ADVENTURES IN TIME AND SPACE

NORMAN DANNER^a AND JAMES S. ROYER^b

^a Department of Mathematics and Computer Science, Wesleyan University, Middletown, CT 06459 USA

 $e\text{-}mail\ address:$ ndanner@wesleyan.edu

 b Department of Electrical Engineering and Computer Science, Syracuse University, Syracuse, NY 13210 USA

e-mail address: royer@ecs.syr.edu

ABSTRACT. This paper investigates what is essentially a call-by-value version of PCF under a complexity-theoretically motivated type system. The programming formalism, ATR, has its first-order programs characterize the polynomial-time computable functions, and its second-order programs characterize the type-2 basic feasible functionals of Mehlhorn and of Cook and Urquhart. (The ATR-types are confined to levels 0, 1, and 2.) The type system comes in two parts, one that primarily restricts the sizes of values of expressions and a second that primarily restricts the time required to evaluate expressions. The size-restricted part is motivated by Bellantoni and Cook's and Leivant's implicit characterizations of polynomial-time. The time-restricting part is an affine version of Barber and Plotkin's DILL. Two semantics are constructed for ATR. The first is a pruning of the naïve denotational semantics for ATR. This pruning removes certain functions that cause otherwise feasible forms of recursion to go wrong. The second semantics is a model for ATR's time complexity relative to a certain abstract machine. This model provides a setting for complexity recurrences arising from ATR recursions, the solutions of which yield second-order polynomial time bounds. The time-complexity semantics is also shown to be sound relative to the costs of interpretation on the abstract machine.

1. INTRODUCTION

A Lisp programmer knows the value of everything, but the cost of nothing. — Alan Perlis

Perlis' quip is an overstatement—but not by much. Programmers in functional (and objectoriented) languages have few tools for reasoning about the efficiency of their programs. Almost all tools from traditional analysis of algorithms are targeted toward roughly the first-order fragment of C. What tools there are from formal methods are interesting, but piecemeal and preliminary.

DOI:10.2168/LMCS-3 (1:9) 2007

²⁰⁰⁰ ACM Subject Classification: F.3.3, F.1.3, F.3.2.

Key words and phrases: type systems, compositional semantics, implicit computational complexity, higher-type computation, basic feasible functionals.

This paper is an effort to fill in part of the puzzle of how to reason about the efficiency of programs that involve higher types. Our approach is, roughly, to take PCF and its conventional denotational semantics [Plo77, Win93] and, using types, restrict the language and its semantics to obtain a higher-type "feasible fragment" of both PCF and the PCF computable functions. Our notion of higher-type feasibility is based on the *basic feasible functionals* (BFFs) [CU93, Meh76], a higher-type analogue of polynomial-time computability, and Kapron and Cook's [KC96] machine-based characterization of the typelevel 2 BFFs.¹ Using a higher-type notion of computational complexity as the basis of our work provides a connection to the basic notions and tools of traditional analysis of algorithms (and their lifts to higher types). Using types to enforce feasibility constraints on PCF provides a connection to much of the central work in formal methods.

Our approach is in contrast to the work of [BNS00, Hof03, LM93] which also involves higher-type languages and types that guarantee feasibility. Those programming formalisms are feasible in the sense that they have polynomial-time normalization properties and that the type-level 1 functions expressible by these systems are guaranteed to be (ordinary) polynomial-time computable. The higher-type constructions of these formalisms are essentially aides for type-level 1 polynomial-time programming. As of this writing, there is scant analysis of what higher-type functions these systems compute.²

For a simple example of a feasible higher-type function, consider $C: (\mathbb{N} \to \mathbb{N}) \to (\mathbb{N} \to \mathbb{N}) \to (\mathbb{N} \to \mathbb{N})$ with $C f g = f \circ g$. (Convention: \mathbb{N} is always interpreted as $\{0, 1\}^*$, i.e., 0-1-strings.) In our setting, a reasonable implementation of C has a run-time bound that is a second-order polynomial (see §2.12) in the complexities of arbitrary f and g; in particular, if f and g are polynomial-time computable, so is C f g. Such a combinator C can be considered as part of the "feasible glue" of a programming environment—when used with other components, its complexity contribution is (higher-type) polynomially-bounded in terms of the complexity of the other components and the combined complexity can be expressed in a natural, compositional way. More elaborate examples of feasible functionals include many of the deterministic black-box constructions from cryptography. Chapter 3 in Goldreich [Gol01] has detailed examples, but a typical such construct takes one pseudorandom generator, g, and builds another, \tilde{g} , with better cryptographic properties but with not much worse complexity properties than the original g. Note that these g's and \tilde{g} 's may be feasible only in a probabilistic- or circuit-complexity sense.³

While our notion of feasibility is based on the BFFs, our semantic models allow our formalism to compute more than just the standard BFFs. For example, consider prn: ($N \rightarrow$

¹Mehlhorn [Meh76] originally discovered the class of type-2 BFFs in the mid-1970s. Later Cook and Urquhart [CU93] independently discovered this class and extended it to all finite types over the full settheoretic hierarchy. **N.B.** If one restricts attention to continuous models, then starting at type-level 3 there are alternative notions of "higher-type polynomial-time" [IKR02]. Dealing with type-level 3 and above involves some knotty semantic and complexity-theoretic issues beyond the scope of this paper, hence our restriction of ATR types to orders 2 and below.

²The work of [BNS00, Hof03] and of this paper sit on different sides of an important divide in highertype computability between notions of *computation over computable data* (e.g., [BNS00, Hof03, LM93]) and notions of *computation over continuous data* (e.g., this paper) [Lon04, Lon05].

³See [KC96, IKR01] for a more extensive justification that the BFFs provide a sensible type-2 analogue of the polynomial-time computable functions.

 $N \rightarrow N$) $\rightarrow N \rightarrow N$ with:

$$\begin{array}{cccc} \operatorname{prn} f \epsilon & \longrightarrow & f \epsilon \epsilon. \\ \operatorname{prn} f (\mathbf{a} \oplus y) & \longrightarrow & f (\mathbf{a} \oplus y) (\operatorname{prn} f y). \end{array} \right\}$$
(1.1)

(*Conventions:* \oplus denotes string concatenation and $\mathbf{a} \in \{\mathbf{0}, \mathbf{1}\}$.) So, prn is a version of Cobham's [Cob65] *primitive recursion on notation* (or alternatively, a string-variant of foldr). It is well-known that prn is **not** a BFF: starting with polynomial-time primitives, prn can be used to define any primitive recursive function. However as Cobham noted, if one modifies (1.1) by adding the side-condition:

$$(\exists p_f, \text{ a polynomial})(\forall x)[|\text{prn } f x| \leq p_f(|x|)],$$

this modified **prn** produces definitions of just polynomial-time computable functions from polynomial-time computable primitives. Bellantoni and Cook [BC92] showed how get rid of *explicit* use of such a side condition through what amounts to a typing discipline. However, their approach (which has been in large part adopted by the implicit computational complexity community, see Hofmann's survey [Hof00]), requires that prn be a "special form" and that the f in (1.1) must be ultimately given by a purely syntactic definition. We, on the other hand, want to be able to define prn within ATR (see Figure 13) and have the definition's meaning given by a conventional, higher-type denotational semantics. We thus use Bellantoni and Cook's [BC92] (and Leivant's [Lei95]) ideas in both syntactic and se*mantic* contexts. That is, we extract the growth-rate bounds implicit in the aforementioned systems, extend these bounds to higher types, and create a type system, programming language, and semantic models that work to enforce these bounds. As a consequence, we can define prn (with a particular typing) and be assured that, whether the f corresponds to a purely syntactic term or to the interpretation a free variable, prn will not go wrong by producing something of huge complexity. The language and its model thus implicitly incorporate side-conditions on growth via types.⁴ Handling constructs like prn as first class functions is important because programmers care more about such combinators than about most any standard BFF.

Outline. Our ATR formalism is based on Bellantoni and Cook [BC92] and Leivant's [Lei95] ideas on using "data ramification" to rein in computational complexity. §3 puts these ideas in a concrete form of BCL, a simple type-level 1 programming formalism, and sketches the proofs of three basic results on BCL: (i) that each BCL expression is *polynomial size-bounded*, (ii) that computing the value of a BCL expression is *polynomial time-bounded*, and (iii) each polynomial-time computable function is denoted by some BCL-expression. Most of this paper is devoted to showing the analogous results for ATR. §4 discusses how one might change BCL into a type-2 programming formalism, some of the problems one encounters, and our strategies for dealing with these problems. ATR, our type-2 system, is introduced in §5 along with its type system, typing rules, and basic syntactic properties. The goal of §§6–10 is to show (type-2) polynomial size-boundedness for ATR. This is complicated by the fact (described in §6) that the naïve semantics for ATR permits exponential blow-ups. §§7–9 show how to prune back the naïve semantics to obtain a setting in which we can

⁴Incorporating side-conditions in models is nothing new. A fixed-point combinator has the implicit sidecondition that its argument is continuous (or at least monotone) so that, by Tarski's fixed-point theorem [Win93], we know the result is meaningful. Models of languages with fixed-point combinators typically have continuity built-in so the side-condition is always implicit.

prove polynomial size-boundedness, which is shown in §10. The goal of §§11–15 is to show (type-2) polynomial time-boundedness for ATR. Our notion of the cost of evaluating ATR expressions is based on a particular abstract machine (described in §11.1) that implements an ATR-interpreter and the costs we assign to this machine's steps (described in §11.2). §12 and §13 set up a time-complexity semantics for ATR⁻ expressions (where ATR⁻ consists of ATR without its recursion construct) and establish that this time-complexity semantics is: (i) sound for the abstract machine's cost model (i.e., the semantics provides upper bounds on these costs), and (ii) *polynomial time-bounded*, that is that the time-complexity each ATR expression e has a second-order polynomial bound over the time-complexities of e's free variables. §16 shows that ATR can compute each type-2 basic feasible functional. §17 considers possible extensions of our work. We begin in §2 which sets out some basic background definitions with §§2.8–2.14 covering the more exotic topics.

Acknowledgments. Thanks to Susan Older and Bruce Kapron for repeatedly listening to the second author describe this work along its evolution. Thanks to Neil Jones and Luke Ong for inviting the second author to Oxford for a visit and for some extremely helpful comments on an early draft of this paper. Thanks to Syracuse University for hosting the first author during September 2005. Thanks also to the anonymous referees of both the POPL version of this paper [DR06] and the present paper for many extremely helpful comments. Finally many thanks to Peter O'Hearn, Josh Berdine, and the Queen Mary theory group for hosting the second author's visit in the Autumn of 2005 and for repeatedly raking his poor type-systems over the coals until something reasonably simple and civilized survived the ordeals. This work was partially supported by EPSRC grant GR/T25156/01 and NSF grant CCR-0098198.

2. NOTATION AND CONVENTIONS

2.1. Numbers and strings. We use two representations of the natural numbers: dyadic and unary. Each element of \mathbb{N} is identified with its dyadic representation over $\{0, 1\}$, i.e., $0 \equiv \epsilon, 1 \equiv 0, 2 \equiv 1, 3 \equiv 00$, etc. We freely pun between $x \in \mathbb{N}$ as a number and a 0-1-string. Each element of ω is identified with its unary representation over $\{0\}$, i.e., $0 \equiv \epsilon, 1 \equiv 0, 2 \equiv 00, 3 \equiv 000$, etc. The elements of \mathbb{N} are used as numeric/string values to be computed over. The elements of ω are used as tallies to represent lengths, run times, and generally anything that corresponds to a size measurement. *Notation:* For each natural number $k, \underline{k} = \mathbf{0}^k$. Also $x \oplus y =$ the concatenation of strings x and y.

