdash_c e_{2t} \mid (C_2, P_2) \quad l \text{ fresh}}{\Phi \vdash_c \text{for } x \text{ in } e_{1t} \text{ do } e_{2t} \mid (C_1 + l \cdot C_2, l \cdot P_2)} \\
\text{LET} \frac{\Phi \vdash_c e_{1t} \mid (C_1, P_1) \quad \Phi, x \mapsto P_1 \vdash_c e_{2t} \mid (C_2, P)}{\Phi \vdash_c \text{let } x = e_{1t} \text{ in } e_{2t} \mid (C_1 + C_2, P)} \\
\text{LETREC} \frac{n = |S| \quad \Phi, x \mapsto x \vdash_c e_{1t} \mid (C_1, P_1) \quad S = \text{Apps}(x, P_1) \quad P' = P_1 \ominus S \quad l \text{ fresh} \quad \Phi, x \mapsto (l \cdot n^l \cdot P') \vdash_c e_{2t} \mid (C_2, P_2)}{\Phi \vdash_c \text{letrec } x = e_{1t} \text{ in } e_{2t} \mid (C_1 + C_2, P_2)} \\
\text{CAST} \frac{\Phi \vdash_c e_t \mid (C_1, P) \quad \text{cost}(\lceil G_2 \Leftarrow G_1 \rceil) = C_2}{\Phi \vdash_c \lceil G_2 \Leftarrow G_1 \rceil e_t \mid (C_1 + C_2, P)}
\end{array}$$

Fig. 7. Cost semantics.

In rule APP, the potential of an application is the application of the corresponding potentials. The cost of an application is the cost of the two subexpressions, plus the cost of these two casts, and the cost of the potential application. The term  $C_1 + C_2 \oplus (P_1 P_2)$  is used to pairwise add the cost of the potential application after it evaluates. For example:

$$1 \oplus (\Lambda x.(1, x) 0) = 1 \oplus [0/x](1, x) = (1 + 1, 0)$$

This term is sometimes left unevaluated. For example, the approximation for a function with a higher-order argument, such as  $\lambda x.x \ 1$ , is  $\Lambda x.(0 \oplus (x \ 0))$ , which will evaluate after we substitute a corresponding potential of the higher-order argument in a function call.

In rule FOR, a fresh loop label  $l$  is introduced to represent the number of loop iterations. The cost of the loop is then the cost of evaluating  $e_1$ , plus  $l$  times the cost of evaluating  $e_2$ . The potential is constructed similarly. For let expressions in rule LET, the cost and potential is computed similarly to FOR except that the body is evaluated only once.

The rule LETREC assigns costs to recursive expressions and bindings. Since  $e_{1t}$  refers to  $x$ , the potential  $P_1$  for  $e_{1t}$  must contain potential applications that apply  $x$  to some other potentials. Moreover, these applications are connected by  $\oplus$ . We use  $\text{Apps}(x, P_1)$  to collect all such potential applications into  $S$ . We then use  $n$  to measure the cardinality of  $S$ . For example, the value of  $n$  in a naive recursive definition of a function to compute the Fibonacci sequence would be 2, since it involves two recursive function calls. We also use the operation  $\ominus$  to remove all the applications in  $S$  from  $P_1$  and assign the result to  $P'$ . As a result,  $P'$  contains no further potential applications applying  $x$  to some term. The recursion is then estimated to have the cost  $l * n^l$ , where  $l$  is a fresh label estimating the size of the input. We then use this cost to compute the cost for  $e_{2t}$ , the overall cost for the whole construct. In general, static cost analysis for recursive programs is difficult and an area receiving significant recent research [Danner et al. 2015; Hoffmann et al. 2017]. This difficulty is further exacerbated in our case since the type information that is available in other static cost analyses is unavailable in the gradual type setting. Our costs for recursions are a coarse upper bound, and we assume the input to recursive calls strictly decrease in size, similar to previous



785 cost analyses [Danner et al. 2015]. Nevertheless, this cost estimation works quite well in practice  
 786 because the bounds for recursion do not affect relative cost comparisons, in particular the bounds  
 787 are shared for nearby type configurations in the cost lattice.

788 Finally, the rule `CAST` assigns a cost to each expression being cast, whose cost is that of the  
 789 underlying expression plus that of the cast, according to the *cost* function from Section 4.1. Rules for  
 790 approximating the costs of references and higher-order casts are handled by costs for abstractions  
 791 and applications after using the transformation procedure in Figure 6.

792

793

#### 4.6 Properties

794 In the following lemmas and theorems, we use the judgment form  $\Phi; \Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid A$ , which is  
 795 equivalent to the judgment  $\Gamma \vdash_G e_s \rightsquigarrow e_t : G$  followed by  $\Phi \vdash_c \searrow (e_t) \mid A$ . Essentially, the judgment  
 796 can be read as: under  $\Gamma$  and  $\Phi$ ,  $e_s$  has type  $G$ , is translated to  $e_t$ , and has the cost approximation  $A$ .

797 Before we present the most important properties of our cost semantics, we present some simple  
 798 lemmas relating terms in the source language to terms in the cost language. The first lemma states  
 799 that a bound variable, if referenced, affects the potential of its abstraction.

800

801

LEMMA 4.1. *If  $\Phi; \Gamma \vdash_{GC} \lambda x. e_s \rightsquigarrow \lambda x. e_t : G \mid (C, P)$  and  $x \in \text{vars}(e_s)$ , then  $x \in \text{vars}(P)$ .*

802 This lemma confirms that the potential costs of an argument, which may be a function with its  
 803 own costs, are reflected in the cost approximation of the function that uses it.

804 The second lemma states that substitution in the source language corresponds to substitution in  
 805 the cost language.

806

807

808

LEMMA 4.2. *Let  $\Phi; \Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid (C, P)$  where  $x \in \text{vars}(e_s)$ ,  $x \mapsto x \in \Phi$ , and  $x : G'$ .  
 If  $\Phi; \Gamma \vdash_{GC} e_s' \rightsquigarrow e_t' : G' \mid (C', P')$ , then the potential for  $[e_s'/x]e_s$  is  $[P'/x]P$ .*

809 This establishes how potential applications can insert the overhead of casts in the argument terms  
 810 into the body of a potential abstraction. Together, these two lemmas establish the correspondence of  
 811 function abstraction and application at the source level and the cost level. Moreover, since we also  
 812 reason about the overhead of using references by translating them into lambdas and then applying  
 813 our cost semantics, it is imperative that abstractions and applications be modeled correctly.

814 We now establish the most important properties of our cost semantics. The following theo-  
 815 rem states that our cost semantics terminates, provided that the program after translation ( $e_t$ ) is  
 816 terminating.

817

818

819

THEOREM 4.3. *[Cost Derivation Termination] For any terminating program  $e_s$ ,  $\Phi; \Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid A$  terminates and produces the approximation  $A$ .*

820 Finally, the most important result is that our cost semantics bounds the number of casts of the  
 821 program, provided that its recursions, if any, are applied to smaller arguments.

822

823

824

825

THEOREM 4.4. *[Costs are Upper Bounds] If  $\Phi; \Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid (C, P)$ , then  $C$  is greater than or  
 equal to the number of casts performed when executing  $e_t$  after replacing loop labels by the number of  
 iterations performed.*

826 The proofs of these two theorems are produced by induction over the rules in Figures 5, referencing  
 827 the rules in Figure 7. The theorems also need to connect with the dynamic semantics for the language,  
 828 which is a fairly standard semantics with a few extensions for loops and other constructs.

829

830

## 5 COSTS FOR ALL CONFIGURATIONS

831 In this section we extend the formalization in Section 4 to make it variational, enabling us to  
 832 efficiently estimate costs for all possible type configurations of a program. Thus, instead of assigning  
 833

|     |                   |   |
|-----|-------------------|---|
| 834 | Target expr.      | $e_t^v ::= \dots \mid [M \Leftarrow M]e_t^v \mid d\langle [M \Leftarrow M], [M \Leftarrow M] \rangle e_t^v$ |
| 835 | Costs             | $C^v ::= \dots \mid d\langle C^v, C^v \rangle$  |
| 836 | Potentials        | $P^v ::= \dots \mid d\langle P^v, P^v \rangle$  |
| 837 | Approximations    | $A^v ::= \dots \mid (C^v, P^v) \mid C^v \oplus P^v$   |
| 838 | Variational types | $V ::= \gamma \mid V \rightarrow V \mid [V] \mid \text{ref } V \mid d\langle V, V \rangle$                  |
| 839 | Migrational types | $M ::= \gamma \mid M \rightarrow M \mid [M] \mid \text{Dyn} \mid \text{ref } M \mid d\langle M, M \rangle$  |
| 840 | Configuration     | $K ::= \emptyset \mid K, x \mapsto G$   |
| 841 | Choice env.       | $\Omega ::= \emptyset \mid \Omega, x \mapsto M$   |
| 842 |                   |   |

843 Fig. 8. Syntax for variational cost analysis. Definitions of variational artifacts (e.g.  $e_t^v$ ) extend the syntax for  
 844 non-variational counterparts (e.g.  $e_t$ ) in Figure 4.

846 a gradual type and a cost to each program, the rules in this section assign a variational gradual  
 847 type, called a *migrational type* [Campora et al. 2018] and a variational cost to each program. The  
 848 migrational type of a program represents the type of all possible type configurations, while its  
 849 variational cost encodes the cost lattice of migrating to each configuration, as illustrated in Figure 2.

## 851 5.1 Syntax

852 In Figure 8, we extend the syntax defined in Figure 4 to accommodate the variational analysis. In  
 853 particular, we extend costs and potentials with a choice construct, as described in Section 2.2. This  
 854 enables the representation of variational costs and variational potentials. The addition of choices  
 855 to types yields two new domains, *variational types* (static types with choices) and *migrational types*  
 856 (gradual types with choices), ranged over by  $V$  and  $M$ , respectively. Casts in the target language  
 857 are made variational by allowing casting to a migrational type.

858 Figure 8 also defines *configurations*, ranged over by  $K$ , and *choice environments*, ranged over by  
 859  $\Omega$ . A configuration is a mapping from initially dynamic parameters to gradual types, denoting  
 860 their type assignments for a particular configuration. A choice environment is a mapping from  
 861 parameters to migrational types and keeps track of the type assignments used to generate the  
 862 full configuration space for a given program. For example, if a program has dynamic parameters  
 863  $x$  and  $y$ , then the choice environment  $\Omega = \{x \mapsto B\langle \text{Dyn}, \text{Int} \rangle, y \mapsto D\langle \text{Dyn}, \text{Int} \rangle\}$  encodes four  
 864 configurations (each of  $x$  and  $y$  can be  $\text{Dyn}$  or  $\text{Int}$ ). For simplicity, we assume that all parameter  
 865 names in the program are all unique. Configurations and choice environments are used to establish  
 866 the correctness of our variational cost analysis by making explicit the relation between the rules in  
 867 this section and those in Section 4.3.

868 Although migrational types now appear in casts in the target language, we do not extend its  
 869 dynamic semantics since in the implementation programs with variational casts are not executed,  
 870 only analyzed. After the analysis, the programmer would select a particular type configuration by  
 871 applying a complete decision to the program that eliminates migrational casts and yields a runnable  
 872 program.

## 874 5.2 Variational Cost Analysis

875 The cost semantics defined in Section 4 requires some changes to accommodate casts to migrational  
 876 types. First, we extend the cost function with the following new rules, which essentially maps the  
 877 cost calculation over cast choices, preserving the results in cost choices.

$$879 \text{cost}(\llbracket d\langle M_1, M_2 \rangle \Leftarrow M \rrbracket) = d\langle \text{cost}(\llbracket M_1 \Leftarrow M \rrbracket), \text{cost}(\llbracket M_2 \Leftarrow M \rrbracket) \rangle$$

$$880 \text{cost}(\llbracket M \Leftarrow d\langle M_1, M_2 \rangle \rrbracket) = d\langle \text{cost}(\llbracket M \Leftarrow M_1 \rrbracket), \text{cost}(\llbracket M \Leftarrow M_2 \rrbracket) \rangle$$

883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931

$$\begin{array}{c}
\text{VARV} \frac{x \mapsto M \in \Gamma}{\Gamma \vdash^v x \rightsquigarrow x : M \mid \Omega} \qquad \text{ABSV} \frac{\Gamma, x \mapsto T \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega}{\Gamma \vdash^v \lambda x : T.e_s \rightsquigarrow \lambda x.e_t^v : T \rightarrow M \mid \Omega} \\
\text{ABSDYNV} \frac{\Gamma, x \mapsto d\langle \text{Dyn}, V \rangle \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega \quad d \text{ fresh}}{\Gamma \vdash^v \lambda x.e_s \rightsquigarrow \lambda x.e_t^v : d\langle \text{Dyn}, V \rangle \rightarrow M \mid \Omega \cup \{x \mapsto d\langle \text{Dyn}, V \rangle\}} \\
\text{APPV} \frac{\Gamma \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad \Gamma \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M_2 \mid \Omega_2 \quad \text{dom}(M_1) \approx M_2}{\Gamma \vdash^v e_{1s} e_{2s} \rightsquigarrow \llbracket \text{dom}(M_1) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket e_{1t}^v \llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket e_{2t}^v : \text{cod}(M_1) \mid \Omega_1 \cup \Omega_2} \\
\text{FORV} \frac{\Gamma \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad \text{ext}_L(M_1) = M \quad \Gamma, x \mapsto M \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M_2 \mid \Omega_2}{\Gamma \vdash^v \text{for } x \text{ in } e_{1s} \text{ do } e_{2s} \rightsquigarrow \text{for } x \text{ in } \llbracket [M] \Leftarrow M_1 \rrbracket e_{1t}^v \text{ do } e_{2t}^v : M_2 \mid \Omega_1 \cup \Omega_2} \\
\text{LETV} \frac{\Gamma \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad \Gamma, x \mapsto M_1 \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M \mid \Omega_2}{\Gamma \vdash^v \text{let } x = e_{1s} \text{ in } e_{2s} \rightsquigarrow \text{let } x = e_{1t}^v \text{ in } e_{2t}^v : M \mid \Omega_1 \cup \Omega_2} \\
\text{LETREC} \frac{\Gamma, x \mapsto M_1 \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad \Gamma, x \mapsto M_1 \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M \mid \Omega_2}{\Gamma \vdash^v \text{letrec } x = e_{1s} \text{ in } e_{2s} \rightsquigarrow \text{letrec } x = e_{1t}^v \text{ in } e_{2t}^v : M \mid \Omega_1 \cup \Omega_2} \\
\text{REFV} \frac{\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega}{\Gamma \vdash^v \text{ref } e_s \rightsquigarrow \text{ref } e_t^v : \text{ref } M \mid \Omega} \qquad \text{DEREFV} \frac{\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M_1 \mid \Omega \quad M = \text{ext}_R(M_1)}{\Gamma \vdash^v !e_s \rightsquigarrow !\llbracket \text{ref } M \Leftarrow M_1 \rrbracket e_t^v : M \mid \Omega} \\
\text{ASSIGNV} \frac{\Gamma \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad M = \text{ext}_R(M_1) \quad \Gamma \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M_2 \mid \Omega_2 \quad M \approx M_2}{\Gamma \vdash^v e_{1s} := e_{2s} \rightsquigarrow \llbracket \text{ref } M \Leftarrow M_1 \rrbracket e_{1t}^v := \llbracket M \Leftarrow M_2 \rrbracket e_{2t}^v : \text{ref } M \mid \Omega_1 \cup \Omega_2} \\
\begin{array}{lll}
\text{dom}(M_1 \rightarrow M_2) = M_1 & \text{dom}(\text{Dyn}) = \text{Dyn} & \text{dom}(d\langle M_1, M_2 \rangle) = d\langle \text{dom}(M_1), \text{dom}(M_2) \rangle \\
\text{cod}(M_1 \rightarrow M_2) = M_2 & \text{cod}(\text{Dyn}) = \text{Dyn} & \text{cod}(d\langle M_1, M_2 \rangle) = d\langle \text{cod}(M_1), \text{cod}(M_2) \rangle \\
\text{ext}_L([M]) = M & \text{ext}_L(\text{Dyn}) = \text{Dyn} & \text{ext}_L(d\langle M_1, M_2 \rangle) = d\langle \text{ext}_L(M_1), \text{ext}_L(M_2) \rangle \\
\text{ext}_R(\text{ref } M) = M & \text{ext}_R(\text{Dyn}) = \text{Dyn} & \text{ext}_R(d\langle M_1, M_2 \rangle) = d\langle \text{ext}_R(M_1), \text{ext}_R(M_2) \rangle
\end{array}
\end{array}$$

Fig. 9. Cast insertion rules after adding variational types to our system. The operations  $\text{dom}$ ,  $\text{cod}$ ,  $\text{ext}_L$ , and  $\text{ext}_R$  are undefined for cases that are not listed here.

Next, we similarly extend the transforming procedure ( $\searrow$ ) to push the translation into choices. Following the same idea, we extend the cost relation  $\vdash_c$  in Figure 7 to deal with variations. We name these extended relations  $\searrow^v$  and  $\vdash_c^v$ , respectively. An interesting bit in  $\vdash_c^v$  is that our cost function now yields variational costs, and so we must define how arithmetic works on such values. Intuitively, any basic operation on a variational cost can be performed by pushing the operations into the choices. For example,  $c + B\langle 2c, 3c \rangle = B\langle c + 2c, c + 3c \rangle = B\langle 3c, 4c \rangle$ . We omit a formal definition here since it is straightforward.

In Figure 9, we present the revised set of type-directed cast insertion rules to support migrational types. The judgment has the form  $\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega$ , which states that under type environment  $\Gamma$ , the source expression  $e_s$  has type  $M$  and is translated to the target expression  $e_t^v$  after casts are inserted, where type change information is recorded in  $\Omega$ . The cast insertion rules for variables,

932 loops, let expressions, and references are nearly identical to the corresponding rules in Figure 5,  
 933 except that they now use the extended syntax.

934 There are now two rules for abstractions, one for statically typed parameters ( $\text{AbsV}$ ) and the  
 935 other for dynamically typed parameters ( $\text{AbsDynV}$ ). The  $\text{AbsDynV}$  rule is where variation is injected  
 936 into the target program so that it captures the whole migration space. The choice type in  $\text{AbsDynV}$   
 937 represents the fact that there are two possibilities during program migration: leaving the parameter  
 938 dynamically typed, or changing it to a static type. Correspondingly, when we carry out our cost  
 939 analysis there will be two different costs for the body. For example, in  $\lambda x.x + 1$ , the reference to  $x$   
 940 in the body must be cast using  $\lceil \text{Int} \leftarrow d\langle \text{Dyn}, \text{Int} \rangle \rceil$ , which will be assigned the cost  $d\langle c, 0 \rangle$  since  
 941 we incur a cast when  $x$  has type  $\text{Dyn}$  and no cast when it has type  $\text{Int}$ . In  $\text{AbsDynV}$ , we also extend  
 942  $\Omega$  to record that we assigned the type  $d\langle \text{Dyn}, V \rangle$  to the parameter.

943 Variation complicates the treatment of function applications. First, the *dom* and *cod* functions are  
 944 extended to support migrational types. Second, the type consistency relation used in the original  
 945  $\text{APP}$  rule of Figure 5 is replaced by the *compatibility* relation ( $\approx$ ) defined in [Campora et al. 2018],  
 946 which extends type consistency to also support variational type equivalence (see Section 2.2). Two  
 947 types  $M_1$  and  $M_2$  are compatible, written as  $M_1 \approx M_2$ , if they are consistent or compatible under  
 948 any selection. For example,  $B\langle \text{Int}, \text{Dyn} \rangle \approx B\langle \text{Dyn}, D\langle \text{Bool}, \text{Int} \rangle \rangle$ , since selecting  $B.1$  in both types  
 949 yields  $\text{Int} \sim \text{Dyn}$ , while selecting  $B.2$  in both yields  $\text{Dyn} \approx D\langle \text{Bool}, \text{Int} \rangle$ .

### 951 5.3 Properties

952 In this subsection, we prove the correctness of our variational cost analysis by showing that it is  
 953 equivalent to generating all type configurations and applying the cost analysis from Section 4 to  
 954 each one individually. We introduce a new judgment  $\Phi; \Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid A^v \mid \Omega$  as shorthand for  
 955  $\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega$  followed by  $\Phi \vdash_c^v \searrow^v (e_t^v) \mid A^v$ .