2.2. Simple types. Below, **b** (with and without decorations) ranges over base types and *B* ranges over nonempty sets of base types. The *simple types* over *B* are given by: $T ::= B \mid (T \to T)$. As usual, \to is right associative and unnecessary parentheses are typically dropped in type expressions, e.g., $(\sigma_1 \to (\sigma_2 \to \sigma_3)) = \sigma_1 \to \sigma_2 \to \sigma_3$. A type $\sigma_1 \to \cdots \to \sigma_k \to \mathbf{b}$ is often written as $(\sigma_1, \ldots, \sigma_k) \to \mathbf{b}$ or, when $\sigma = \sigma_1 = \cdots = \sigma_k$, as $(\sigma^k) \to \mathbf{b}$. The *simple product types* over *B* are given by: $T ::= B \mid T \to T \mid () \mid T \times T$, where () is the type of the empty product. As usual, $\sigma^1 = \sigma$, $\sigma^{k+1} = \sigma^k \times \sigma$, \times is left associative, and unnecessary parentheses are typically dropped in type expressions. The *level* of a simple (product) type is given by: $level(\mathbf{b}) = level(()) = 0$; $level(\sigma \times \tau) = \max(level(\sigma), level(\tau))$; and $level(\sigma \to \tau) = \max(1 + level(\sigma), level(\tau))$. In this paper types are always interpreted over cartesian closed categories; hence, the two types $(\sigma_1, \ldots, \sigma_k) \to \tau$ and $\sigma_1 \times \cdots \times \sigma_k \to \tau$ may be identified. By convention, we identify $() \to \tau$ with τ and $\lambda() \cdot e$ with e.

2.3. Subtyping. Suppose \leq : is a reflexive partial order on *B*. Then \leq : can be extended to a reflexive partial order over the simple types over *B* by closing under:

$$\sigma_1 \leq : \sigma_0 \& \tau_0 \leq : \tau_1 \iff \sigma_0 \to \tau_0 \leq : \sigma_1 \to \tau_1.$$

$$(2.1)$$

We read " $\sigma \leq \tau$ " as σ is a *subtype* of τ ; and write $\tau \geq \sigma$ for $\sigma \leq \tau$ and $\sigma \leq \tau$ for $[\sigma \leq \tau]$ and $\sigma \neq \tau$.

2.4. Type contexts and judgments. A type context Γ is a finite (possibly empty) mapping of variables to types; these are usually written as a list: $v_1: \sigma_1, \ldots, v_k: \sigma_k$. Γ, Γ' denotes the union of two type contexts with disjoint preimages. $\Gamma \cup \Gamma'$ denotes the union of two consistent type contexts, that is, $\Gamma(x)$ and $\Gamma'(x)$ are equal whenever both are defined. The type judgment $\Gamma \vdash_{\mathcal{F}} e: \sigma$ asserts that the assignment of type σ to expression e follows from the type assignments of Γ under the typing rules for formalism \mathcal{F} . We typically omit the subscript in $\vdash_{\mathcal{F}}$ when \mathcal{F} is clear from context.

2.5. Semantic conventions. For a particular semantics S for a formalism \mathcal{F} , $S[\![\,\cdot\,]\!]$ is the semantic map that takes an \mathcal{F} -syntactic object to its S-meaning. $S[\![\tau]\!]$ is the collection of things named by a type τ under S. For a type context $\Gamma = x_1: \tau_1, \ldots, x_n: \tau_n$, $S[\![\Gamma]\!]$ is the set of all finite maps $\{x_1 \mapsto a_1, \ldots, x_n \mapsto a_n\}$, where $a_1 \in S[\![\tau_1]\!], \ldots, a_n \in S[\![\tau_n]\!]$; i.e., environments. Convention: ρ (with and without decorations) ranges over environments and $\{\}$ = the empty environment. $S[\![\Gamma \vdash e:\tau]\!]$ is the map from $S[\![\Gamma]\!]$ to $S[\![\tau]\!]$ such that $S[\![\Gamma \vdash e:\tau]\!] \rho$ denotes the element of $S[\![\tau]\!]$ that is the S-meaning of expression e when e's free-variables have the meanings given by ρ . Conventions: $S[\![e]\!]$ is typically written in place of $S[\![\Gamma \vdash e:\tau]\!]$ since the type judgment is usually understood from context. When e is closed, $S[\![e]\!]$ is sometimes written in place of $S[\![e]\!]$ {}. Also, $e_0 =_S e_1$ means $S[\![e_0]\!] = S[\![e_1]\!]$.

2.6. Syntactic conventions. Substitutions (e.g., e[x := e']) are always assumed to be capture avoiding. Terms of the form $(x e_1 \dots e_k)$ are sometimes written as $x(e_1, \dots, e_k)$.

2.7. Call-by-value PCF. The syntax of our version of PCF is given in Figure 1, where the syntactic categories are: constants (K), raw-expressions (E), variables (X), and typeexpressions (T) and where $\mathbf{a} \in \{\mathbf{0}, \mathbf{1}\}$. Figure 2 states PCF's typing rules, where **op** stands for any of $\mathbf{c_0}$, $\mathbf{c_1}$, \mathbf{d} , $\mathbf{t_0}$, and $\mathbf{t_1}$ and where $E ::= e, e_0, e_1, e_2, K ::= k, T ::= \sigma, \tau$, and X ::= x. For emphasis we may write $\lambda x : \sigma \cdot e$ instead of $\lambda x \cdot e$, but the type of x can always be inferred from any type judgement in which $\lambda x \cdot e$ occurs. The intended interpretation of N is $\mathbb{N} \ (\cong \{\mathbf{0}, \mathbf{1}\}^*)$. The reduction rules are essentially the standard ones for call-by-value PCF (see [Plo75, Pie02]). In particular, the reduction rules for $\mathbf{c_0}$, $\mathbf{c_1}$, \mathbf{d} , $\mathbf{t_0}$, $\mathbf{t_1}$, down, ifthen-else, and fix are given in Figure 3. Note that in if-then-else tests, ϵ corresponds to *false* and elements of ($\mathbb{N} - \{\epsilon\}$) correspond to *true*. In tests, we use $x \neq \epsilon$ for syntactic sugar for x and use $|e_0| \leq |e_1|$ as syntactic sugar for (down $\mathbf{c_0}(e_0) \mathbf{c_0}(e_1)$).⁵ An operational semantics for PCF is provided by the CEK-machine given in §11.1. We take \mathcal{V} (for *value*) to be a

⁵We will see in §4 and §5 why down (as opposed to $|\cdot| \leq |\cdot|$) is a primitive.

$$E ::= K | (\mathbf{c_a} E) | (\mathbf{d} E) | (\mathbf{t_a} E) | (\operatorname{down} E E)$$
$$| X | (E E) | (\lambda X \cdot E) | (\operatorname{if} E \operatorname{then} E \operatorname{else} E) | (\operatorname{fix} E)$$
$$K ::= \{\mathbf{0}, \mathbf{1}\}^* \qquad T ::= \operatorname{the simple types over N}$$

Figure	1:	PCF	SI	vntax

Const-I: $\overline{\Gamma \vdash k: N}$	op- <i>I</i> : $\frac{\Gamma \vdash e: N}{\Gamma \vdash (op \ e): N}$	down- <i>I</i> : $\frac{\Gamma_0 \vdash e_0: N \qquad \Gamma_1 \vdash e_1: N}{\Gamma_0 \cup \Gamma_1 \vdash (down \ e_0 \ e_1): N}$
<i>Id-I:</i> $\overline{\Gamma, x: \sigma \vdash x: \sigma}$	$\rightarrow -I: \frac{\Gamma, x: \sigma \vdash e: \tau}{\Gamma \vdash (\lambda x \cdot e): \sigma \rightarrow \tau}$	$\rightarrow -E: \frac{\Gamma_0 \vdash e_0: \sigma \rightarrow \tau \Gamma_1 \vdash e_1: \sigma}{\Gamma_0 \cup \Gamma_1 \vdash (e_0 \ e_0): \tau}$
<i>If-I:</i> $\Gamma_0 \vdash e_0: N$ $\Gamma_0 \cup \Gamma_1 \cup \Gamma$	$\Gamma_1 \vdash e_1 : N \qquad \Gamma_2 \vdash e_2 : N$ $\Gamma_2 \vdash (if \ e_0 \ then \ e_1 \ else \ e_2) : N$	$ \begin{array}{ll} {\rm fix-}I: & \frac{\Gamma \vdash (\lambda x \cdot e) : \sigma \to \sigma}{\Gamma \vdash ({\rm fix} \; (\lambda x \cdot e)) : \sigma} \end{array} \end{array} $

Figure 2: The PCF typing rules

$(c_{\mathbf{a}} v) \longrightarrow \mathbf{a} \oplus v.$	$(d (a \oplus u))$	$v)) \longrightarrow v.$	$(d \epsilon)$ —	$\rightarrow \epsilon$.
$ (\mathbf{t_a} \ v) \longrightarrow \begin{cases} 0, & \text{if } v \text{ begins with} \\ \epsilon, & \text{otherwise.} \end{cases} $	th $\mathbf{a};$	$(down\; v_0\; v_1)$ –	$\longrightarrow \begin{cases} v_0, \\ \epsilon, \end{cases}$	if $ v_0 \le v_1 $; otherwise.
$(\text{if } v_0 \text{ then } v_1 \text{ else } v_2) \longrightarrow \begin{cases} v_1, \\ v_2, \end{cases}$	$ if \ v_0 \neq \epsilon; \\ if \ v_0 = \epsilon. $	fix $(\lambda x \cdot e)$	$\longrightarrow e[x]$	$:= (fix\;(\lambda x \boldsymbol{.} e)].$

Figure 3: The PCF reduction rules for c_a , d, t_a , down, if-then-else, and fix

conventional denotational semantics for PCF [Win93]. Standard arguments show that our operational semantics corresponds to \mathcal{V} .

2.8. Total continuous functionals. Let σ and τ be simple product types over base type N. Inductively define: $\mathbf{TC}_{()} = \star$; $\mathbf{TC}_{\mathsf{N}} = \mathbb{N}$; $\mathbf{TC}_{\sigma \times \tau} = \mathbf{TC}_{\sigma} \times \mathbf{TC}_{\tau}$; $\mathbf{TC}_{\sigma \to \tau} =$ the *Kleene/Kreisel total continuous functions* from \mathbf{TC}_{σ} to \mathbf{TC}_{τ} ; the \mathbf{TC}_{σ} 's together form a cartesian closed category \mathbf{TC} .⁶ This paper is concerned with only the type-level 0, 1, and 2 portions of \mathbf{TC} from which we construct models of our programming formalisms.

2.9. Total monotone continuous functionals. Let σ and τ be simple product types over base type T (for *tally*). Inductively define the \mathbf{MC}_{σ} sets and partial orders \leq_{σ} by: $\mathbf{MC}_{\mathsf{T}} = \omega$ and \leq_{T} = the usual ordering on ω ; $\mathbf{MC}_{()} = \star$ and $\star \leq_{()} \star$; $\mathbf{MC}_{\sigma \times \tau} = \mathbf{MC}_{\sigma} \times \mathbf{MC}_{\tau}$ and $(a, b) \leq_{\sigma \times \tau} (a', b') \iff a \leq_{\sigma} a'$ and $b \leq_{\tau} b'$; and $\mathbf{MC}_{\sigma \to \tau} =$ the Kleene/Kreisel total continuous functions from \mathbf{MC}_{σ} to \mathbf{MC}_{τ} that are monotone (w.r.t. \leq_{σ} and \leq_{τ}), and $\leq_{\sigma \to \tau}$ is the point-wise ordering on $\mathbf{MC}_{\sigma \to \tau}$. (E.g., $\mathbf{MC}_{\mathsf{T} \to \mathsf{T}} = \{f: \omega \to \omega \mid f(0) \leq f(1) \leq f(2) \leq \cdots \}$.) The \mathbf{MC}_{σ} 's turn out to form a cartesian closed category \mathbf{MC} . As with \mathbf{TC} , our

⁶For background on the Kleene/Kreisel total continuous functions and **TC**, see the historical survey of Longley [Lon05] and the technical surveys of Normann [Nor99] and Schwichtenberg [Sch96].

$$P ::= K \mid (\lor P P) \mid (+P P) \mid (*P P) \mid V \mid (P P) \mid (\lambda V \cdot P)$$

 $K ::= \underline{0} \mid \underline{1} \mid \underline{2} \mid \ldots$ T ::= the level 0, 1, and 2 simple types over T

Figure 4: The syntax for second-order polynomials and their standard types

Const-I:
$$\underline{\Sigma \vdash k: \mathsf{T}}$$
 \odot -I: $\underline{\Sigma \vdash p_0: \mathsf{T} \quad \Sigma \vdash p_1: \mathsf{T}}_{\underline{\Sigma} \vdash (\odot \ p_0 \ p_1): \mathsf{T}}$ $(\odot = *, +, \lor)$

Figure 5: The additional typing rules for the second-order polynomials

concern is with only the type-level 0, 1, and 2 portions of **MC** from which we construct our models of size and time bounds. *Convention:* We typically omit the subscript in \leq_{σ} when the σ is clear from context.

2.10. Lengths. For $v \in \mathbb{N}$, let $|v| = \underline{k}$, where k is the length of the dyadic representation of v (e.g., $|\mathbf{110}| = \underline{3} = \mathbf{000}$). For $f \in \mathbf{TC}_{(\mathbb{N}^k) \to \mathbb{N}}$, define $|f| \in \mathbf{MC}_{(\mathbb{T}^k) \to \mathbb{T}}$ by:

$$|f|(\vec{\ell}) = \max\{|f(\vec{v})| \mid |v_1| \le \ell_1, \dots, |v_k| \le \ell_k\}.$$
(2.2)

(This is Kapron and Cook's [KC96] definition.) For each σ , a simple type over N, let $|\sigma| = \sigma[N := T]$ (e.g., $|N \to N| = T \to T$). So by the above, $|v| \in \mathbf{MC}_{|\sigma|}$ when $level(\sigma) \leq 1$ and $v \in \mathbf{TC}_{\sigma}$. Here is a type-level 2 notion of length that suffices for this paper. For $\gamma = (\sigma_1, \ldots, \sigma_k) \to N$ of level-2, $F \in \mathbf{TC}_{\gamma}$, and $\ell_1 \in \mathbf{MC}_{|\sigma_1|}, \ldots, \ell_k \in \mathbf{MC}_{|\sigma_k|}$, define

$$|F|(\vec{\ell}) = \max\{ |F(\vec{v})| \mid |v_1| \le_{|\sigma_1|} \ell_1, \dots, |v_k| \le_{|\sigma_k|} \ell_k \}.$$
(2.3)

|F| as defined above turns out to be an element of $\mathbf{MC}_{|\gamma|}$.⁷

2.11. Maximums and polynomials. Let $v_1 \vee v_2 = \max(\{v_1, v_2\})$ and let $\bigvee_{i=1}^k v_i = \max(\{v_1, \dots, v_k\})$ for $v_1, \dots, v_k \in \omega$. By convention, $\max(\emptyset) = 0$. We allow \vee as another arithmetic operation in polynomials; \vee binds closer than either multiplication or addition. Coefficients in polynomials will always be nonnegative; hence polynomials denote monotone nondecreasing functions, i.e., type-level 1 elements of **MC**.

2.12. Second-order polynomials. We define the second-order polynomials [KC96] as a type-level 2 fragment of the simply typed λ -calculus over base type T with arithmetic operations \vee , +, and *. Figure 4 gives the syntax, where the syntactic categories are: constants (K), raw expressions (P), and type expressions (T). We often write \vee -, +-, and *-expressions in infix form. The typing rules are Id-I, \rightarrow -I, and \rightarrow -E from Figure 2 together with the rules in Figure 5. Moreover, the only variables allowed are those of of type levels 0 and 1. Our semantics \mathcal{L} (for *length*) for second-order polynomials is: $\mathcal{L}[\![\sigma]\!] = \mathbf{MC}_{\sigma}$ for each σ , a simple type over T, and $\mathcal{L}[\![\Sigma \vdash p:\sigma]\!] =$ the standard definition. The *depth* a second-order polynomial q is the maximal depth of nesting of applications in q's β -normal

⁷When γ is type-level 2 and $\ell \in \mathbf{MC}_{|\gamma|}$, generally $\{F \in \mathbf{TC}_{\gamma} \mid |F| \leq_{|\gamma|} \ell\}$ fails to be compact in the appropriate topology. Consequently, the type-3 analogue of (2.3) fails to yield lengths that are total. There are alternative notions of type-2 length that avoid this problem; [IKR02] investigates two of these.

form, e.g., $g_0((g_0(2 * y * g_1(y^2)) \vee 6)^3)$ has depth 3. There is a special case for variables of higher type: type-level ℓ variables are assigned depth ℓ .⁸ For second-order polynomials, depth plays something like the role degree does for ordinary polynomials.

2.13. Time complexity. The CEK machine (§11.1) provides an operational semantics for PCF as well as for the formalisms BCL (§3) and ATR (§5). Since this paper concerns the evaluation of expressions and the associated costs, we use the CEK machine as our standard model of computation and use the CEK cost model (§11.2) as our standard notion of time complexity. As discussed in §11.2, Schönhage's *storage modification machine* [Sch80] is roughly the standard complexity-theoretic model of computation and cost underlying our CEK model. Storage modification machines and Turing machines are polynomially-related models of computation [Sch80]. Our CEK machine handles oracles (type-1 functions over \mathbb{N}) as the values of particular variables in the initial environment for an evaluation. As with Kapron and Cook's answer-length cost model for oracle Turing machines [KC96], part of the CEK-cost of querying an oracle includes the length of the answer.

2.14. **Basic feasibility.** Suppose $\tau = (\sigma_1, \ldots, \sigma_k) \to \mathsf{N}$ is a simple type over N of level 1 or 2 and that $f \in \mathcal{V}[\![\tau]\!]$. $(\mathcal{V}[\![\cdot]\!]$ was introduced in §2.7.) We say that f is a *basic feasible functional* (or BFF) when there is a closed, type- τ PCF-expression e_f and a second-order polynomial function q_f such that (i) $\mathcal{V}[\![e_f]\!] = f$ and (ii) for all $v_i \in \mathcal{V}[\![\sigma_1]\!], \ldots, v_k \in \mathcal{V}[\![\sigma_k]\!]$, CEK-time $(e_f, v_1, \ldots, v_k) \leq q_f(|v_1|, \ldots, |v_k|)$, where CEK-time is introduced in Definition 48 of §11.2. For level-1 τ , this gives us the usual notion of type-1 polynomial-time computability. The original definitions and characterizations of the type-2 BFFs [Meh74, CU93, CK90] were all in terms of programming formalisms. The definition here is based on Kapron and Cook's machine-based characterization of the type-2 BFFs [KC96].

3. The BCL formalism

The programming formalisms of this paper are built on work of Bellantoni and Cook [BC92] and Leivant [Lei95]. Bellantoni and Cook's paper takes a programming formalism for the primitive recursive functions, imposes certain intensionally-motivated constraints, and obtains a formalism for the polynomial-time computable functions. To explain these constraints and how they rein in computational strength, we sketch both BCL, a simple type-1 programming formalism based on Bellantoni and Cook's and Leivant's ideas, and BCL's properties.⁹ This sketch provides an initial framework for this paper's formalisms.

BCL has the same syntax as PCF (§2.7) with three changes: (i) fix is replaced with prn (for *primitive recursion on notation* [Cob65]) that has the reduction rule given by (1.1), (ii) the only variables allowed are those of base type, and (iii) the type system is altered as described below. If we were to stay with the simple types over N and the PCF-typing rules (Figure 2 and with prn: $(N \rightarrow N \rightarrow N) \rightarrow N \rightarrow N)$, the resulting formalism would compute exactly the primitive recursive functions. Instead we modify the types and typing as follows. N is replaced with two base types, N_{norm} (normal values) and N_{safe} (safe values), subtype ordered N_{norm} \leq : N_{safe}. The BCL types are just the type-level 0 and 1 simple types

⁸Since, for example, for $f: \mathsf{T} \to \mathsf{T}$, $f \equiv_{\eta} \lambda x \cdot f(x)$ and $depth(\lambda x \cdot f(x)) = 1$.

⁹BCL is much closer to Leivant's formalism [Lei95], which uses a ramified type system, than Bellantoni and Cook's, which does not use a conventional type system.

 $E ::= \dots | (prn E)$ T ::= the level 0 and 1 types over N_{norm} and N_{safe}

Figure 6: BCL syntax		
Zero-I: $\frac{\Gamma \vdash e: \sigma}{\Gamma \vdash e: \tau}$ ($\sigma \leq : \tau$)		
$d\text{-}I': \frac{\Gamma \vdash e: N_{\mathrm{norm}}}{\Gamma \vdash (d \ e): N_{\mathrm{norm}}} \qquad prn\text{-}I: \frac{\Gamma \vdash e: N_{\mathrm{norm}} \to N_{\mathrm{safe}} \to N_{\mathrm{safe}}}{\Gamma \vdash (prn \ e): N_{\mathrm{norm}} \to N_{\mathrm{safe}}}$		
Figure 7: Additional BCL typing rules		
$\begin{array}{ll} \operatorname{cat:} N_{\operatorname{norm}} \to N_{\operatorname{safe}} \to N_{\operatorname{safe}} = & //\operatorname{cat} w \ x = w \oplus x. \ So, \ \operatorname{cat} w \ x = w + x . \\ \lambda w, x \text{. let } f: N_{\operatorname{norm}} \to N_{\operatorname{safe}} \to N_{\operatorname{safe}} = \\ & \lambda y, z \text{. if } t_{0}(y) \ \operatorname{then} c_{0}(z) \ \operatorname{else} \ \operatorname{if} t_{1}(y) \ \operatorname{then} c_{1}(z) \ \operatorname{else} x \\ & \operatorname{in} \ \operatorname{prn} \ f \ w \end{array}$ $\begin{array}{l} dup: N_{\operatorname{norm}} \to N_{\operatorname{norm}} \to N_{\operatorname{safe}} = & //\operatorname{dup} \ w \ x = \underbrace{x \oplus \cdots \oplus x}_{x}. \ So, \ \operatorname{dup} \ w \ x = w \cdot x . \\ \lambda w, x \text{. let} \ g: N_{\operatorname{norm}} \to N_{\operatorname{safe}} \to N_{\operatorname{safe}} = \lambda y, z \text{. if } y \neq \epsilon \ \operatorname{then} \ (\operatorname{cat} x \ z) \ \operatorname{else} \epsilon \\ & \operatorname{in} \ \operatorname{prn} \ g \ w \end{array}$		

Figure 8: Two sample BCL programs

over N_{norm} and N_{safe} . Both base types have intended interpretation N. The point of the two base types is to separate the roles of N-values: a N_{norm} -value can be used to drive a recursion, but cannot be the result of a recursion, whereas a N_{safe} -value can be the result of a recursion, but cannot be used to drive a recursion. These intentions are enforced by the BCL typing rules, consisting of: ID-I, $\rightarrow -I$, and $\rightarrow -E$ from Figure 2; Const-I, c_0-I , c_1-I , d-I, t_0-I , t_1-I , down-I, and If-I also from Figure 2 where each N is changed to N_{safe} ; and the rules in Figure 7. (Zero-I and d-I' are needed to make the prn reduction rules type-correct.) Figure 8 contains two sample BCL programs. For the sake of readability, we use the let construct as syntactic sugar.¹⁰

Propositions 1, 2, and 3 state the key computational limitations and capabilities of BCL. In the following \vec{x} : N_{norm} abbreviates x_1 : N_{norm}, \ldots, x_m : N_{norm} and \vec{y} : N_{safe} abbreviates y_1 : N_{safe}, \ldots, y_n : N_{safe} . Recall from §2.13 that our standard notion of time complexity is the time cost model of the CEK-machine (Definition 48(a)).

Proposition 1 (BCL polynomial size-boundedness). Suppose $\vec{x}: N_{norm}, \vec{y}: N_{safe} \vdash e: b$.

(a) If $\mathbf{b} = \mathsf{N}_{\text{norm}}$, then for all values of $\vec{x}, \vec{y}, |e| \leq \bigvee_{i=1}^{m} |x_i|$.

(b) If $\mathbf{b} = \mathsf{N}_{\text{safe}}$, then there is a polynomial p over over $|x_1|, \ldots, |x_m|$ such that, for all values of $\vec{x}, \vec{y}, |e| \le p + \bigvee_{j=1}^n |y_j|$.