956 We first extend selection, defined in Section 2.2, to choice environments by applying selections to  
 957 its range. For example, let  $\Omega_a = \{x \mapsto B\langle \text{Dyn}, \text{Int} \rangle, y \mapsto D\langle \text{Dyn}, \text{Int} \rangle\}$ , then  $\lfloor \Omega_a \rfloor_{\{B.1, D.2\}} = K_a$ , where  
 958  $K_a = \{x \mapsto \text{Dyn}, y \mapsto \text{Int}\}$ . We also need a way to apply a configuration to a typing process in order  
 959 to conveniently type different individual type configurations. We write  $\Phi; K\Gamma \vdash_{GC} e_{1s} \rightsquigarrow e_{1t} : G \mid A$   
 960 to express that the cost analysis (Section 4.6) is directed by the configuration  $K$ . The configuration  $K$   
 961 overrides the type assignments in the environment  $\Gamma$ , so for each variable reference  $x$ , if  $x \mapsto G \in K$ ,  
 962 then  $x$  has type  $G$ , otherwise it retrieves the type from  $\Gamma$ . Since we assume all parameters have  
 963 unique names, this can be achieved without ambiguity. For example, let  $\text{add} = \lambda x : \text{Dyn}. \lambda y : \text{Dyn}. x + y$   
 964 and  $\Phi; \Gamma \vdash_{GC} \text{add} \rightsquigarrow e_{1t} : \text{Int} \mid (2c, P_1)$ , then  $e_{1t}$  includes two casts  $\lceil \text{Int} \leftarrow \text{Dyn} \rceil$  to the parameters  
 965  $x$  and  $y$ , whereas in  $\Phi; K_a\Gamma \vdash_{GC} \text{add} \rightsquigarrow e_{2t} : \text{Int} \mid (c, P_2)$ , then  $e_{2t}$  includes a cast  $\lceil \text{Int} \leftarrow \text{Dyn} \rceil$  to the  
 966 parameter  $x$  only. The reason is that  $K_a$  forces the parameter  $y$  to have type  $\text{Int}$ .

967 We can now state the correctness of the rules in Figure 9 in two steps. First, Theorem 5.1 states  
 968 that the cost of any configuration can be obtained through some variational cost calculation.

970 **THEOREM 5.1 (VARIATIONAL COST COMPLETENESS).** *For any  $K$ , if  $\Phi; K\Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid A$ ,  
 971 then there exists some variational cost analysis  $\Phi; \Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid A^v \mid \Omega$  such that  $e_t = \lfloor e_t^v \rfloor_\delta$ ,  
 972  $A = \lfloor A^v \rfloor_\delta$ ,  $G = \lfloor M \rfloor_\delta$ , and  $K = \lfloor \Omega \rfloor_\delta$ , where  $\delta$  can be decided by  $e_s$  and  $K$ .*

974 Given an expression  $e$  and a configuration  $K$  for  $e$ , the corresponding decision  $\delta$  can be deduced by  
 975 considering which alternative of each choice in  $e$  must be selected. For example, in the *add* example  
 976 above, the decision for  $K_a$  is  $\{B.1, D.2\}$ .

977 According to Theorem 5.1, it's possible that we need to use different variational cost analyses to  
 978 obtain costs for different configurations. Fortunately, the following theorem shows that we can  
 979 find appropriate variational types for dynamic parameters such that the costs for all configurations  
 980

981 can be obtained through just one variational cost analysis. We capture this idea in the following  
 982 theorem.

983 **THEOREM 5.2 (VARIATIONAL COSTS SOUNDNESS AND EFFICIENCY).** *Given any  $e_s$  for which the*  
 984 *entire configuration space is well-typed, there exists some  $\Omega$  for  $\Phi; \Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid A^v \mid \Omega$  such*  
 985 *that  $\Phi; \lfloor \Omega \rfloor_\delta \Gamma \vdash_{GC} e_s \rightsquigarrow \lfloor e_t^v \rfloor_\delta : \lfloor M \rfloor_\delta \mid \lfloor A^v \rfloor_\delta$  for any  $\delta$ .*

986  
 987 The proof of both theorems follows by induction over the rules in Figure 9, connecting to the  
 988 rules in Figure 5. The  $\Omega$  in Theorem 5.2 can be found by migrational type inference [Campora  
 989 et al. 2018]. Moreover, we can lift the restriction about the configuration space being well-typed by  
 990 employing the pattern-constrained judgments introduced in that work, and our implementation  
 991 uses this approach. In fact, the inference process and variational cost analysis can be combined and  
 992 computed together, and we refer to this combined process as *Cost Space Typing* (CST).

#### 993 5.4 Uses of Variational Costs

994  
 995 In Section 1.2, we outlined several scenarios programmers might encounter when migrating  
 996 gradually typed programs. This subsection will briefly describe how, given a source program  $e_s$ ,  
 997 the CST  $\Phi; \Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid A^v \mid \Omega$  can be used to solve the problems posed by these scenarios.

998 Supporting the scenarios often involves comparing costs that refer to different loop variables  
 999 and so are not comparable in principle. As a pragmatic solution, we simply instantiate all loop  
 1000 variables by the same large integer and then directly compare the resulting values. For example,  
 1001 given  $2l_1 + 4$  and  $3l_2 + 3$ , we can instantiate both  $l_1$  and  $l_2$  by 100 and conclude that the first cost  
 1002 is cheaper than the second since  $204 < 303$ . The evaluation in Section 6 shows that this simple  
 1003 solution yields good results, although future work could attempt to constrain loop variables more  
 1004 precisely through program analysis.

1005 In scenario S1, the goal is to maximize static safety first, then optimize for performance. Given the  
 1006 migrational typing, we can apply the method described by Campora et al. [2018] to extract the set  
 1007 of decisions corresponding to the most static migrations. Each decision  $\delta_i$  in this set characterizes a  
 1008 maximal set of parameters to be annotated in  $e_s$ . Finally, we select the  $\delta_i$  that minimizes the value  
 1009 of  $\lfloor A^v \rfloor_{\delta_i}$  after instantiating loop variables in the approximation  $A^v$  as described above.

1010 In scenario S2, the goal is to maximize performance, before maximizing the presence of static  
 1011 types for safety. To do this, we instantiate the loop variables then search  $A^v$  for the decisions  
 1012 with minimum cost. The search procedure is a straightforward recursive function that keeps track  
 1013 of the lowest cost encountered as it builds up the decision(s) corresponding to the lowest costs  
 1014 encountered. We can allow the user to require certain functions or parameters to be annotated by  
 1015 simply selecting  $A^v$  with the corresponding selectors before searching. This allows programmers  
 1016 to optimize performance while still enforcing certain typing goals. To find incremental migrations  
 1017 where at most  $c$  out of  $n$  parameters are migrated, we first compute the set of  $\binom{n}{c}$  decisions where  
 1018  $c$  parameters are annotated, then minimize  $A^v$  under this set of decisions.

1019 In scenario S3, the goal is to increase type-based safety guarantees without hurting performance.  
 1020 This can be achieved by searching  $A^v$  for all costs that are lower than the current program. We can  
 1021 efficiently represent the result of such a search as a mask that can be applied to the typing results,  
 1022 similar to the treatment of type errors in Chen et al. [2012]. This would effectively produce a set of  
 1023 migrations that increase performance, after which we could maximize the static checking done in  
 1024 the remaining migrations, as described in the solution to S1.

1025 Finally, in scenario S4, the goal is to understand why a particular migration exhibits degraded  
 1026 performance and to make informed decisions about the performance of different migrations.  
 1027 Supporting this scenario requires selecting the relevant migrations from  $A^v$ , then translating the  
 1028 symbolic costs in each into explanations. For example, when considering configurations ① and ② in  
 1029

the program in Figure 1, we can report that ② has worse performance than ① since it inserts casts in the  $l_1$  and  $l_2$  loops in `part_A_times_u` and `part_At_times_u` that were not there before. Currently, HERDER contains a basic implementation capable of reporting the differences in the set of loop labels between two configurations based on their costs.

## 6 IMPLEMENTATION AND EVALUATION

In this section we discuss HERDER, a tool for carrying out variational cost analysis for Reticulated Python programs. In Section 6.1 we give an overview of its implementation. We then evaluate how well HERDER realizes the capabilities described in Section 1.3, which are needed to support the scenarios described in Section 1.2. Since capability C1, migrational type inference, is provided by earlier work [Campora et al. 2018], we focus on C2 and C3. Specifically, in Sections 6.2 and 6.3, we evaluate how effectively and efficiently, respectively, HERDER identifies performant configurations.

### 6.1 Implementation

HERDER is implemented on top of Reticulated’s guarded semantics [Vitousek et al. 2014], extending it with migrational type inference [Campora et al. 2018], cost analysis of casts, and variations. In addition to the constructs from Figure 4, our implementation supports conditionals and other Python constructs. We use the `sympy` package—which supports creating, algebraically manipulating, and substituting symbolic variables—to implement the approximation language.

### 6.2 Evaluation of Effectiveness

We evaluate the effectiveness of HERDER by showing how often it can pick the most performant configurations among all valid configurations. Our sample programs are drawn from the Python benchmarks suite,<sup>3</sup> and are largely a subset of those used by Vitousek et al. [2017] to evaluate Reticulated, plus some programs from the `scimark` benchmark. We also generated three synthetic benchmarks (`syn`) to stress-test the performance of HERDER, when the cost analysis contains deeply nested choices, a large number of choices, or a large number of code lines, respectively.

Among the 15 programs we evaluated, 3 of them, namely `fankuch`, `nqueens`, and `pyflate`, have very uniform performance, meaning that the addition of type annotations has very little effect on performance. As expected, HERDER correctly identifies the most static configuration as having the best performance. For this reason, in the rest of this section, we will not discuss them in detail.

The metrics and evaluation results for the remaining 12 programs, are given in Figure 10. The first column group lists the name of each benchmark, its lines of code, and the total number of function parameters it contains. The number of parameters describes the size of the search space for our analysis. Since each parameter can either be statically annotated by the inferred type or not, a program with  $p$  parameters contains up to  $2^p$  configurations.

The second column group reports the results of running HERDER on each of the benchmarks. The Analysis column reports the time to perform the analysis in seconds. The Rec column reports the number of static type annotations that HERDER recommended be added to the program. Finally, the Run column reports the time in seconds needed to execute the target program after adding the recommended type annotations.

To evaluate the effectiveness of HERDER at recommending performant configurations, we must compare the recommended configuration to other possibilities. For each benchmark, if there are fewer than 100 potential configurations, we generate and time every one. If there are more than 100 potential configurations (i.e. if the number of parameters is greater than 6), we randomly sample a

<sup>3</sup><http://pyperformance.readthedocs.io/benchmarks.html>

|      | Bench    | LOC   | #P   | Analysis | Rec  | Run   | Dynamic |       | Worst |       | Best  |       | Top 3 |
|------|----------|-------|------|----------|------|-------|---------|-------|-------|-------|-------|-------|-------|
|      |          |       |      |          |      |       | Time    | Ratio | Time  | Ratio | Time  | Ratio |       |
| 1079 |          |       |      |          |      |       |         |       |       |       |       |       |       |
| 1080 | float    | 64    | 6    | 1.67     | 3    | 69.7  | 103.6   | 1.49  | 128.7 | 1.83  | 61.5  | 0.88  | ✓     |
| 1082 | meteor   | 254   | 27   | 12.92    | 2    | 27.0  | 19.0    | 0.70  | 134.2 | 4.96  | 19.0  | 0.70  | ✓     |
| 1083 | nbody    | 157   | 16   | 3.34     | 6    | 64.3  | 65.6    | 1.02  | 174.0 | 2.70  | 59.6  | 0.93  | ✓     |
| 1084 | pidigits | 68    | 5    | 0.71     | 4    | 8.1   | 47.0    | 5.83  | 45.7  | 5.65  | 8.0   | 0.99  | ✓     |
| 1085 | raytrace | 448   | 66   | 63.60    | 47   | 30.1  | 56.7    | 1.88  | 109.3 | 3.63  | 30.1  | 1     | ✓     |
| 1086 | sm(FFT)  | 140   | 14   | 1.95     | 8    | 34.4  | 35.0    | 1.02  | 48.0  | 1.40  | 29.1  | 0.85  | ✗     |
| 1087 | sm(MC)   | 97    | 5    | 0.96     | 4    | 40.8  | 7.4     | 0.18  | 48.9  | 1.2   | 7.4   | 0.18  | ✓     |
| 1088 | sm(SOR)  | 110   | 16   | 7.19     | 16   | 37.8  | 97.9    | 2.59  | 120.7 | 3.19  | 37.8  | 1     | ✓     |
| 1089 | spectral | 85    | 4    | 0.82     | 2    | 17.9  | 48.8    | 2.73  | 93.8  | 5.24  | 17.9  | 1     | ✓     |
| 1090 | syn_1    | 363   | 41   | 139.49   | 34   | 77.0  | 128.6   | 1.67  | 270.9 | 3.51  | 77.0  | 1     | ✓     |
| 1091 | syn_2    | 3710  | 787  | 1438.47  | 716  | 10.0  | 30.8    | 3.08  | 53.6  | 5.36  | 10.0  | 1     | ✓     |
| 1092 | syn_3    | 15213 | 2505 | 14607.84 | 2444 | 139.8 | 131.3   | 0.94  | 452.9 | 3.24  | 131.3 | 0.94  | ✓     |

Fig. 10. Evaluation of HERDER. All timing results in the table are in seconds. The first column group provides basic stats about each benchmark: its name, lines of code, and the number of parameters it contains. The second group describes the results: the time to perform the analysis, the number of recommended static annotations, and the runtime of the resulting program. The third group compares the runtime of the resulting program to other potential configurations: the fully dynamic program, and the worst and best configurations identified in our sample. The final column indicates whether the configuration recommended by HERDER is among the top 3 of most performant configurations (✓) or not (✗).

set of 100 configurations, generating and timing the variant program for each. As a baseline, we also measure the runtime of the fully dynamically typed version of each benchmark.

The third column group reports results from this comparison. The first pair of columns reports the runtime of the fully dynamic version of the benchmark and the ratio of the dynamic version over our recommended configuration. The next two pairs of columns report the runtime and corresponding ratios for the worst and best configuration in our sample for each benchmark. Finally, in the last column of the table, we report how the runtime of our recommended configuration ranked amongst the runtimes of all variants in the sample.

The results demonstrate that HERDER is effective at finding performant configurations. In 11/12 cases, it recommends one of the top three configurations. When the recommendation does not achieve the best performance, its recommendation is usually within 15% of the optimal configuration. Moreover, we observe that the worst configuration is often 3-5x slower than HERDER's pick. Some of the ratio results along with the top 3 result might seem odd. For example, sm(MC) has a poor "Best" ratio, yet HERDER's pick remains in top 3, while sm(FFT) has a better "Best" ratio, but the pick is not in top 3. The reason for the poor "Best" ratio for sm(MC) is that the performance for the dynamic configuration was an outlier and Herder's pick had the fastest time compared to the rest. With respect to sm(FFT), most of the configurations had similar times, so the pick happened to fall out of top 3, even though the difference between the third best time and the pick's time was small.

Three interesting cases are the scimark Monte Carlo benchmark, sm(MC); the meteor\_contest benchmark; and the syn\_3 benchmarks, where the fully dynamic version of the program is faster than any gradually typed version. Currently, HERDER only introduces choices to reason about alternative function parameter types but it infers return types directly. This means it does not consider the fully dynamic version of the program as a potential configuration. This limitation can be easily remedied by extending HERDER to reason variationally about return types in addition to

1128 parameter types, allowing return types to remain dynamic when needed, so that accurate costs can  
1129 be provided for these benchmarks.

1130 Notice that the fully dynamic configuration does not necessarily have near-best performance, as  
1131 is typical in Racket [Takikawa et al. 2016]. This is due to the fact that Python programs frequently  
1132 unpack parameters via tuple assignment, which causes casts to be inserted. Typed literals can also  
1133 cause some overhead in programs. Consequently, in the Rec column, we observe a fair amount of  
1134 variability in whether the best configurations are more or less static. In some benchmarks, such  
1135 as nboddy, the recommended configuration has fewer type annotations, and in sm(MC), we know  
1136 the best configuration contains no annotations. On the other hand, the best configuration for the  
1137 scimark successive over-relaxation benchmark, sm(SOR), adds annotations to all parameters.

1138 **Effectiveness of supporting program migration scenarios** We next evaluate how well Herder  
1139 can support each migration scenario we outlined in Section 1.2. For scenarios S1 and S2 we use  
1140 Figure 10 as evidence. Specifically, Figure 10 shows how well HERDER can find globally performant  
1141 configurations via the cost analysis, and consequently it directly supports S2.

1142 To support S1, we first need to identify the migration space that contains all different migrations  
1143 that maximize static type safety, which can be realized efficiently through the method described  
1144 in Campora et al. [2018]. Performance optimization in this scenario then amounts to locating the  
1145 most performant migration within the identified migration space, which is much smaller than  
1146 the global space HERDER searches for in Figure 10. Therefore, we believe that the effectiveness  
1147 demonstrated in Figure 10 carries over to this scenario. As evidence of this argument, in our  
1148 evaluation for S3 in the next paragraph, HERDER similarly has to search a small space, and it  
1149 achieved good results.

1150 In scenario S3, we are considering the case where a recently added type annotation hurt per-  
1151 formance, and asking HERDER whether it is preferable to remove the annotation or add more  
1152 annotations to restore performance. To evaluate this scenario, from the 9 non-synthetic programs  
1153 in Figure 10, we generated 45 configurations by randomly adding type annotations. For 15 of these  
1154 configurations, the performance is better than their corresponding untyped configurations. For  
1155 the remaining 30 configurations, we asked HERDER whether adding more type annotations (and if  
1156 so, which ones) or removing the annotations yields better performance. In 28 out of 30, HERDER  
1157 generates correct recommendations, yielding an accuracy of 93.3%.

1158 Finally, for S4, it is hard to empirically verify HERDER's ability to support this scenario without  
1159 a user study. Still, Figure 10 and our evaluation for S3 show that HERDER is effective at both  
1160 finding performant configurations and directly comparing the performance of two configurations.  
1161 Consequently, explanations generated by HERDER are likely to help users understand the perfor-  
1162 mance bottleneck present in the slower configuration. We leave a user study evaluating HERDER's  
1163 effectiveness for supporting S4 to future work.

### 1164 6.3 Evaluation of Efficiency

1165 In this subsection we evaluate how HERDER scales with the size and complexity of the source  
1166 program. First, we consider how HERDER scales as the number of type parameters in the program  
1167 increases. Each parameter effectively doubles the size of the search space since it can either be  
1168 annotated by a static type or not. For our evaluation, we artificially generate a set of programs with  
1169 an increasing number of parameters. Programs with 2–12 parameters were produced by repeating  
1170 a small arithmetic function with two parameters 1–6 times; programs with more parameters were  
1171 produced by copying, pasting, and renaming several functions from the FFT and nboddy benchmarks.  
1172 For each program, we measure the runtime of HERDER to type and cost the entire configuration  
1173 space and produce a recommendation. As baselines for comparison, we also measure the runtime  
1174  
1175  
1176



|      | Params | HERDER  | Brute-force | Dynamic | LOC  | HERDER | Brute-force | Dynamic |
|------|--------|---------|-------------|---------|------|--------|-------------|---------|
| 1177 |        |         |             |         |      |        |             |         |
| 1178 | 2      | 0.07    | 0.27        | 0.03    | 89   | 0.83   | 7.81        | 0.44    |
| 1179 | 4      | 0.12    | 1.48        | 0.05    | 130  | 1.20   | 12.40       | 0.73    |
| 1180 | 6      | 0.18    | 7.47        | 0.09    | 225  | 2.08   | 21.70       | 1.10    |
| 1181 | 8      | 0.27    | 36.97       | 0.12    | 1090 | 10.66  | 105.18      | 5.38    |
| 1182 | 10     | 0.38    | 128.75      | 0.14    | 2146 | 25.63  | 296.72      | 14.16   |
| 1183 | 12     | 0.531   | 2545.33     | 0.17    | 5010 | 80.60  | 605.70      | 30.82   |
| 1184 | 51     | 13.03   | -           | 4.12    |      |        |             |         |
| 1185 | 99     | 40.41   | -           | 13.09   |      |        |             |         |
| 1186 | 195    | 181.86  | -           | 46.76   |      |        |             |         |
| 1187 | 387    | 675.64  | -           | 178.94  |      |        |             |         |
| 1188 | 807    | 2896.76 | -           | 773.20  |      |        |             |         |
| 1189 |        |         |             |         |      |        |             |         |

Fig. 11. Runtime of HERDER as the number of parameters increases (left) and the number of LOC increases but keep the number of parameters the same (right), compared to the runtime of a brute-force search and the time to type and cost a single all-dynamic variant. All times are in seconds.

of a corresponding brute-force search of all configurations (that is, generating each configuration and typing/costing each one separately), and also the runtime of typing and costing a single configuration—in this case, the configuration where all parameters are dynamically typed.

The results of the evaluation are shown in the Figure 11 (left). In the table, observe that the brute-force approach scales poorly since the search space grows exponentially with the number of parameters. In contrast, HERDER scales approximately linearly with the number of parameters, taking 2–4 times as long as typing and costing a single variant.

Next we consider how HERDER scales as the size of the source program increases. For our evaluation, we generate programs of increasing size but with a fixed number of 4 parameters. We do this by starting with a subset of the scimark FFT benchmark with four parameters, then increase the length of its core function by repeating its body multiple times. As before, we compare HERDER with a brute-force search and also with the time to type and cost the all-dynamic variant.