Proposition 1's proof is an induction on e's syntactic structure, where the prn-case is the crux of the argument. Here is a sketch of a mild simplification of that case. (This sketch is the model for several key subsequent arguments.) Suppose $e = \operatorname{prn} e' x$, where $x_0: N_{\operatorname{norm}}, \vec{x}: N_{\operatorname{norm}}, y_0: N_{\operatorname{safe}}, \vec{y}: N_{\operatorname{safe}} \vdash (e' x_0 y_0): N_{\operatorname{safe}}$ and $x \in \{x_1, \ldots, x_m\}$. Also suppose

¹⁰Where (let x = e' in e) $\stackrel{\text{def}}{\equiv} e[x := e']$. This permits naming defined functions.

that, for all values of $x_0, \ldots, x_m, y_0, \ldots, y_n$, $|e' x_0 y_0| \leq p'(|x_0|) + \bigvee_{j=0}^n |y_j|$ where p' is a polynomial over $|x_0|$ (explicitly) and $|x_1|, \ldots, |x_m|$ (implicitly). Fix the values of x_1, \ldots, y_n , where in particular x has the value $\mathbf{a}_1 \ldots \mathbf{a}_k$ for $\mathbf{a}_1, \ldots, \mathbf{a}_k \in \{\mathbf{0}, \mathbf{1}\}$. We determine bounds for $|\mathsf{prn} \ e' \ \epsilon|$, $|\mathsf{prn} \ e' \ \mathbf{a}_k|$, $|\mathsf{prn} \ e' \ \mathbf{a}_{k-1}\mathbf{a}_k|, \ldots, |\mathsf{prn} \ e' \ \mathbf{a}_1 \ldots \mathbf{a}_k|$ in turn. First, $|\mathsf{prn} \ e' \ \epsilon| = |e' \ \epsilon \ \epsilon| \leq p'(\underline{0}) + \bigvee_{j=1}^n |y_j|$. Next,

$$\begin{aligned} |\mathsf{prn} \ e' \ \mathbf{a}_k| &= |e' \ \mathbf{a}_k \ (\mathsf{prn} \ e' \ \epsilon)| \leq p'(\underline{1}) + |\mathsf{prn} \ e' \ \epsilon| \lor \bigvee_{j=1}^n |y_j| \leq \\ p'(\underline{1}) + (p'(\underline{0}) + \bigvee_{j=1}^n |y_j|) \lor \bigvee_{j=1}^n |y_j| \leq p'(\underline{0}) + p'(\underline{1}) + \bigvee_{j=1}^n |y_j|. \end{aligned}$$

Continuing, we end up with $|\text{prn } e' x| \le p'(\underline{0}) + p'(\underline{1}) + \dots + p'(\underline{k}) + \bigvee_{j=1}^{n} |y_j| \le (|x| + \underline{1}) * p'(|x|) + \bigvee_{j=1}^{n} |y_j|$. So, p = (|x| + 1) * p'(|x|) suffices for this case.

Proposition 2 (BCL polynomial time-boundedness). Given \vec{x} : N_{norm} , \vec{y} : $N_{safe} \vdash e$: $(\mathbf{b}_1, \ldots, \mathbf{b}_{\ell}) \rightarrow \mathbf{b}$, there is a polynomial q over over $|w_1|, \ldots, |w_{\ell}|, |x_1|, \ldots, |x_m|, |y_1|, \ldots, |y_n|$ such that, for all values of w_1, \ldots, y_n , q bounds the CEK-cost of evaluating $(e \ w_1 \ \ldots \ w_{\ell})$.

Proposition 2's proof rests on three observations: (i) evaluating (prn e e') takes |e'|many (top-level) recursions, (ii) by the first observation and the details of CEK costs, the time-cost of a CEK evaluation of a BCL expression can be bounded by a polynomial over the lengths of base type values involved, and (iii) Proposition 1 provides polynomial bounds on all these lengths. Proposition 2 thus follows through a straightforward induction on the syntactic structure of e. Proposition 3's proof is mostly an exercise in programming.

Proposition 3 (BCL polynomial-time completeness). For each polynomial-time computable $f \in ((\mathbb{N}^{\ell}) \to \mathbb{N})$, there is an $\vdash e_f: (\mathbb{N}_{norm}^{\ell}) \to \mathbb{N}_{safe}$ such that $\mathcal{V}\llbracket e_f \rrbracket = f$.

BCL is \leq :-predicative in the sense that no information about a N_{safe}-value can ever make its way into a N_{norm}-value. For example:

Proposition 4. Suppose $\vdash e: (\mathsf{N}_{norm}, \mathsf{N}_{safe}) \to \mathsf{N}_{norm}$. Then $e \equiv_{\alpha\beta} \lambda w, x \cdot e'$ with $e' = \epsilon$ or else $e' = (\mathsf{d}^{(k)} w)$ for some $k \ge 0$, where $(\mathsf{d}^{(0)} w) = w$ and $(\mathsf{d}^{(k+1)} w) = (\mathsf{d} (\mathsf{d}^{(k)} w))$.

BCL's \leq :-predicativity plays a key role in proving the polynomial size-bounds of Proposition 1, but plays no direct (helpful) role in the other proofs.

4. BUILDING A BETTER BCL

Our definition of ATR in the next section can be thought of as building an extension of BCL that: (i) computes the type-2 BFFs, (ii) replaces prn with something closer to fix, and (iii) admits reasonably direct complexity theoretic analyses. This section motivates some of the differences between BCL and ATR.

Types and depth. We want to extend BCL's type system to allow definitions of functions as such $F_0 = \lambda f \in \mathbb{N} \to \mathbb{N}, x \in \mathbb{N}$. f(f(x)), a basic feasible functional. A key question then is how to assign types to functional parameters such as f above. Under $f: \mathbb{N}_{norm} \to \mathbb{N}_{safe}$, F_0 fails to have a well-typed definition. Under any of $f: \mathbb{N}_{norm} \to \mathbb{N}_{norm}, f: \mathbb{N}_{safe} \to \mathbb{N}_{safe}$, and $f: \mathbb{N}_{safe} \to \mathbb{N}_{norm}, F_0$ has a well-typed definition, but then so does $F_1 = \lambda f \in \mathbb{N} \to \mathbb{N}, x \in \mathbb{N}$. $f^{(|x|)}(x)$ which is *not* basic feasible. Thus some nontrivial modification of the BCL types seems necessary for any extension to type-level 2.

$$\begin{array}{lll} \operatorname{\mathsf{down}} I' \colon & \frac{\Gamma_0 \vdash e_0 \colon \mathsf{N}_{\mathrm{safe}} & \Gamma_1 \vdash e_1 \colon \mathsf{N}_{\mathrm{norm}}}{\Gamma_0 \cup \Gamma_1 \vdash (\operatorname{\mathsf{down}} e_0 \ e_1) \colon \mathsf{N}_{\mathrm{norm}}} \\ \\ \mathit{If} \text{-} I' \colon & \frac{\Gamma_0 \vdash e_0 \colon \mathsf{N}_{\mathrm{safe}} & \Gamma_1 \vdash e_1 \colon \mathsf{N}_{\mathrm{norm}} & \Gamma_2 \vdash e_2 \colon \mathsf{N}_{\mathrm{norm}}}{\Gamma_0 \cup \Gamma_1 \cup \Gamma_2 \vdash (\mathrm{if} \ e_0 \ \mathrm{then} \ e_1 \ \mathrm{else} \ e_2) \colon \mathsf{N}_{\mathrm{norm}}} \end{array}$$

Figure 9: Additional rules for BCL'

We sketch a naïve extension that uses of an informal notion of the depth of an expression (based on second-order polynomial depth, see §2.12). Let the *naïve depth* of an expression (in normal form) be the depth of nesting of applications of type-level 1 variables. For example, given $f: \mathbb{N} \to \mathbb{N} \to \mathbb{N}$, then $f(\mathbf{c_0}(f(\mathbf{c_0}(x), y)), \mathbf{c_1}(\mathbf{d}(y)))$ has naïve depth 2. We can regard the values of x and y as (depth-0) inputs and the values of $\mathbf{c_0}(x)$ and $\mathbf{d}(y)$ as the results of polynomial-time computations over those inputs. Taking the type-level 1 variables as representing oracles, the value of $f(\mathbf{c_0}(x), y)$ can then be regarded as a depth-1 input (that is an input that is in response to a depth-0 query); hence, $\mathbf{c_0}(f(\mathbf{c_0}(x), y))$ is the result of a polynomial-time computation over a depth-1 input. Similarly, the value of $f(\mathbf{c_0}(f(\mathbf{c_0}(x), y)), \mathbf{c_1}(\mathbf{d}(y)))$ can be regarded as a depth-2 input. Thus, our naïve extension amounts to having, for each $d \in \omega$, depth-d versions of both \mathbb{N}_{norm} and \mathbb{N}_{safe} and treating all arrow types as "depth polymorphic" so, for instance, the type of f as above indicates that f takes depth-d safe values to depth-(d + 1) normal values, for each $d \in \omega$. This permits a well-typed definition for F_0 , but not for F_1 .

The naïvete of the above is shown by another example. Let

$$F_2 = \lambda f \in \mathbb{N} \to \mathbb{N}, \ y \in \mathbb{N} \, \left[g^{(|y|)}(y), \ \text{where} \ g = \lambda w \in \mathbb{N} \, \left(f(w) \ \text{mod} \ (y+1) \right) \right]. \tag{4.1}$$

 F_2 is basic feasible, $|F_2(f, y)| \leq |y|$, but it is reasonable to think of $F_2(f, y)$ having unbounded naïve depth.

Our solution to this problem is to use a more relaxed version of \leq :-predictivity than that of BCL. To explain this let us consider BCL', which is the result of adding rules of Figure 9 to BCL. (The rewrite rule for down is given in Figure 3.) These typing rules allow information about N_{safe} values to flow into N_{norm} values, but only in very controlled ways. In down-I', the controlling condition is that the length of this N_{safe} information is bounded by the length of some prior N_{norm} value. In *If-I'*, essentially only one bit of information about a N_{safe} value is allowed to influence the N_{norm} value of the expression. Because of these controlling conditions, the proofs of Propositions 1, 2, and 3 go through for BCL' with only minor changes, but in place of Proposition 4 we have:

Proposition 5. $\{\mathcal{V}\llbracket e \rrbracket : \vdash_{\mathsf{BCL}'} e : \mathsf{N}_{norm} \to \mathsf{N}_{safe} \to \mathsf{N}_{norm}\} = the set of polynomial-time computable <math>f \in \mathbb{N} \to \mathbb{N} \to \mathbb{N}$ such that $|f(x, y)| \leq |x|$ for all x and y.

Each BCL' type γ has a quantitative meaning in the sense that every element of $\{\mathcal{V}[\![e]\!] \mid \Gamma \vdash_{\mathsf{BCL}'} e: \gamma\}$ has a polynomial size-bound of a particular form. ATR has rules analogous to *If-I'* and down-*I'* and, consequently, functions such as F_2 have well-typed definitions. Moreover, each ATR type γ has a quantitative meaning in the sense that $\{\mathcal{V}[\![e]\!] \mid \vdash_{\mathsf{ATR}} e: \gamma\}$ = the set of all ATR-computable functions having second-order polynomial size-bounds of a form dictated by γ . In particular, for each γ , a $d_{\gamma} \in \omega$ can be read off such that all the bounding polynomials for type- γ objects can be of depth $\leq d_{\gamma}$. This is the (non-naïve)

connection of ATR's type-system to the notion of depth. The above glosses over the issue of the "depth polymorphic" higher types which are discussed in §5.

Truncated fixed points. For PCF, fix is thought as expressing general recursion. It would be ever so convenient if one could replace fix with some higher-type polynomial-time construct and obtain "the" feasible version of PCF in which all (and only) the polynomial-time recursion schemes are expressible. However, because of some basic limitations of subrecursive programming formalisms [Mar72, Roy87], it is unlikely that there is *any* finite collection of constructs through which one can express all and only such recursion schemes.

Our goals are thus more modest. We make use of the programming construct crec, for *clocked recursion*. The crec construct is a descendant of Cobham's [Cob65] bounded recursion on notation and *not* a true fixed-point constructor. The reduction rule for crec is:

$$\operatorname{crec} a (\lambda_r f \cdot e) \longrightarrow \lambda \vec{x} \cdot (\operatorname{if} |a| \le |x_1| \operatorname{then} (e' \vec{x}) \operatorname{else} \epsilon)$$

$$\operatorname{with} e' = e[f := (\operatorname{crec} (\mathbf{0} \oplus a) (\lambda_r f \cdot e))],$$

$$(4.2)$$

where a is a constant and $\vec{x} = x_1, \ldots, x_k$ is a sequence of variables. Roughly, |a| acts as the tally of the number of recursions thus far and $\mathbf{0} \oplus a$ is the result of a tick of the clock. The value of x_1 is the program's estimate of the total number of recursions it needs to do its job. Typing constraints will make sure that each **crec**-recursion terminates after polynomiallymany steps. Without these constraints, **crec** is essentially equivalent to fix. Clocking the fixed point process is a strong restriction. However, results on clocked programming systems ([RC94, Chapter 4]) suggest that clocking, whether explicit or implicit, is needed to produce programs for which one can determine explicit run-time bounds.