The results of the evaluation are shown in Figure 11 (right). In the table, observe that all three measurements scale approximately linearly with respect to the size of the program. The runtime of the brute-force approach is approximately 20 times the time to type and cost a single variant. This is as expected since the size of the search space is  $2^4 = 16$  and there is some overhead associated with generating the variants and identifying the cheapest configuration. More significantly, observe that HERDER also scales linearly and takes 2–3 times the single-variant time. This demonstrates that the size of the program does not induce an unexpected performance hit in our approach.

Together, the evaluations in Sections 6.2 and 6.3 demonstrate that HERDER can accurately and efficiently recommend performant migrations of gradually typed Python programs.

## 7 RELATED WORK

### 7.1 Static Cost Analysis

There has been substantial work on developing static analyses to infer bounds on the resource usage of programs. For example, automatic cost analyses have been defined over the ML family of languages that can bound the usage of time, memory, or energy [Hoffmann et al. 2017; Hoffmann and Hofmann 2010; Hoffmann and Shao 2015]. Similarly, our work is to statically measure the overhead of a certain resource, namely the number of casts. The methodology used in those works relies on rich type information and uses linear constraint solving to generate the resource

1226 bounds. Since gradually typed programs typically lack the rich type information to support these  
1227 methodologies, our approach is primarily syntax driven. While a main focus of those works is to  
1228 derive asymptotically tight upper bounds in their analysis for a given program, the main purpose  
1229 of our cost semantics is to compare the costs of different configurations by observing the worst  
1230 case behaviors possible for these programs. As Section 6.2 shows, our cost analysis serves this goal  
1231 well, accurately finding configurations with little overhead from inserted casts.

1232 Our approach of translating source programs into a corresponding cost language was inspired  
1233 by Danner et al. [2015, 2013]. Their work infers worst-case bounds on the number of evaluation  
1234 steps to execute a program; we infer worst case bounds on the number of casts performed in the  
1235 evaluation of a gradual program. We adapt the syntactic generation of approximations from their  
1236 work in order to produce costs for an appropriate formal model for languages like Reticulated.  
1237 Unlike theirs, our cost semantics is fully automatic. This means that we can apply our cost analysis  
1238 to existing Reticulated Python programs without users providing any additional information to  
1239 help us extract costs.

1240 Relational cost analysis [Çiçek et al. 2017] fulfills a very similar role to variational cost analysis.  
1241 In relational cost analysis the goal is to reason about the difference in cost between two similar  
1242 programs or two similar runs of the same program. For example, Çiçek et al. [2016] discuss  
1243 a relational cost analysis measuring the overhead of a rerun on a program with incremental  
1244 computation based on changes to the input. While a main goal of relational cost analyses is to relate  
1245 the cost of two runs of the same program with different inputs, that of variational cost analysis  
1246 in this work is to relate costs of two similar programs that differ in program structures. However,  
1247 variational cost analysis still bears much resemblance to unary typing in Çiçek et al. [2017] in the  
1248 way both reason about unrelated parts of two programs through the idea of sharing. It would be  
1249 interesting to see how precisely variational techniques can interact with relational cost analysis.

1250

## 1251 7.2 Gradual Typing Performance

1252 Since Takikawa et al. [2016] reported that sound gradually typed programs can incur massive  
1253 slowdown when mixing dynamic and static code, there have been several efforts to address the  
1254 problem.

1255 The Nom programming language [Muehlboeck and Tate 2017] addresses it by designing the  
1256 static and dynamic semantics for a gradually typed language from the ground up instead of adding  
1257 gradual typing to an existing language. This helps reduce the overhead of running casts. Nom’s  
1258 design allows programs to gain performance benefits as more types are added to a program.

1259 Bauman et al. [2017] focus on improving performance for Racket by using a tracing JIT compiler,  
1260 Pycket, with new representations for contracts that eliminate much of their overhead. Since Pycket  
1261 is a tracing JIT, it can effectively optimize untyped and typed boundaries upon observing them.  
1262 Though Pycket makes the interaction between typed and untyped code less costly, overhead still  
1263 remains at certain code boundaries. It would be interesting to see how variational cost analysis and  
1264 inference can be used to recommend types that eliminate the introduction of typed border crossing  
1265 in code that iterates heavily, which would in turn help Pycket.

1266 Richards et al. [2017] also help performance by designing intrinsic object contracts for implemen-  
1267 tations on a virtual machine, allowing the shape checks used by the VM’s JIT to act as the runtime  
1268 type-safety checks for gradual typing. It improves performance by not creating new allocations  
1269 when a contract is applied to an object that shares the shape of another object that already had the  
1270 contract applied to it. This design means that the performance of their approach is not determined  
1271 by the interaction of typed and untyped code, although they can still incur significant memory over-  
1272 head in certain typing configurations. It would be interesting to see if a variational cost semantics  
1273 can be created to reason about how different typings cause large memory consumption.

1274

1275 Overall, our approach is orthogonal to these approaches in quite a few ways. First, our method-  
1276 ology is a cost semantics and not a change to a runtime environment. Therefore its goal is not  
1277 solely optimizing performance but instead reasoning about how types affect performance and  
1278 can thus be used to support the different migration scenarios outlined in Section 1.2. Second, our  
1279 semantics is not appropriate for such specialized implementations, where the interaction of typed  
1280 and untyped code does not incur significant overhead. Instead it is intended for languages that  
1281 employ the traditional approach of translating gradually typed programs into untyped programs.  
1282 This makes it appropriate for the implementation of Typed Racket used in practice, which has  
1283 a traditional contract-based guarded semantics, or for many of the different semantics available  
1284 for Reticulated. Since many existing gradual languages work by translating typed programs into  
1285 untyped programs with proxies, we feel that our cost semantics is widely applicable. Finally, the  
1286 implementation of the cost analysis is relatively simple, with relatively small additions to the type  
1287 system and cast insertion rules. In contrast, the engineering effort for designing or modifying a  
1288 JIT compiler may take substantially more work, and codesigning a language's type system and  
1289 semantics is not an applicable strategy for many existing gradually typed languages.

1290 There are other approaches using type-based techniques to improve program performance.  
1291 Rastogi et al. [2012] use gradual type inference to help soundly migrate programs toward more  
1292 precise static types. In their system more precise typing is correlated with performance so they  
1293 want to infer as many types as possible while preserving soundness. Similarly our system can be  
1294 adapted to other type inference systems, such as the flow based one they present, so that we can  
1295 infer many possible most-static types, then find one optimizing static checking and performance.

1296 Finally, confined gradual typing [Allende et al. 2014] introduces constructs to explicitly manage  
1297 data flow between static and dynamic parts of the program, preventing costly boundary crossings.  
1298 Our cost analysis can be thought of as a way to reason about what parts of the programs contain these  
1299 costly boundary crossings and what type assignments can limit these crossings. This potentially  
1300 allows our system to help programmers by automatically finding and diagnosing these boundaries.

### 1301 7.3 Migrational and Variational Typing

1302 The design of the variational cost semantics in Section 5 builds on the technical machinery of  
1303 migrational typing [Campora et al. 2018], which in turn builds on the prior work on variational  
1304 typing [Chen et al. 2012, 2014]. This machinery is important to being able to efficiently infer types  
1305 and perform cost analysis for the entire space of potential migrations. Since migrating gradual  
1306 types [Campora et al. 2018] focused on the static type safety of gradual programs, it does not  
1307 contain variational cast insertion, which is necessary in this work in order to measure the costs of  
1308 the inserted casts. Aside from integrating the cost semantics, the implementation of HERDER also  
1309 required extending that work to support new constructs, such as loops, that are not present in the  
1310 original calculus but widely used in Python.

## 1311 8 CONCLUSION

1312 Gradual typing promises the reconciliation of static and dynamic typing. However, a major practical  
1313 limitation of current implementations is that the interfaces between dynamically and statically  
1314 typed code can have a huge runtime overhead. Different assignments of type annotations have a  
1315 significant affect on these costs, but it is hard to predict how to assign types to improve performance.  
1316 This leaves programmers stuck wondering how to migrate programs to types fulfilling the safety  
1317 and performance goals they desire for their program.

1318 To address this issue, we have presented a variational cost semantics for gradually typed programs  
1319 that approximates the runtime costs for all possible type configurations of a program. The cost  
1320 semantics provides a systematic way to create tools that help programmers identify performant  
1321  
1322  
1323

1324 migrations and understand how typing affects performance as they migrate their programs. We  
 1325 have implemented our semantics on top of Reticulated Python in HERDER, and our evaluation  
 1326 shows that HERDER is effective, efficient, and can be used to aid programmers in many different  
 1327 migration scenarios.

1328 Our approach is amenable to many gradually typed languages in which (partial) type inference  
 1329 is possible, and where inserted casts incur noticeable performance overhead. Overall, this makes  
 1330 variational cost analysis a viable approach to reasoning about the complex interaction of typing and  
 1331 performance during gradual program development. In the future, we will investigate how the ideas  
 1332 in this paper can help address the performance problem in other gradual language implementations,  
 1333 such as Typed Racket.  
 1334

1335

1336

## 1337 REFERENCES

- 1337 Esteban Allende, Johan Fabry, Ronald Garcia, and Éric Tanter. 2014. Confined Gradual Typing. *SIGPLAN Not.* 49, 10 (Oct. 2014), 251–270. <https://doi.org/10.1145/2714064.2660222>
- 1338 Spenser Bauman, Carl Friedrich Bolz-Tereick, Jeremy Siek, and Sam Tobin-Hochstadt. 2017. Sound Gradual Typing: Only  
 1339 Mostly Dead. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 54 (Oct. 2017), 24 pages. <https://doi.org/10.1145/3133878>
- 1340 John Campora, Sheng Chen, Martin Erwig, and Eric Walkingshaw. 2018. Migrating Gradual Types. In *Proceedings of the  
 1341 45th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL '18)*. ACM, New York, NY, USA.
- 1342 Ezgi Çiçek, Gilles Barthe, Marco Gaboardi, Deepak Garg, and Jan Hoffmann. 2017. Relational Cost Analysis. *SIGPLAN Not.*  
 1343 52, 1 (Jan. 2017), 316–329. <https://doi.org/10.1145/3093333.3009858>
- 1344 Ezgi Çiçek, Zoe Paraskevopoulou, and Deepak Garg. 2016. A Type Theory for Incremental Computational Complexity with  
 1345 Control Flow Changes. *SIGPLAN Not.* 51, 9 (Sept. 2016), 132–145. <https://doi.org/10.1145/3022670.2951950>
- 1346 Sheng Chen, Martin Erwig, and Eric Walkingshaw. 2012. An Error-tolerant Type System for Variational Lambda Calculus.  
 1347 In *Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming (ICFP '12)*. ACM, New York,  
 1348 NY, USA, 29–40. <https://doi.org/10.1145/2364527.2364535>
- 1349 Sheng Chen, Martin Erwig, and Eric Walkingshaw. 2014. Extending Type Inference to Variational Programs. *ACM Trans.  
 1350 Program. Lang. Syst.* 36, 1, Article 1 (March 2014), 54 pages. <https://doi.org/10.1145/2518190>
- 1351 Sheng Chen, Martin Erwig, and Eric Walkingshaw. 2016. A Calculus for Variational Programming. In *European Conf. on  
 1352 Object-Oriented Programming (ECOOP) (LIPICs)*, Vol. 56. 6:1–6:26.
- 1353 Norman Danner, Daniel R. Licata, and Ramyaa Ramyaa. 2015. Denotational Cost Semantics for Functional Languages with  
 1354 Inductive Types. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP  
 1355 2015)*. ACM, New York, NY, USA, 140–151. <https://doi.org/10.1145/2784731.2784749>
- 1356 Norman Danner, Jennifer Paykin, and James S. Royer. 2013. A Static Cost Analysis for a Higher-order Language. In  
 1357 *Proceedings of the 7th Workshop on Programming Languages Meets Program Verification (PLPV '13)*. ACM, New York, NY,  
 1358 USA, 25–34. <https://doi.org/10.1145/2428116.2428123>
- 1359 Martin Erwig and Eric Walkingshaw. 2011. The Choice Calculus: A Representation for Software Variation. *ACM Trans.  
 1360 Softw. Eng. Methodol.* 21, 1, Article 6 (Dec. 2011), 27 pages. <https://doi.org/10.1145/2063239.2063245>
- 1361 Martin Erwig and Eric Walkingshaw. 2013. Variation Programming with the Choice Calculus. In *Generative and Transfor-  
 1362 mational Techniques in Software Engineering IV (GTTSE 2011), Revised and Extended Papers (LNCS)*, Vol. 7680. 55–99.
- 1363 Ronald Garcia and Matteo Cimini. 2015. Principal Type Schemes for Gradual Programs. In *Proceedings of the 42Nd Annual  
 1364 ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '15)*. ACM, New York, NY, USA,  
 1365 303–315. <https://doi.org/10.1145/2676726.2676992>
- 1366 Ronald Garcia, Alison M. Clark, and Éric Tanter. 2016. Abstracting Gradual Typing. In *Proceedings of the 43rd Annual ACM  
 1367 SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '16)*. ACM, New York, NY, USA, 429–442.  
 1368 <https://doi.org/10.1145/2837614.2837670>
- 1369 Jan Hoffmann, Ankush Das, and Shu-Chun Weng. 2017. Towards Automatic Resource Bound Analysis for OCaml. In  
 1370 *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL 2017)*. ACM, New York,  
 1371 NY, USA, 359–373. <https://doi.org/10.1145/3009837.3009842>
- 1372 Jan Hoffmann and Martin Hofmann. 2010. Amortized Resource Analysis with Polynomial Potential - A Static Inference  
 of Polynomial Bounds for Functional Programs. In *Proceedings of the 19th European Symposium on Programming  
 (ESOP'10) (Lecture Notes in Computer Science)*, Vol. 6012. Springer, 287–306.
- Jan Hoffmann and Zhong Shao. 2015. Automatic Static Cost Analysis for Parallel Programs. In *Proceedings of the 24th  
 European Symposium on Programming on Programming Languages and Systems - Volume 9032*. Springer-Verlag New York,  
 Inc., New York, NY, USA, 132–157. [https://doi.org/10.1007/978-3-662-46669-8\\_6](https://doi.org/10.1007/978-3-662-46669-8_6)

1372

- 1373 Fabian Muehlboeck and Ross Tate. 2017. Sound Gradual Typing is Nominally Alive and Well. In *OOPSLA*. ACM, New York,  
 1374 NY, USA. <https://doi.org/10.1145/3133880>
- 1375 Aseem Rastogi, Avik Chaudhuri, and Basil Hosmer. 2012. The Ins and Outs of Gradual Type Inference. *SIGPLAN Not.* 47, 1  
 1376 (Jan. 2012), 481–494. <https://doi.org/10.1145/2103621.2103714>
- 1377 Gregor Richards, Ellen Artica, and Alexi Turcotte. 2017. The VM Already Knew That: Leveraging Compile-time Knowledge  
 1378 to Optimize Gradual Typing. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 55 (Oct. 2017), 27 pages. <https://doi.org/10.1145/3133879>
- 1379 Jeremy Siek, Ronald Garcia, and Walid Taha. 2009. Exploring the Design Space of Higher-Order Casts. In *Programming  
 1380 Languages and Systems*, Giuseppe Castagna (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 17–31.
- 1381 Jeremy G. Siek and Walid Taha. 2006. Gradual Typing for Functional Languages. In *IN SCHEME AND FUNCTIONAL  
 1382 PROGRAMMING WORKSHOP*. 81–92.
- 1383 Jeremy G. Siek and Manish Vachharajani. 2008. Gradual Typing with Unification-based Inference. In *Proceedings of the 2008  
 1384 Symposium on Dynamic Languages (DLS '08)*. ACM, New York, NY, USA, Article 7, 12 pages. <https://doi.org/10.1145/1408681.1408688>
- 1385 Asumu Takikawa, Daniel Feltey, Ben Greenman, Max S. New, Jan Vitek, and Matthias Felleisen. 2016. Is Sound Gradual  
 1386 Typing Dead?. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming  
 1387 Languages (POPL '16)*. ACM, New York, NY, USA, 456–468. <https://doi.org/10.1145/2837614.2837630>
- 1388 Sam Tobin-Hochstadt and Matthias Felleisen. 2006. Interlanguage Migration: From Scripts to Programs. In *Companion to  
 1389 the 21st ACM SIGPLAN Symposium on Object-oriented Programming Systems, Languages, and Applications (OOPSLA '06)*.  
 ACM, New York, NY, USA, 964–974. <https://doi.org/10.1145/1176617.1176755>
- 1390 Sam Tobin-Hochstadt, Matthias Felleisen, Robert Findler, Matthew Flatt, Ben Greenman, Andrew M. Kent, Vincent St-Amour,  
 1391 T. Stephen Strickland, and Asumu Takikawa. 2017. Migratory Typing: Ten Years Later. In *2nd Summit on Advances in  
 1392 Programming Languages (SNAPL 2017) (Leibniz International Proceedings in Informatics (LIPIcs))*, Benjamin S. Lerner,  
 1393 Rastislav Bodik, and Shriram Krishnamurthi (Eds.), Vol. 71. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl,  
 Germany, 17:1–17:17. <https://doi.org/10.4230/LIPIcs.SNAPL.2017.17>
- 1394 Guido van Rossum, Jukka Lehtosalo, and Lukasz Langa. 2014. PEP 484 – Type Hints. (2014). <https://www.python.org/dev/peps/pep-0484/#rationale-and-goals>
- 1395 Michael M. Vitousek, Andrew M. Kent, Jeremy G. Siek, and Jim Baker. 2014. Design and Evaluation of Gradual Typing for  
 1396 Python. (2014), 45–56.
- 1397 Michael M. Vitousek, Cameron Swords, and Jeremy G. Siek. 2017. Big Types in Little Runtime: Open-world Soundness and  
 1398 Collaborative Blame for Gradual Type Systems. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of  
 1399 Programming Languages (POPL 2017)*. ACM, New York, NY, USA, 762–774. <https://doi.org/10.1145/3009837.3009849>

## 1400 A PROOFS

1401 Before we start on the proofs for the lemmas and theorems from Section 4 and 5, we assert some  
 1402 lemmas whose proofs are straightforward. First, is a lemma stating that variable references are  
 1403 preserved by the cast insertion procedure.  
 1404

1405 LEMMA A.1 (TRANSLATION PRESERVES VARIABLES). *If  $x \in e_s$  and  $\Gamma \vdash_G e_s \rightsquigarrow e_t : G$  then  $x \in e_t$ .*

1406 The proof of this lemma follows inductively from the cast insertion rules in Figure 5. This next  
 1407 lemma states that our transformation for cost analysis also preserves variable references.  
 1408

1409 LEMMA A.2 (TRANSFORMATION PRESERVES VARIABLES). *If  $x \in e_t$  then  $x \in \searrow (e_t)$ .*

1410 The proof of this lemma is from case analysis on Figure 6, and we assume that any inserted  
 1411 variables from the procedure do not capture existing variables. Finally, we assume that an inversion  
 1412 rule exists for our cost semantics. We only state the cases used in proofs of lemmas and theorems.  
 1413

1414 LEMMA A.3 (INVERSION OF THE COST RELATION).

1415 (1) *If  $\Phi \vdash_c \lambda x. e_t \mid A$  then  $A = (0, \Lambda x. A')$  for some  $A'$  and where  $\Phi, x \mapsto x \vdash_c e_t \mid A'$ .*

1416 The proof of this immediately follows from the definition of the cost semantics.  
 1417 With these lemmas in hand, we can provide the proof of Lemma 4.1.  
 1418  
 1419  
 1420  
 1421

1422 PROOF OF LEMMA 4.1. We are given:

$$1423 \quad \Phi; \Gamma \vdash_{GC} \lambda x. e_s \rightsquigarrow \lambda x. e_t : G \mid A$$

1425 The beginning of the proof begins by using Lemma A.1 and Lemma A.2 to derive  $x \in \searrow (\lambda x. e_t)$ .  
1426 Note that  $\searrow (\lambda x. e_t) = \lambda x. \searrow (e_t)$ . Next, we use Lemma A.3 to derive:

$$1427 \quad A = (0, \Lambda x. A') \quad \Phi, x \mapsto x \vdash_c e_t \mid (0, \Lambda x. A')$$

1429 Since we know that  $x \in e_t$ , then VAR must have been used when producing the cost for  $e_t$  to derive:

$$1430 \quad \Phi, x \mapsto x \vdash_c x \mid x$$

1432 Thus, we know  $x \in A'$  and consequently  $x \in (0, \Lambda x. A')$ .

1433  $\square$

1434 For the proof of Lemma 4.2, we will only consider the most interesting cases, namely those  
1435 involving references.