Along with clocking, we impose two other restrictions on recursions.

One use. In any expression of the form (crec a ($\lambda_r f \cdot e$)), we require that f has at most one use in e. Operationally this means that, in any possible evaluation of e, at most one application of f takes place. One consequence of this restriction is that no free occurrence of f is allowed within any inner crec expression. (Even if f occurs but once in an inner crec, the presumption is that f may be used many times.) Affine typing constraints enforce this one-use restriction. Note that prn is a one-use form of recursion.

The motivation for the one-use restriction stems from the recurrence equations that come out of time-complexity analyses of recursions. Under the one-use restriction, bounds on the cost of m steps of a **crec** recursion are provided by recurrences of the form $T(m, \vec{n}) \leq$ $T(m-1, \vec{n}) + q(\vec{n})$, where \vec{n} represents the other parameters and q is a (second-order) polynomial. Such T's grow polynomially in m. Thus, a polynomial bound on the depth of a **crec** recursion implies a polynomial bound on the recursion's total cost. If, say, two uses were allowed, the recurrences would be of the form $T(m, \vec{n}) \leq 2 \cdot T(m-1, \vec{n}) + q(\vec{n})$ and such T's can grow exponentially in m.

Tail recursions. We restrict **crec** terms to expressing just tail recursions. *Terminology:* The *tail terms* of an expression e consist of: (i) e itself, (ii) e', when $(\lambda x \cdot e')$ is a tail term, and (iii) e_1 and e_2 , when (if e_0 then e_1 else e_2) is a tail term. A *tail call* in e is a tail term of the form $(f e_1 \ldots e_k)$. Informally, a tail recursive definition is a function definition in

 $E ::= \dots | (\operatorname{crec} K (\lambda_r X \cdot E)) \qquad L ::= (\Box \diamond)^* | \diamond (\Box \diamond)^*$ $T_0 ::= \mathsf{N}_L \qquad T ::= \text{ the level } 0, 1, \text{ and } 2 \text{ simple types over } T_0$

which every recursive call is a tail call. Formally, we say that $(\operatorname{crec} a \ (\lambda_r f \cdot e))$ expresses a tail recursion when each occurrence of f in e is as the head of a tail call in e^{11} .

Simplicity is the foremost motivation for the restriction to tail recursions as they are easy to work with from both programming and complexity-theoretic standpoints. Additionally, tail recursion is a well-studied and widely-used universal form of recursion: there are *continuation passing style* translations of many program constructs into pure tail-recursive programs. (Reynolds [Rey93] provides a nice historical introduction.) Understanding the complexity theoretic properties of tail-recursive programs should lead to an understanding of a much more general set of programs.

5. Affine tiered recursion

Syntax. ATR (for *affine tiered recursion*) has the same syntax as PCF with three changes: (i) fix is replaced with **crec** as discussed in the previous section, (ii) the only variables allowed are those of type-levels 0 and 1, and (iii) the type system is altered as described below.

Types. The ATR types consist of *labeled base types* (T_0 from Figure 10) and the level 1 and 2 simple types over these base types. We first consider labels (L from Figure 10).

Labels. Labels are strings of alternating \diamond 's and \Box 's in which the rightmost symbol of a nonempty label is always \diamond . A label $\mathbf{a}_k \dots \mathbf{a}_0$ can be thought of as describing programoracle conversations: each symbol \mathbf{a}_i represents an action (\Box = an oracle action, \diamond = a program action) with the ordering in time being \mathbf{a}_0 through \mathbf{a}_k . Terminology: ε = the empty label, $\ell \leq \ell'$ means label ℓ is a suffix of label ℓ' , and $\ell \lor \ell'$ is the \leq -maximum of ℓ and ℓ' . Also let $succ(\ell)$ = the successor of ℓ in the \leq -ordering, $depth(\ell)$ = the number of \Box 's in ℓ , and, for each $d \in \omega$, $\Box_d = (\Box \diamond)^d$ and $\diamond_d = \diamond(\Box \diamond)^d$. Note: $depth(\Box_d) = depth(\diamond_d) = d$.

Labeled base types. The ATR base types are all of the form N_{ℓ} , where ℓ is a label. These base types are subtype-ordered by: $N_{\ell} \leq : N_{\ell'} \iff \ell \leq \ell'$. We thus have the linear ordering: $N_{\varepsilon} \leq : N_{\diamond} \leq : N_{\Box \diamond} \leq : N_{\diamond \Box \diamond} \leq : \cdots$, or equivalently, $N_{\Box_0} \leq : N_{\diamond_0} \leq : N_{\Box_1} \leq : N_{\diamond_1} \leq : \cdots$. Define $depth(N_{\ell}) = depth(\ell)$. N_{\Box_d} and N_{\diamond_d} are the depth-*d* analogues of the BCL'-types N_{norm} and N_{safe} , respectively. These types can be interpreted as follows.

- A N_{ε} -value is an ordinary base-type input or else is bounded by some prior (i.e., previously computed) N_{ε} -value.
- A N_{\diamond_d} -value is the result of a (type-2) polynomial-time computation over N_{\Box_d} -values or else is bounded by some prior N_{\diamond_d} -value.

¹¹Because of the one-use restriction, this simple definition of tail recursion suffices for this paper. For details on the more general notion see [Rey98, FWH01].

• A $N_{\Box_{d+1}}$ -value is the answer to a query made to a type-1 input on N_{\diamond_d} -values or else is bounded by some prior $N_{\Box_{d+1}}$ -value.

The N_{\Box_d} types are called *oracular* and the N_{\diamond_d} 's are called *computational*.

The ATR arrow types. These are just the level 1 and 2 simple types over the N_{ℓ} 's. The subtype relation \leq : is extended to these arrow types as in (2.1). Terminology: Let $shape(\sigma) =$ the simple type over N resulting from erasing all the labels. The *tail* of a type is given by:

$$tail(N_{\ell}) = N_{\ell}.$$
 $tail(\sigma \to \tau) = tail(\tau).$

Let $depth(\sigma) = depth(tail(\sigma))$. When $tail(\sigma)$ is oracular, we also call σ oracular and let $side(\sigma) = \Box$. When $tail(\sigma)$ is computational, we call σ computational and let $side(\sigma) = \diamond$.

Definition 6 (Predicative, impredicative, flat, and strict types). An ATR type γ is *pred*icative when γ is a base type or when $\gamma = (\sigma_1, \ldots, \sigma_k) \rightarrow \mathsf{N}_\ell$ and $tail(\sigma_i) \leq \mathsf{N}_\ell$ for each *i*. A type is *impredicative* when it fails to be predicative. An ATR type $(\sigma_1, \ldots, \sigma_k) \rightarrow \mathsf{N}_\ell$ is *flat* when $tail(\sigma_i) = \mathsf{N}_\ell$ for some *i*. A type is *strict* when it fails to be flat.

Examples: $N_{\varepsilon} \to N_{\diamond}$ is predicative whereas $N_{\diamond} \to N_{\varepsilon}$ is impredicative, and both are strict. Both $N_{\diamond} \to N_{\diamond}$ and $N_{\diamond} \to N_{\Box\diamond} \to N_{\diamond}$ are flat, but the first is predicative and the second impredicative. Recursive definitions tend to involve flat types.

Example 23 below illustrates that values of both impredicative and flat types require special restrictions in any sensible semantics of ATR. Our semantic restrictions for these types are made precise in §7 and §9 below. Here we give a quick sketch of these restrictions as they figure in definition of \propto , the *shifts-to* relation, used in the typing rules. For each impredicative type $(\vec{\sigma}) \rightarrow N_{\ell}$: if $\vdash f:(\vec{\sigma}) \rightarrow N_{\ell}$, then the value of $|f(\vec{x})|$ is essentially independent of the values of the $|x_i|$'s with $tail(\sigma_i) :\geq N_{\ell}$. For each flat type $(\vec{\sigma}) \rightarrow N_{\ell}$ (that for simplicity here we further restrict to be a level-1 computational type): if $\vdash f:(\vec{\sigma}) \rightarrow N_{\ell}$, then $|f(\vec{x})| \leq p + \bigvee\{|x_i| \mid tail(\sigma_i) = N_{\ell}\}$, where p is a second-order polynomial over elements of $\{|x_i| \mid tail(\sigma_i) \leq N_{\ell}\}$. (Compare this to the bound of Proposition 1(b).)

Typing rules. The ATR-typing rules are given in Figure 11. The rules Zero-I, Const-I, Int-Id-I, Subsumption, op-I, \rightarrow -I, and \rightarrow -E are essentially lifts from BCL (with one subtlety regarding \rightarrow -E discussed below). The if-I and down-I rules were motivated in §4. The remaining three rules Aff-Id-I and crec-I (that relate to recursions and the split type contexts) and Shift (that coerces types) require some discussion.

Affinely restricted variables and crec. Each ATR type judgment is of the form $\Gamma; \Delta \vdash e; \gamma$ where each type context is separated into two parts: a *intuitionistic zone* (Γ) and an *affine* zone (Δ). Γ and Δ are simply finite maps (with disjoint preimages) from variables to ATRtypes. By convention, "_" denotes an empty zone. Also by convention we shall restrict our attention to ATR type judgments in which each affine zone consists of at most one type assignment. (See Scholium 7(a).) In reading the rules of Figure 11, think of a variable in an affine zone as destined to be the recursor variable in some crec expression. An intuitionistic zone can be thought of as assigning types to each of the mundane variables.

Terminology: A variable f is said to be affinely restricted in $\Gamma; \Delta \vdash e; \sigma$ if and only if f is assigned a type by Δ or is λ_r -abstracted over in e.

$$\begin{aligned} & Zero-I: \quad \overline{\Gamma; \Delta \vdash \epsilon: \mathsf{N}_{\varepsilon}} \qquad Const-I: \quad \overline{\Gamma; \Delta \vdash k: \mathsf{N}_{\diamond}} \\ & Int-Id-I: \quad \overline{\Gamma; x: \sigma; \Delta \vdash x: \sigma} \qquad Aff-Id-I: \quad \overline{\Gamma; x: \gamma \vdash x: \gamma} \\ & Shift: \quad \frac{\Gamma; \Delta \vdash e: \sigma}{\Gamma; \Delta \vdash e: \tau} \quad (\sigma \propto \tau) \qquad Subsumption: \quad \frac{\Gamma; \Delta \vdash e: \sigma}{\Gamma; \Delta \vdash e: \tau} \quad (\sigma \leq : \tau) \\ & \mathsf{op}-I: \quad \frac{\Gamma; \Delta \vdash e: \mathsf{N}_{\diamond_d}}{\Gamma; \Delta \vdash (\mathsf{op} \ e): \mathsf{N}_{\diamond_d}} \qquad \mathsf{down}-I: \quad \frac{\Gamma; \Delta_0 \vdash e_0: \mathsf{N}_{\ell_0} \quad \Gamma; \Delta_1 \vdash e_1: \mathsf{N}_{\ell_1}}{\Gamma; \Delta_0, \Delta_1 \vdash (\mathsf{down} \ e_0 \ e_1): \mathsf{N}_{\ell_1}} \\ & \rightarrow -I: \quad \frac{\Gamma, x: \sigma; \Delta \vdash e: \tau}{\Gamma; \Delta \vdash (\lambda x \cdot e): \sigma \rightarrow \tau} \qquad \rightarrow -e: \quad \frac{\Gamma; \Delta \vdash e_0: \sigma \rightarrow \tau \quad \Gamma; _\vdash e_1: \sigma}{\Gamma; \Delta \vdash (e_0 \ e_1): \tau} \\ & \mathsf{if}-I: \quad \frac{\Gamma; _\vdash e_0: \mathsf{N}_{\ell} \quad \Gamma; \Delta_1 \vdash e_1: \mathsf{N}_{\ell'} \quad \Gamma; \Delta_2 \vdash e_2: \mathsf{N}_{\ell'}}{\Gamma; \Delta_1 \cup \Delta_2 \vdash (\mathsf{if} \ e_0 \ \mathsf{then} \ e_1 \ \mathsf{else} \ e_2): \mathsf{N}_{\ell'}} \\ & \mathsf{crec}-I: \quad \frac{\vdash K: \mathsf{N}_{\diamond} \quad \Gamma; f: \gamma \vdash e: \gamma}{\Gamma; _\vdash (\mathsf{crec} \ K \ (\lambda_r f \cdot e)): \gamma} \quad (\gamma \in \mathcal{R} \ \mathsf{and} \ TailPos(f, e)) \end{aligned}$$

where:

 $\mathcal{R} \stackrel{\text{def}}{=} \{ (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k) \to \mathbf{b} \mid \mathbf{b}_1 \text{ and each } \mathbf{b}_i \leq :\mathbf{b}_1 \text{ is oracular} \}.$ $TailPos(f, e) \stackrel{\text{def}}{=} [\text{ Each occurrence of } f \text{ in } e \text{ is as the head of a tail call }].$

Figure 11: ATR typing rules

The use of split type contexts is adapted from Barber and Plotkin's DILL [Bar96, BP97],¹² a linear typing scheme that permits a direct description of \rightarrow , the intuitionistic arrow of the conventional simple types. The key rule borrowed from DILL is $\rightarrow -E$ which forbids free occurrences of affinely restricted variables in the operand position of any intuitionistic application. This precludes the typing of crec-expressions containing subterms such as $\lambda_r f_{\bullet}(\lambda g_{\bullet}(g(g \epsilon)) f) \equiv_{\beta} \lambda_r f_{\bullet}(f(f \epsilon))$ where f is used multiple times.

The crec-*I* rule forbids any free occurrence of an affinely restricted variable; if such a free occurrence was allowed, it could be used any number of times through the crec-recursion. The crec-*I* rule requires that the recursor variable have a type $\gamma \in \mathcal{R}$ which in turn becomes the type of the crec-expression. The restrictions in \mathcal{R} 's definition (in Figure 11) are a more elaborate version of the typing restrictions for prn-expressions in BCL. When $\gamma = (N_{\Box_d}, \mathbf{b}_2, \ldots, \mathbf{b}_k) \rightarrow \mathbf{b} \in \mathcal{R}$, it turns out that \mathcal{R} 's restrictions limit a type- γ crec-expression to at most *p*-many recursions, where *p* is some fixed, depth-*d* second-order polynomial (Theorem 43). Excluding $N_{\diamond}, \ldots, N_{\diamond_{d-1}}$ in γ forbids depth $0, \ldots, d-1$ analogues of N_{safe} -parameters from figuring in the recursion, and consequently, the recursion cannot accumulate information that could change the value of *p* unboundedly.

Scholium 7.

(a) Judgments with with multiple type assignments in their affine zone are derivable. However, such a judgment is a dead end in the sense that crec-*I*, the only means to eliminate an affine-zone variable, requires a singleton affine zone.

¹²The discussion of DILL in [O'H03] is quite helpful.

(b) ATR has no explicit --types. *Implicitly*, a $(\lambda_r f \cdot e)$ subexpression is of type $\gamma - \gamma \gamma$ and crec-*I* plays roles of both --*I* and --*E*. ATR's very restricted use of affinity permits this --bypass.

(c) As mentioned in §4, the restriction to tail recursions in crec-I is in the interest of simplicity. In a follow-up to the present paper, we show how to relax this restriction to allow a broader range of affine recursions in ATR programs [DR07]. Dealing with this broader range of recursions turns out to require nontrivial extensions of the techniques of §§12–15 below.

Shift. The Shift rule covariantly coerces the type of a term to be deeper. Before stating the definition of the shifts-to relation (\propto), we first consider the simple case of shifting types of shape $\mathbb{N} \to \mathbb{N}$. The core idea is simply: $(\mathbb{N}_{\ell_1} \to \mathbb{N}_{\ell_0}) \propto (\mathbb{N}_{\ell'_1} \to \mathbb{N}_{\ell'_0})$ when $depth(\mathbb{N}_{\ell'_0}) - depth(\mathbb{N}_{\ell_0}) = depth(\mathbb{N}_{\ell'_1}) - depth(\mathbb{N}_{\ell_1}) \geq 0$. The motivation for this is that if pand q are second-order polynomials of depths d_p and d_q , respectively, and x is a base-type variable appearing in p that is treated as representing a depth- d_x value (with $d_x \leq d_q$), then p[x := q] is, in the worst case, of depth $d_p + (d_q - d_x)$. The full story for shifting level-1 types has to account of arbitrary arities, the sides of the component types, and impredicative and flat types, but even so it is still not too involved. Shifting level-2 types involves a new set of issues that we discuss after dealing with the level-1 case. Recall that $\max(\emptyset) = 0$.

Definition 8 (\propto , the shifts-to relation).

(a) We inductively define \propto by: $\mathsf{N}_{\Box_d} \propto \mathsf{N}_{\Box_{d'}}$ and $\mathsf{N}_{\diamond_d} \propto \mathsf{N}_{\diamond_{d'}}$ when $d \leq d'$; and $(\sigma_1 \rightarrow \cdots \rightarrow \sigma_k \rightarrow \mathsf{N}_{\ell_0}) \propto (\sigma'_1 \rightarrow \cdots \rightarrow \sigma'_k \rightarrow \mathsf{N}_{\ell'_0})$ when

(i) $\mathsf{N}_{\ell_0} \propto \mathsf{N}_{\ell'_0}, \, \sigma_1 \propto \sigma'_1, \, \dots, \, \sigma_k \propto \sigma'_k,$

(ii) $tail(\sigma_i) \stackrel{\circ}{=} \mathsf{N}_{\ell_0}$ implies $tail(\sigma'_i) = \mathsf{N}_{\ell'_0}$ for $i = 1, \ldots, k$, and

(iii) $depth(\mathsf{N}_{\ell'_0}) - depth(\mathsf{N}_{\ell_0}) \ge D((\vec{\sigma}) \xrightarrow{\circ} \mathsf{N}_{\ell_0}, \vec{\sigma}').$

(b) $D((\vec{\sigma}) \to \mathsf{N}_{\ell_0}, \vec{\sigma}') \stackrel{\text{def}}{=} \max\{ \operatorname{depth}(\sigma'_i) - \operatorname{depth}(\sigma_i) \mid \sigma_i \leq :\mathsf{N}_{\ell_0} \}, \text{ for } \sigma_1 \propto \sigma'_1, \ldots, \sigma_k \propto \sigma'_k \text{ where each } \sigma_i \text{ and } \sigma'_i \text{ is a base type. (See Definition 9 for the general definition of D.)}$

For base types: $N_{\ell} \propto N_{\ell'}$ if and only if $depth(N_{\ell}) \leq depth(N_{\ell'})$ and $side(N_{\ell}) = side(N_{\ell'})$. It follows from this and condition (i) that no type (or component of a type) can change sides as a result of a shift.

For level-1 types: Condition (i) says that the component types on the right are either the same as or else deeper versions of the corresponding types on the left. Condition (ii) preserves flatness (which is critical in level-2 shifting). Condition (iii) is just the core idea stated above. Note that the max in Definition 8(b) includes only types $\leq: N_{\ell_0}$. This is because as remarked above, if $\sigma_i :\geq N_{\ell_0}$, then the *i*-th argument has essentially no effect on the size of the $N_{\ell'_0}$ -result.

Example: Consider the problem: $\Gamma; _\vdash f(f(x)):?$, where $\Gamma = f: \mathsf{N}_{\diamond} \to \mathsf{N}_{\Box\diamond}, x: \mathsf{N}_{\diamond}$. Using $\to -E$ and *Subsumption*, we derive $\Gamma; _\vdash f(x): \mathsf{N}_{\diamond\Box\diamond}$. Using *Shift* we derive $\Gamma; _\vdash f: \mathsf{N}_{\diamond\Box\diamond} \to \mathsf{N}_{\Box\diamond\Box\diamond}$. Using $\to -E$ again we obtain $\Gamma; _\vdash f(f(x)): \mathsf{N}_{\Box\diamond\Box\diamond}$ as desired.

Now let us consider shifting level-2 types. Suppose we want to shift $(\mathsf{N}_{\Box_0} \to \mathsf{N}_{\Box_1}) \to \mathsf{N}_{\Box_3}$ to some type of the form $(\mathsf{N}_{\Box_0} \to \mathsf{N}_{\Box_2}) \to \mathsf{N}_{\Box_d}$. What should the value of d be? Suppose $f: \mathsf{N}_{\Box_0} \to \mathsf{N}_{\Box_1}$. Without using subsumption, building a term of type N_{\Box_3} from f requires nesting applications of f (using type-1 shifts). The longest chain of such depth-increasing

$$undo((\mathsf{N}_{\ell_1},\ldots,\mathsf{N}_{\ell_k})\to\mathsf{N}_{\ell_0},\mathsf{N}_{\ell}) \stackrel{\text{def}}{=}$$

 $\begin{cases} \text{undefined,} & \text{if } (\mathsf{N}_{\ell_1}, \dots, \mathsf{N}_{\ell_k}) \to \mathsf{N}_{\ell_0} \text{ is flat or } \ell_0 > \ell; \\ \mathsf{N}_{\ell' \oplus \ell''}, & \text{otherwise, where } \ell' = \max\{\,\ell_i \mid \ell_i < \ell_0\,\} \text{ and } \ell'' \text{ is the} \\ & \text{suffix of } \ell \text{ following the leftmost occurrence of } \ell_0 \text{ in } \ell. \end{cases}$

Figure 12: The definition of undo

applications is 3.¹³ When the argument type $N_{\Box_0} \rightarrow N_{\Box_1}$ is shifted to $N_{\Box_0} \rightarrow N_{\Box_2}$, each application of this argument now ups the depth by an additional +1. So, the largest depth that can result from the change is $d = 3 + 3 \cdot 1 = 6$. When shifting $(\vec{\sigma}) \rightarrow N_{\ell}$ to some $(\vec{\sigma}') \rightarrow N_{\ell'}$ with each σ_i and σ'_i a level-1 type, to determine ℓ' we must: (a) determine all the ways a N_{ℓ} value could be built by a chain of depth-increasing applications of arguments of the types $\vec{\sigma}$, (b) for each of these ways, figure the increase in the depth of the N_{ℓ} -value when each σ_i -argument is replaced by its σ'_i version, and (c) compute the maximum of these increases. To help in this, we introduce *undo* in Figure 12. *Example:* For d > 0, $undo(N_{\Box_0} \rightarrow N_{\Box_1}, N_{\Box_d}) = undo(N_{\Box_0} \rightarrow N_{\Box_1}, N_{\diamond_d}) = N_{\Box_{d-1}}$. To compute $undo(\tau, N_{\ell})$, one determines if a type- τ argument could be used in a chain of depth-increasing applications that build a N_{ℓ} value, and if so, one figures (in terms of ℓ) where a leftmost application of such an argument could occur, and returns the \leq :-largest type of the arguments of this application. (It is straightforward to prove that *undo* behaves as claimed.) **N.B.** If $undo(\gamma, N_{\ell})$ is defined, then $undo(\gamma, N_{\ell}) \lneq N_{\ell}$. We now define:

Definition 9 (*D* for level-2 types). Suppose $\sigma_1 \propto \sigma'_1, \ldots, \sigma_k \propto \sigma'_k$.

(a) $D((\vec{\sigma}) \to \mathsf{N}_{\ell}, \vec{\sigma}') \stackrel{\text{def}}{=} \max(\{ depth(\sigma'_i) - depth(\sigma_i) + D((\vec{\sigma}) \to undo(\sigma_i, \mathsf{N}_{\ell}), \vec{\sigma}') \}$ $undo(\sigma_i, \mathsf{N}_{\ell}) \text{ is defined }\}), \text{ when each } \sigma_i \text{ and } \sigma'_i \text{ is level-1.}$

(b) $D((\vec{\sigma}) \to \mathsf{N}_{\ell}, \vec{\sigma}') \stackrel{\text{def}}{=} D((\vec{\sigma})_0 \to \mathsf{N}_{\ell}, (\vec{\sigma}')_0) + D((\vec{\sigma})_1 \to \mathsf{N}_{\ell}, (\vec{\sigma}')_1)$, when the σ_i 's contain both level-0 and level-1 types and where $(\vec{\gamma})_i$ denotes the subsequence of level-*i* types of $\vec{\gamma}$.

The recursion of Definition 9(a) determines maximum increase in depth as outlined above. Since applications amount to simultaneous substitutions, the contributions of the level-0 and level-1 argument shifts are independent. Thus Definition 9(b)'s formula suffices for the general case. *Example:* See the discussion below of *fcat* from Figure 13.

Now let us consider the reason behind condition (ii) in Definition 8(a). A term of a flat type can be used an *arbitrary* number of times in constructing a value. Consequently, if Definition 8(a) had allowed flat level-1 types (which increase the depth by 0) to be shifted to strict level-1 types (which increase the depth by a positive amount), then it would have been impossible to bound the depth increase of shifts involving arguments of flat types.

Some examples. Figure 13 contains five sample programs. These examples use the syntactic sugar of the let and letrec constructs.¹⁴ The first three programs and their typing are

¹³Note that the outer two of these three applications must involve shifting the type of the argument. Also, informally, in $f(f(f(\mathsf{down}(f(f(f(f(\epsilon)))), \epsilon))))$ only the outer three applications of f count as a chain of depth-increasing applications because of the drop in depth caused by the down. Formally, no shadowed (Definition 29) application can be in a depth-increasing chain.

¹⁴Where (letrec f = e' in e) $\stackrel{\text{def}}{\equiv} e[f := (\operatorname{crec} \underline{0} (\lambda_r f \cdot e'))]$ and let is as in footnote 10.

$$\begin{aligned} & \text{reverse: } \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} = // \text{ reverse } \mathbf{a}_{1} \dots \mathbf{a}_{k} = \mathbf{a}_{k} \dots \mathbf{a}_{1}. \\ & \lambda w. \text{ letrec } f: \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = \\ & \lambda b, x, r. \text{ if } (\mathsf{f}_{0} x) & \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{0} r) \\ & \text{ else } r & \text{ in } f w w \epsilon \end{aligned}$$

$$\begin{aligned} & prn: (\mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond}) \to \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} = // \text{ See } (1.1). \\ & \lambda e, y. \text{ letrec } f: \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = \\ & \lambda b, x, z, r. \text{ if } (\mathsf{t}_{0} x) & \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{0} z) (e (\mathsf{c}_{0} z) r) \\ & \text{ else } if (\mathsf{t}_{1} x) \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{1} z) (e (\mathsf{c}_{1} z) r) \\ & \text{ else } if (\mathsf{t}_{1} x) \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{1} z) (e (\mathsf{c}_{1} z) r) \\ & \text{ else } if (\mathsf{t}_{1} x) \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{1} z) (e (\mathsf{c}_{1} z) r) \\ & \text{ else } if (\mathsf{t}_{1} x) \text{ then } f b (\mathsf{d} x) (\mathsf{c}_{1} z) (e (\mathsf{c}_{1} z) r) \\ & \text{ else } r \\ & \text{ in } f y (reverse y) \epsilon (e \epsilon \epsilon) \end{aligned}$$

$$\begin{aligned} & cat: \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = // cat w x = w \oplus x \text{ as before.} \\ & \lambda w, x. \text{ let } f: \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = \\ & \lambda y, z. \text{ if } (\mathsf{t}_{0} y) \text{ then } (\mathsf{c}_{0} z) \text{ else } if (\mathsf{t}_{1} y) \text{ then } (\mathsf{c}_{1} z) \text{ else } x \\ & \text{ in } prn \ f w \end{aligned}$$

$$\begin{aligned} & fcat: (\mathsf{N}_{\diamond} \to \mathsf{N}_{\Box}) \to \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond \Box} = // fcat f \ a_{1} \dots a_{k} = (f \ a_{1} \dots a_{k}) \oplus \\ & \lambda f, x. \text{ let } e: \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond \Box} \to \mathsf{N}_{\diamond \Box} \to \mathsf{N}_{\diamond \Box} = // (f \ a_{2} \dots a_{k}) \oplus \dots \oplus (f \ a_{k}) \oplus (f \epsilon) \\ & \lambda y, r. (cat (f y) r) \\ & \text{ in } prn \ e x \end{aligned}$$

$$\begin{aligned} & findk: (\mathsf{N}_{\diamond} \to \mathsf{N}_{\Box}) \to \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\varepsilon} = // See \ (5.1) \\ & \lambda f, x. \text{ letree } h: \mathsf{N}_{\Box} \to \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\varepsilon} = // Invariant: k \leq len(m) \text{ and } |m| \leq |f|(|x|) \\ & \lambda m, k. \text{ if } k = x \text{ then } k \\ & \text{ else if } k = (len m) \text{ then } k \\ & \text{ else if } k = (len m) \text{ then } k \\ & \text{ else } h (max (f (k + 1)) m) (\text{ down } (k + 1) x) \\ & \text{ in } h (f \epsilon) \epsilon \end{aligned}$$

$$\mathsf{Figure } 13: \mathsf{ATR} \text{ versions of } reverse, prn, \ cat, \ fcat, \ and \ findk \end{aligned}$$

all straightforward. For the typing of *fcat*, *cat*'s type is shifted to $N_{\Box_1} \rightarrow N_{\diamond_1} \rightarrow N_{\diamond_1}$ and *prn*'s type is shifted to $(N_{\diamond_0} \rightarrow N_{\diamond_1} \rightarrow N_{\diamond_1}) \rightarrow N_{\Box_0} \rightarrow N_{\diamond_1}$. The final program computes

$$\lambda f \in (\mathbb{N} \to \mathbb{N}), x \in \mathbb{N} \, \left\{ \begin{array}{l} (\mu k < x) \left[k = \max_{i \le k} len(f(i)) \right], & \text{if such a } k \text{ exists;} \\ x, & \text{otherwise;} \end{array} \right.$$
(5.1)

where len(z) = the dyadic representation of the length of z. This is a surprising and subtle example of a BFF due to Kapron [Kap91] and was a key example that lead to the Kapron-Cook Theorem [KC96]. In *findk*, we assume we have: a type- $(N_{\Box_1} \rightarrow N_{\Box_1} \rightarrow N_{\Box_1})$ definition of $(x, y) \mapsto [\underline{1}$, if x = y; ϵ , otherwise], a type- $(N_{\Box_1} \rightarrow N_{\Box_1})$ definition of *len*, a type- $(N_{\Box_1} \rightarrow N_{\Box_1} \rightarrow N_{\Box_1})$ definition of *max*, and a type- $(N_{\diamond_0} \rightarrow N_{\diamond_0})$ definition of $x \mapsto x+1$. Filling in these missing definitions is a straightforward exercise. A more challenging exercise is to define (5.1) via *prn*'s.

Semantics. The CEK machine of (§11.1) provides an operational semantics of ATR. For a denotational semantics we *provisionally* take the obvious modification of PCF's \mathcal{V} -semantics. (\mathcal{V} was introduced in §2.7.) Example 23 illustrates some inherent difficulties with \mathcal{V} as a

semantics for ATR. We shall circumvent these difficulties by some selective pruning of \mathcal{V} in §7 and §9.

Some syntactic properties.

Definition 10 (Use). If variable x fails to occur free in expression e, then uses(x, e) = 0; otherwise uses(x, e) is given by:

$$\begin{split} uses(x,x) &= 1. \qquad uses(x,(\texttt{op}\ e_0)) = uses(x,(\lambda y \cdot e_0)) = uses(x,e_0). \\ uses(x,(\texttt{down}\ e_0\ e_1)) = uses(x,(e_0\ e_1)) = uses(x,e_0) + uses(x,e_1). \\ uses(x,(\texttt{if}\ e_0\ \texttt{then}\ e_1\ \texttt{else}\ e_2)) = uses(x,e_0) + uses(x,e_1) \lor uses(x,e_2). \\ uses(x,(\texttt{crec}\ K\ (\lambda_r f \cdot e_0)))) = \dagger (\equiv unbounded). \end{split}$$

By convention, $a < \dagger$ and $a + \dagger = \dagger + a = a \lor \dagger = \dagger \lor a = \dagger$ for each $a \in \mathbb{N}$.

Lemma 11 (One-use). If Γ ; $f: \gamma \vdash e: \gamma \text{ or } \Gamma$; $\vdash (\operatorname{crec} k (\lambda_r f \cdot e)): \gamma$, then $uses(f, e) \leq 1$.

Lemma 12 (Subject reduction). If $\Gamma; \Delta \vdash e: \gamma$ and $e \beta \eta$ -reduces to e', then $\Gamma; \Delta \vdash e': \gamma$.

Lemma 13 (Unique typing of subterms). If Γ ; $\Delta \vdash e:\sigma$, then each occurrence of a subterm in e has a uniquely assignable type that is consistent with Γ ; $\Delta \vdash e:\sigma$.

Lemma 14. $\Gamma; \Delta \vdash \lambda \vec{x} \cdot e: (\vec{\sigma}) \to \mathsf{N}_{\ell}$ if and only if $\Gamma, \vec{x}: \vec{\sigma}; \Delta \vdash e: \mathsf{N}_{\ell}$.

Lemma 11 follows from a straightforward structural induction on judgment derivations. The proof of Lemma 12 is an adaptation the argument for [Pie02, Theorem 15.3.4]. The proof of Lemma 13 is also an adaptation of standard arguments. We make frequent, implicit use of Lemma 13 below. Lemma 14 is a reality check on the definition of \propto . The proof of this is a completely standard induction on derivations *except* in the case where the last rule used in deriving $\Gamma; \Delta \vdash \lambda \vec{x}.e: (\vec{\sigma}) \rightarrow N_{\ell}$ is *Shift*. The argument for this case is an induction on the structure of e, where application is the key subcase. There one simply checks that our definition of \propto correctly calculates upper bounds on the increase in depth.

ATR's computational limitations and capabilities. The major goals of the rest of the paper are to establish type-level 2 analogues of Propositions 1, 2, and 3 for ATR. We shall first prove Theorem 43, a polynomial size-boundedness result for ATR. The groundwork for this result will be the investigation of second-order size-bounds in the next few sections.

Remark 15 (Related work). As noted in §1, ramified types based on Bellantoni and Cook's ideas, higher types, and linear types are common features of work on implicit complexity (see Hofmann's survey [Hof00]), but most of that work has focused on guaranteeing complexity of type-level 1 programs. The ATR type system is roughly a refinement of the type systems of [IKR01, IKR02] which were constructed to help study higher-type complexity classes. Also, the type systems of this paper and [IKR01, IKR02] were greatly influenced by Leivant's elegant ramified type systems [Lei95, Lei94]. We note that in [Lei03] Leivant proposes a formalism that uses intersection types to address the same problems dealt with by our *Shift* rule (e.g., how to type f(f(x))).

Figure 14: Additional typing rules for the second-order polynomials under the size types

Remark 16 (Pragmatic predicativity). Many of the formalisms based on Bellantoni and Cook's ideas are predicative in the sense of Proposition 4—no information about "safe values" can influence "normal values." Two principles followed in this paper are: (i) The ramification of data (e.g., the normal/safe distinction) and the complexity it adds to the type system is something we will put up with to control the size of values; (ii) however, if there is a good reason to cut through the ramification while still controlling sizes, then we will happily do so. As a consequence of (i), our type system for second-order polynomial size-bounds is strictly predicative. As a consequence of (ii), ATR's type system includes the if-I and down-I rules and impredicative types to handle examples like F_2 of (4.1).

There is a price for the down construct—its use tends to complicate correctness arguments for algorithms. For example, consider the subexpression (down (k + 1) x) in the ATR-program for *findk* in Figure 13. The purpose of the down is to guarantee to the type system that the subexpression's value is small (e.g., $\leq |x|$). The correctness of the algorithm depends critically on the easy observation that, in any run of the program, the value of the subexpression will always be k + 1. This is common in expressing algorithms in ATR—one knows that a value is small, but an application of down is needed to convince the type-system of this. As a result the correctness proof needs a lemma showing that original value is indeed small and the down expression does not change the value. Thus our use of down and (mild) impredicativity is a compromise between the simplicity, but restrictiveness, of predicative systems and the richer, but more complex, type systems that permit finer reasoning about size.¹⁵

6. Size bounds

6.1. The second-order polynomials under the size types. To work with size bounds, we introduce the *size types* and a typing of second-order polynomials under these types. The size types parallel the intuitionistic part of ATR's type system.