1437 PROOF OF LEMMA 4.2. This proof will follow by induction over the rules in Figure 5, also using  
1438 the rules in Figure 6 and Figure 7. For each case, we are given the following:

$$1439 \quad \Phi; \Gamma \vdash_{GC} e_s \rightsquigarrow e_t : G \mid (C, P) \quad x : G' \quad x \in \text{vars}(e_s) \quad \Phi; \Gamma \vdash_{GC} e'_s \rightsquigarrow e'_t : G' \mid (C', P')$$

1441 Case REF: We have  $e_s = \mathbf{ref} \ e_{s1}$  and  $\searrow (e_t) = \searrow (\mathbf{ref} \ e_{t1}) = \mathbf{let} \ z = \searrow (e_{t1}) \ \mathbf{in} \ \lambda y. z$ . This means  
1442 that the potential must be  $P = \Lambda y. (0, P_1)$  where  $P_1$  is the potential for  $\searrow (e_{t1})$ . Our induction  
1443 hypotheses are:

$$1444 \quad \Phi; \Gamma \vdash_{GC} e_{s1} \rightsquigarrow e_{t1} : G_1 \mid (C_1, P_1) \quad \Phi; \Gamma \vdash_{GC} [e'_s/x]e_{s1} \rightsquigarrow e'_{t1} : G_1 \mid (C'_1, [P'/x]P_1)$$

1445 From the induction hypothesis we directly have:  $\Phi \vdash_c \searrow (e'_{t1}) \mid (C'_1, [P'/x]P_1)$  and we know  
1446  $z$  will be assigned  $\searrow (e'_{t1})$ . We can derive (assuming  $y \notin \text{vars}(\searrow (e'_{t1}))$ ):

$$1447 \quad \Phi, y \mapsto y, z \mapsto [P'/x]P_1 \vdash_c z \mid (0, [P'/x]P_1)$$

1449 Since we know that  $\searrow (\mathbf{ref} \ e_{t1}) = \mathbf{let} \ z = \searrow (e_{t1}) \ \mathbf{in} \ \lambda y. z$ , we can similarly use ABS on the  
1450 substituted program to derive:

$$1451 \quad \Phi, z \mapsto [P'/x]P_1 \vdash_c \lambda y. z \mid (0, \Lambda y. (0, [P'/x]P_1))$$

1453 Finally, we can use LET to derive:

$$1454 \quad \Phi \vdash_c \mathbf{let} \ z = \searrow (e'_{t1}) \ \mathbf{in} \ \lambda y. z \mid (C'_1, \Lambda y. (0, [P'/x]P_1))$$

1456 Case ASSIGN: We have  $e_s = e_{s1} := e_{s2}$  and  $\searrow (e_t) = \searrow (e_{t1} := e_{t2}) = \searrow (e_{t1}) \ \searrow (e_{t2})$ . This means  
1457 that the potential must be  $P = (P_1 \ P_2)$  where  $P_1$  is the potential of  $\searrow (e_{t1})$  and  $P_2$  is the  
1458 potential of  $\searrow (e_{t2})$ . Our induction hypotheses are:

$$1460 \quad \Phi; \Gamma \vdash_{GC} e_{s1} \rightsquigarrow e_{t1} : G_1 \mid (C_1, P_1) \quad \Phi; \Gamma \vdash_{GC} e_{s2} \rightsquigarrow e_{t2} : G_2 \mid (C_2, P_2)$$

$$1461 \quad \Phi; \Gamma \vdash_{GC} [e'_s/x]e_{s1} \rightsquigarrow e'_{t1} : G_1 \mid (C'_1, [P'/x]P_1) \quad \Phi; \Gamma \vdash_{GC} [e'_s/x]e_{s2} \rightsquigarrow e'_{t2} : G_2 \mid (C'_2, [P'/x]P_2)$$

1462 Directly from our induction hypotheses we have  $\Phi \vdash_c e'_{t1} \mid (C'_1, [P'/x]P_1)$  and  $\Phi \vdash_c$   
1463  $e'_{t2} \mid (C'_2, [P'/x]P_2)$ . Notice that due to the structure of  $\searrow$  we have  $\searrow (e'_{t1} \ e'_{t2}) = \searrow (e'_{t1}) \ \searrow$   
1464  $(e'_{t2})$ . We can then use APP with these hypotheses to conclude:

$$1465 \quad \Phi \vdash_c \searrow (e'_{t1} \ e'_{t2}) \mid C'_1 + C'_2 \oplus ([P'/x]P_1 \ [P'/x]P_2)$$

1467 From the definition  $\oplus$ , the potential of this approximation is solely formed from  
1468  $([P'/x]P_1 \ [P'/x]P_2)$  and via a usual definition of substitution, we see that the potential  
1469 of  $[e'_s/x]e_s$  is  $[P'/x](P_1 \ P_2)$

1470

Case Deref: We have  $e_s = ! e_{s1}$  and  $\searrow (e_t) = \searrow (! e_{t1}) = \searrow (e_{t1})()$ . This means that the potential must be  $P = (P_1 0)$  where  $(P_1)$  is the potential for  $\searrow (e_{t1})$  and  $0$  is the potential for  $()$ . Our induction hypotheses are:

$$\Phi; \Gamma \vdash_{GC} e_{s1} \rightsquigarrow e_{t1} : G_1 \mid (C_1, P_1) \quad \Phi; \Gamma \vdash_{GC} [e_s'/x]e_{s1} \rightsquigarrow e'_{t1} : G_1 \mid (C'_1, [P'/x]P_1)$$

From the induction hypothesis we have  $\Phi \vdash_c e'_{t1} \mid [P'/x]P_1$ . We know that  $\searrow (! e'_{t1}) = \searrow (e'_{t1})()$  and by definition we know  $0 = [e_s'/x]0$ . With this, we can use APP to derive:

$$\Phi \vdash_c \searrow (e'_{t1}) () \mid C'_1 \oplus [e_s'/x]P_1 [e_s'/x]0$$

We can then use the definition of substitution to see that the potential for  $[e_s'/x]e_{s1}$  is  $[e_s'/x](P_1 0)$ .

The rest of the cases follow easily from the induction hypothesis, since the transformation algorithm does not change the structure of the top level target expressions.  $\square$

For the proof Theorem 4.3, we will discuss the most interesting terms only, namely function application and recursive let bindings. To state termination, we informally use an evaluation relation without presenting a full dynamic semantics for the language, since the input language is fairly standard. The only aspect of termination that complicates the combined cast insertion and costing process is the actual cost analysis derivations in  $\vdash_c$ , so we ignore presenting termination arguments for the cast insertion and transformation algorithm.

PROOF OF THEOREM 4.3. The proof proceeds by induction over the target language expressions, with the induction hypothesis  $e_t \Downarrow v \Rightarrow \Phi \vdash_c \searrow (e_t) \mid A$ .

Case APP: We have the following induction hypotheses:

$$e_{t1} \Downarrow \lambda x. e'_{t1} \Rightarrow \Phi \vdash_c \searrow (e_{t1}) \mid (C_1, \Lambda x. A'_1) \quad e_{t2} \Downarrow v_2 \Rightarrow \Phi \vdash_c \searrow (e_{t2}) \mid (C_2, P_2)$$

From the statement of the theorem, we know that  $e_{t1} e_{t2} \Downarrow v$ . This implies that:

$$\lambda x. e'_{t1} v_2 = [v_2/x]e'_{t1} = v$$

To generate the cost for  $\searrow (e_{t1}) \searrow (e_{t2})$  we use APP to derive:

$$\Phi \vdash_c \searrow (e_{t1}) \searrow (e_{t2}) \mid C_1 + C_2 \oplus (\Lambda x. A'_1 P_2)$$

The only term that causes concerns about non-termination comes from  $(\Lambda x. A'_1 P_2)$ . By definition, this term equals:

$$[P_2/x]A'_1$$

Since each subterm of  $[v_2/x]e'_{t1}$  evaluates, we can apply the induction hypothesis to each corresponding cost subterm generated in  $[P_2/x]A'_1$  and consequently the evaluation of this cost analysis term terminates.

Case LETREC: We assume that letrec is defined in terms of a fixed point operator, that is used to eventually evaluate  $e_{t1}$  to a value. We then have the following induction hypotheses:

$$e_{t1} \Downarrow v_1 \Rightarrow \Phi \vdash_c \searrow (e_{t1}) \mid (C_1, P_1) \quad [v_1/x]e_{t2} \Downarrow v_2 \Rightarrow \Phi, x \mapsto l \cdot n^l \cdot P' \vdash_c \searrow (e_{t2}) \mid (C_2, P_2)$$

By definition we have  $\searrow (\mathbf{letrec} \ x = e_{t1} \ \mathbf{in} \ e_{t2}) = \mathbf{letrec} \ x = \searrow (e_{t1}) \ \mathbf{in} \ \searrow (e_{t2})$ . Directly from the induction hypotheses, we know that producing  $P_1$  terminates. The other premises for LETREC calculate  $S$ ,  $n$ ,  $P'$ , and  $l$  and none of these calculations can cause non-termination. From this and the induction hypotheses, we know that production of  $(C_2, P_2)$  completes. Consequently, know LETREC terminates and derives:

$$\Phi \vdash_c \searrow (\mathbf{letrec} \ x = e_{t1} \ \mathbf{in} \ e_{t2}) \mid (C_1 + C_2, P_2)$$

$\square$

In the following proofs for theorem 5.1 and theorem 5.2 we will only relate the typing and cast insertion of the two systems and will ignore relating the costs, because we would need to invoke the rules in Figure 7, and overall the structure arguments using variations and decisions are similar to what is done for cast insertion. For the following theorems we only consider cases for more interesting rules, namely variable references, function abstractions, and function applications. The machinery for other cases follow from what is shown for these cases.

**PROOF OF THEOREM 5.1.** This theorem's proof follows from induction over the rules in Figure 5. The general structure of the induction hypothesis is:  $K\Gamma \vdash_G e_s \rightsquigarrow e_t : G \Rightarrow \Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega$  where  $G = [M]_\delta$  and  $K = [\Omega]_\delta$ .

**Case VAR:** Since we have  $x \mapsto G \in \Gamma$  and we can directly use this environment to derive  $\Gamma \vdash^v x \rightsquigarrow x : G \mid \emptyset$  where  $K = \emptyset$  and  $\delta = \emptyset$ .

**Case ABS:** We first focus on the case where the abstraction contains a static annotation.

Subcase  $G=T$  Given that we know:

$$K\Gamma, x \mapsto T \vdash_G e_s \rightsquigarrow e_t : G \quad \Gamma \vdash_G \lambda x : T.e_s \rightsquigarrow \lambda x : T.e_t : T \rightarrow G$$

Our induction hypotheses are:

$$\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega \quad [M]_\delta = G \quad [\Omega]_\delta = K \quad [e_t^v]_\delta = e_t$$

With the induction hypotheses it's very easy to see that we can use **ABSV** to conclude:

$$\Gamma \vdash^v \lambda x : T.e_s \rightsquigarrow \lambda x : T.e_t^v : T \rightarrow M \mid \Omega$$

We can see that from the induction hypothesis that  $[T \rightarrow M]_\delta = [T]_\delta \rightarrow [M]_\delta = T \rightarrow G$  and we also know  $[\Omega]_\delta = K$ . Finally, we need to show  $[\lambda x : T.e_t^v]_\delta = \lambda x : T.e_t$ . We already know from the induction hypotheses that  $[e_t^v]_\delta = [e_t]_\delta$ , and so it's clear that we can expand the selection on the entire abstraction to get the desired result.

Subcase  $G=\text{Dyn}$ : Given that we know:

$$K\Gamma, x \mapsto \text{Dyn} \vdash_G e_s \rightsquigarrow e_t : G_2 \quad \Gamma \vdash_G \lambda x : \text{Dyn}.e_s \rightsquigarrow \lambda x : \text{Dyn}.e_t : G_1 \rightarrow G_2 \quad K\Gamma \vdash_G x \rightsquigarrow x : G_1$$

Our induction hypotheses are:

$$\Gamma \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega \quad [M]_\delta = G \quad [\Omega]_\delta = K \quad [e_t^v]_\delta = e_t$$

From the induction hypotheses, we can directly use **ABSDYNV** to conclude:

$$\Gamma \vdash^v \lambda x : \text{Dyn}.e_s \rightsquigarrow \lambda x : \text{Dyn}.e_t : d\langle \text{Dyn}, V \rangle \rightarrow M \mid \Omega \cup \{x \mapsto d\langle \text{Dyn}, V \rangle\}$$

From the induction hypothesis, we already know  $[M]_\delta = G$ . We just need to show  $[d\langle \text{Dyn}, V \rangle]_\delta = G_1$ . First, notice that since  $[\Omega]_\delta = K$ , we have  $K(x) = [\Omega(x)]_\delta = [d\langle \text{Dyn}, V \rangle]_\delta$ . Finally, we need to show  $[\lambda x : \text{Dyn}.e_t^v]_\delta = \lambda x : \text{Dyn}.e_t$ . We already know from the induction hypotheses that  $[e_t^v]_\delta = [e_t]_\delta$ , and so it's clear that we can expand the selection on the entire abstraction to get the desired result.

**Case APP:** We are given the following:

$$K\Gamma \vdash_G e_{s1} \rightsquigarrow e_{t1} : G_1 \quad K\Gamma \vdash_G e_{s2} \rightsquigarrow e_{t2} : G_2 \quad \text{dom}(G_1) \sim G_2$$

We have the following induction hypotheses:

$$\begin{aligned} \Gamma \vdash^v e_{1s} \rightsquigarrow e_{1t}^v : M_1 \mid \Omega_1 \quad [M_1]_\delta = G_1 \quad [\Omega_1]_\delta = K \quad [e_{1t}^v]_\delta = e_{1t} \\ \Gamma \vdash^v e_{2s} \rightsquigarrow e_{2t}^v : M_2 \mid \Omega_2 \quad [M_2]_\delta = G_2 \quad [\Omega_2]_\delta = K \quad [e_{2t}^v]_\delta = e_{2t} \end{aligned}$$

The main difficulty we face is showing that there is always some arbitrary  $M_1$  and  $M_2$  such that  $M_1 \approx M_2$ . We always have  $G_1 \sim G_2$  and so at the very least, we can have  $M_1 = G_1$  and  $M_2 = G_2$  and then we have  $M_1 \approx M_2$ . The rest of the results for this case follow straightforwardly with  $\delta = \emptyset$ . Otherwise, assume that this is not the case but we have  $M_1 \approx M_2$  and from our



induction hypotheses we have  $\lfloor M_1 \rfloor_\delta = G_1$  and  $\lfloor M_2 \rfloor_\delta = G_2$ . We can use APPV to conclude:

$$\Gamma \vdash^v e_{s1} e_{s2} \rightsquigarrow \llbracket \text{dom}(M) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket e_{t1}^v \llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket e_{t2}^v : \text{cod}(M_1) \mid \Omega_1 \cup \Omega_2$$

Since we already had  $\lfloor \Omega_1 \rfloor_\delta = K$  and  $\lfloor \Omega_2 \rfloor_\delta = K$ , it is clear that

$$\lfloor \Omega_1 \cup \Omega_2 \rfloor_\delta = K$$

Moreover, since we already had  $\lfloor M_1 \rfloor_\delta = G_1$ , it's clear that we have:

$$\lfloor \text{cod}(M_1) \rfloor_\delta = \text{cod}(G_1)$$

Additionally, the induction hypotheses make it clear that  $\llbracket \text{dom}(M_1) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket_\delta = \llbracket \text{dom}(G_1) \rightarrow \text{cod}(G_2) \Leftarrow G_1 \rrbracket$  and  $\llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket_\delta = \llbracket \text{dom}(G_1) \Leftarrow G_2 \rrbracket$ , since the selections on the casts eliminate the variational types. With the induction hypotheses for the target expressions, we can conclude:

$$\begin{aligned} \llbracket \text{dom}(M_1) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket e_{t1}^v \llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket e_{t2}^v \rfloor_\delta = \\ \llbracket \text{dom}(G_1) \rightarrow \text{cod}(G_1) \Leftarrow G_1 \rrbracket e_{t1} \llbracket \text{dom}(G_1) \Leftarrow G_2 \rrbracket e_{t2} \end{aligned}$$

□

PROOF OF THEOREM 5.2. This theorem's proof follows from induction over the rules in Figure 9.

Case VAR: The proof of this case is straightforward in that  $\lfloor \Omega \rfloor_\delta$  is used for the configuration, and if  $x$  is not bound in  $\Omega$  then it was already in  $\Gamma$ .

Case ABS<sub>DYNV</sub>: We are given the following:

$$\begin{aligned} \Gamma, x \mapsto d\langle \text{Dyn}, V \rangle \vdash^v e_s \rightsquigarrow e_t^v : M \mid \Omega' \\ \Gamma \vdash^v \lambda x : \text{Dyn}.e_s \rightsquigarrow \lambda x : \text{Dyn}.e_t^v : d\langle \text{Dyn}, V \rangle \rightarrow M \mid \Omega' \cup \{x \mapsto d\langle \text{Dyn}, V \rangle\} \end{aligned}$$

Let  $\Omega = \Omega' \cup \{x \mapsto d\langle \text{Dyn}, V \rangle\}$ . We have the following induction hypothesis:

$$\lfloor \Omega \rfloor_\delta \Gamma, x \mapsto \text{Dyn} \vdash_G e_s \rightsquigarrow e_t^v : \lfloor M \rfloor_\delta$$

Since we know that  $x \in \lfloor \Omega \rfloor_\delta$ , we know that the type of  $x$  will be updated to use  $\lfloor \Omega(x) \rfloor_\delta = \lfloor d\langle \text{Dyn}, V \rangle \rfloor_\delta$ . We can then use ABS to conclude the following:

$$\lfloor \Omega \rfloor_\delta \Gamma \vdash_G \lambda x : \text{Dyn}.e_s \rightsquigarrow \lambda x : \text{Dyn}.\lfloor e_t^v \rfloor_\delta : \lfloor d\langle \text{Dyn}, V \rangle \rfloor_\delta \rightarrow \lfloor M \rfloor_\delta$$

It is clear that  $\lambda x : \text{Dyn}.\lfloor e_t^v \rfloor_\delta = \lfloor \lambda x : \text{Dyn}.e_t^v \rfloor_\delta$  and  $\lfloor d\langle \text{Dyn}, V \rangle \rfloor_\delta \rightarrow \lfloor M \rfloor_\delta = \lfloor d\langle \text{Dyn}, V \rangle \rightarrow M \rfloor_\delta$ .

Case ABS<sub>V</sub>: This case proceeds quite similarly to the case for ABS<sub>DYNV</sub> except that the steps involving the configuration is unnecessary.

Case APPV: We are given the following:

$$\begin{aligned} \Gamma \vdash^v e_{s1} \rightsquigarrow e_{t1}^v : M_1 \mid \Omega_1 \quad \Gamma \vdash^v e_{s2} \rightsquigarrow e_{t2}^v : M_2 \mid \Omega_2 \\ \Gamma \vdash^v e_{s1} e_{s2} \rightsquigarrow \llbracket \text{dom}(M_1) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket e_{t1}^v \llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket e_{t2}^v : \text{cod}(M_1) \mid \Omega_1 \cup \Omega_2 \end{aligned}$$

We have the following induction hypotheses:

$$\lfloor \Omega_1 \rfloor_\delta \Gamma \vdash_G e_{s1} \rightsquigarrow \lfloor e_{t1}^v \rfloor_\delta : \lfloor M_1 \rfloor_\delta \quad \lfloor \Omega_2 \rfloor_\delta \Gamma \vdash_G e_{s2} \rightsquigarrow \lfloor e_{t2}^v \rfloor_\delta : \lfloor M_2 \rfloor_\delta$$

Note that since we know that we have:

$$\Gamma \vdash^v e_{s1} e_{s2} \rightsquigarrow \llbracket \text{dom}(M_1) \rightarrow \text{cod}(M_1) \Leftarrow M_1 \rrbracket e_{t1}^v \llbracket \text{dom}(M_1) \Leftarrow M_2 \rrbracket e_{t2}^v : \text{cod}(M_1) \mid \Omega_1 \cup \Omega_2$$

we also know that  $\text{dom}(M_1) \approx M_2$ . The definition of compatability implies that we have:

$$\lfloor \text{dom}(M_1) \rfloor_\delta \sim \lfloor M_2 \rfloor_\delta$$

1618 This allows us to use APP to conclude:

1619

$$1620 \llbracket (\Omega_1 \cup \Omega_2) \rrbracket_\delta \Gamma \vdash_G e_{s1} e_{s2} \rightsquigarrow$$

1621

$$1622 \llbracket \llbracket \text{dom}(M_1) \rrbracket_\delta \rightarrow \llbracket \text{cod}(M_1) \rrbracket_\delta \Leftarrow \llbracket M_1 \rrbracket_\delta \rrbracket \llbracket e_{t1}^v \rrbracket_\delta \llbracket \llbracket \text{dom}(M_1) \rrbracket_\delta \Leftarrow \llbracket M_2 \rrbracket_\delta \rrbracket \llbracket e_{t2}^v \rrbracket_\delta : \llbracket \text{cod}(M_1) \rrbracket_\delta$$

1623

□

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666