Definition 17.

(a) For each ATR type σ , let $|\sigma| = \sigma[N := T]$. (E.g., $|N_{\epsilon} \to N_{\diamond}| = T_{\epsilon} \to T_{\diamond}$.) These $|\sigma|$'s are the *size types*. All the ATR-types terminology and operations (e.g., *shape*, *tail*, \leq :, ∞ , etc.) are defined analogously for size types.

¹⁵Hofmann's work on non-size-increasing functions [Hof03, Hof02] provides a nice example of a type system for fine control of sizes, but that system is not helpful in dealing with the F_2 or findk examples.

(b) The typing rules for the second-order polynomials under the size types consist of Id-I, $\rightarrow -I$, and $\rightarrow -E$ from Figure 2 and the rules of Figure 14.

Recall the \mathcal{L} -semantics for second-order polynomials introduced in §2.12. We provisionally take $\mathcal{L}[\![\sigma]\!] = \mathcal{L}[\![shape(\sigma)]\!]$ and define $\mathcal{L}[\![\Sigma \vdash p:\sigma]\!]$ as before. Later, a pruned version of the \mathcal{L} -semantics will end up as our intended semantics for the second-order polynomials to parallel our pruning of the \mathcal{V} -semantics for ATR.

The following definition formalizes what it means for an ATR expression to be polynomially size-bounded. **N.B.** This definition heavily overloads the "length of" notation, $|\cdot|$. In particular, if x is an ATR variable, we treat |x| as a size-expression variable. Definition 18(c) is based on a similar notion from [IKR02].

Definition 18. Suppose Γ ; $\Delta \vdash e: \sigma$ is an ATR-type judgment.

(a) $|\Gamma; \Delta| \stackrel{\text{def}}{=} \{ |x| \mapsto |\sigma| \colon (\Gamma; \Delta)(x) = \sigma \}.$

(b) For each $\rho \in \mathcal{V}[\Gamma; \Delta]$, define $|\rho| \in \mathcal{L}[\Gamma; \Delta]$ by $|\rho|(|x|) = |\rho(x)|$.¹⁶

(c) We say that the second-order polynomial p bounds the size of e (or, p is a size-bound for e) with respect to $\Gamma; \Delta$ when $|\Gamma; \Delta| \vdash p: |\sigma|$ and $|\mathcal{V}[\![e]\!] \rho| \leq_{|\sigma|} \mathcal{L}[\![p]\!] |\rho|$ for all $\rho \in \mathcal{V}[\![\Gamma; \Delta]\!]$. (The "with respect to" clause is dropped when it is clear from context.)

Lemmas 19, 21, and 22 below note a few basic properties of the second-order polynomials under the size types. Lemma 21 connects the depth of a second-order polynomial p and the depths of the types assignable to p. Lemmas 19 and 20 follow by proofs similar to those for Lemmas 12 and 14. Lemma 21's proof is a straightforward induction on judgment derivations, and Lemma 22 is just an observation. *Terminology:* Inductively define $\underline{0}_{\gamma}$ by: $\underline{0}_{\mathsf{T}_{\ell}} = \underline{0}$ and $\underline{0}_{\sigma \to \tau} = \lambda x \cdot \underline{0}_{\tau}$. By abuse of notation, we often write $\underline{0}_{\gamma}$ for $\mathcal{L}\llbracket\vdash \underline{0}_{\gamma}: \gamma \rrbracket$

Lemma 19 (Subject Reduction). Suppose $\Sigma \vdash p: \sigma$ and $p \ \beta \eta$ -reduces to p'. Then $\Sigma \vdash p': \sigma$.

Lemma 20. $\Sigma \vdash \lambda \vec{x} \cdot p: (\vec{\sigma}) \to \mathsf{T}_{\ell}$ if and only if $\Sigma, \vec{x}: \vec{\sigma} \vdash p: \mathsf{T}_{\ell}$.

Lemma 21 (Label Soundness). Suppose $\Sigma \vdash p: \sigma$ has a derivation in which the only types assigned by contexts are from $\{\mathsf{N}_{\varepsilon}\} \cup \{(\mathsf{N}_{\diamond}^k) \to \mathsf{N}_{\Box\diamond} \mid k > 0\}$. Then $depth(p) \leq depth(\sigma)$.

Lemma 22. $\underline{0}_{\gamma}$ is the least element of $\mathcal{L}[\![\gamma]\!]$ under the pointwise ordering.

6.2. Semantic troubles. The naïve (and *false!*) ATR-analogue of Proposition 1 is:

For each $\Gamma; \Delta \vdash e: \sigma$, there is a p_e that bounds the size of e with respect to $\Gamma; \Delta$.

Example 23 illustrates the problems with this. **N.B.** If the definition of BCL had allowed unrestricted free variables of type-level 1, the problems of Example 23 would have occurred in that setting too.

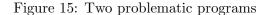
Example 23. Let e_1 and e_2 be as given in Figure 15, let *prn* be as in Figure 13, and let *dup* be an ATR-version of the definition in Figure 8.

(a) Suppose $\Gamma_1 = g_1: \mathbb{N}_{\diamond} \to \mathbb{N}_{\varepsilon}$ and $\rho_1 = \{g_1 \mapsto \lambda z \in \mathbb{N} \cdot z\}$. Then $|\mathcal{V}[\![e_1]\!] \rho_1| = \lambda n \in \omega \cdot n^{2^n}$. Note $|\rho_1(g_1)| = \lambda n \in \omega \cdot n$ is a polynomial function. The problem is that $\rho_1(g_1) = \lambda x \in \mathbb{N} \cdot x$ subverts the intent of the type-system by allowing an unrestricted flow of information about "safe" values into "normal" values.¹⁷

¹⁶**N.B.** The $|\cdot|$ in "|x|" is syntactic, whereas the $|\cdot|$ in " $|\rho|$ " and " $|\rho(x)|$ " are semantic.

 $^{^{17}}$ By using a similar trick and the full power of crec, one can write nonterminating ATR programs.

 $\begin{array}{ll} e_1 \colon \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} &= \\ \lambda w \, . \, \mathsf{let} \, h_1 \colon \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = \lambda x, y \, . \, \mathsf{if} \, x \neq \epsilon \, \mathsf{then} \, (dup \, (g_1 \, y) \, (g_1 \, y)) \, \mathsf{else} \, w \\ \mathsf{in} \, prn \, h_1 \, w \end{array}$ $e_2 \colon \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} = \\ \lambda w \, . \, \mathsf{let} \, h_2 \colon \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} \to \mathsf{N}_{\diamond} = \lambda x, y \, . \, \mathsf{if} \, x \neq \epsilon \, \mathsf{then} \, (g_2 \, y) \, \mathsf{else} \, w \\ \mathsf{in} \, prn \, h_2 \, w \end{array}$



(b) Suppose $\Gamma_2 = g_2: \mathbb{N}_{\diamond} \to \mathbb{N}_{\diamond}$ and $\rho_2 = \{g_2 \mapsto \lambda z \in \mathbb{N} \cdot z \oplus z\}$. Then $|\mathcal{V}[\![e_2]\!] \rho_2| = \lambda n \in \omega \cdot n2^n$. Note $|\rho_2(g_2)| = \lambda n \in \omega \cdot 2n$ is a polynomial function. The problem is that $\rho_2(g_2) = \lambda y \in \mathbb{N} \cdot y \oplus y$ subverts the fundamental restriction on the sizes of "safe" values in growth-rate bounds as in Proposition 1(b).

The problem of Example 23(a) is addressed in §7 by pruning the \mathcal{L} - and \mathcal{V} -semantics to restrict impredicative-type values. The problem of Example 23(b) is addressed in §9 by further pruning to restrict flat-type values.

7. Impredicative types and nearly well-foundedness

Failing to restrict impredicative-type values leads to problems like the one of Example 23(a). These problems can be avoided by requiring that each impredicative-type value have a length that is *nearly well-founded*.

Definition 24. A $t \in \mathcal{L}[\![\gamma]\!]$ is γ -well-founded when $\gamma = \mathsf{T}_{\ell}$ or else $\gamma = (\sigma_1, \ldots, \sigma_k) \to \mathsf{T}_{\ell}$ and, for each *i* with $tail(\sigma_i) :\geq \mathsf{T}_{\ell}$, the function *t* has no dependence on its *i*-th argument. A *t* is nearly γ -well-founded when there is a γ -well-founded *t'* such that $t \leq t'$.

Remark 25. Why *nearly* well-founded? The natural sources of ATR-terms with impredicative types are the if-then-else and down constructs. Let $c = \lambda x, y, z$ (if x then y else z) and $d = \lambda x, y$ (down xy), where $\vdash c: (N_{\ell}, N_{\ell'}, N_{\ell'}) \rightarrow N_{\ell'}$, $\vdash d: (N_{\ell}, N_{\ell'}) \rightarrow N_{\ell'}$, and $\ell > \ell'$. Thus $|c| \in \mathcal{L}[[|(N_{\ell}, N_{\ell'}, N_{\ell'}) \rightarrow N_{\ell'}|]]$ and $|d| \in \mathcal{L}[[|(N_{\ell}, N_{\ell'}) \rightarrow N_{\ell'}|]]$. Neither |c| nor |d| is wellfounded since $|c| = \lambda k, m, n$ (m, if k = 0; n, otherwise) and $|d| = \lambda k, m$ min(k, m). However, both |c| and |d| are nearly well-founded as $|c| \leq \lambda k, m, n \cdot (m \lor n)$ and $|d| \leq \lambda k, m \cdot m$.

Lemma 26. Suppose $\Sigma \vdash p: \sigma$, $\rho \in \mathcal{L}[\![\Sigma]\!]$, and $\rho(x)$ is nearly $\Sigma(x)$ -well-founded for each $x \in \operatorname{preimage}(\Sigma)$. Then $\mathcal{L}[\![p]\!] \rho$ is nearly σ -well-founded.

Lemma 26 follows by a straightforward induction and indicates that a semantics for the second-order polynomials based on nearly well-foundedness will be well defined. *Terminology.* The *restriction* of $f \in (X_1, \ldots, X_k) \to Y$ to $(X'_1, \ldots, X'_k) \to Y$ (where $X'_1 \subseteq X_1, \ldots, X'_k \subseteq X_k$) is $\lambda x_1 \in X'_1, \ldots, x_k \in X'_k \cdot f(x_1, \ldots, x_k)$.

Definition 27 (The nearly well-founded semantics).

(a) Inductively define $\mathcal{L}_{nwf}[\![\gamma]\!]$ by: $\mathcal{L}_{nwf}[\![\mathsf{T}_{\ell}]\!] = \omega$. For $\gamma = (\sigma_1, \ldots, \sigma_k) \to \mathsf{T}_{\ell}$, $\mathcal{L}_{nwf}[\![\gamma]\!]$ is the restriction to $(\mathcal{L}_{nwf}[\![\sigma_1]\!], \ldots, \mathcal{L}_{nwf}[\![\sigma_k]\!]) \to \mathcal{L}_{nwf}[\![\mathsf{T}_{\ell}]\!]$ of the γ -nearly well-founded elements of $\mathcal{L}[\![\gamma]\!]$. Define $\mathcal{L}_{nwf}[\![\Sigma]\!]$ and $\mathcal{L}_{nwf}[\![\Sigma \vdash p:\gamma]\!]$ in the standard way.

(b) Inductively define $\mathcal{V}_{nwf}[\![\gamma]\!]$ by: $\mathcal{V}_{nwf}[\![\mathsf{N}_{\ell}]\!] = \mathbb{N}$. For $\gamma = (\sigma_1, \ldots, \sigma_k) \to \mathsf{N}_{\ell}$, $\mathcal{V}_{nwf}[\![\gamma]\!]$ is the restriction to $(\mathcal{V}_{nwf}[\![\sigma_1]\!], \ldots, \mathcal{V}_{nwf}[\![\sigma_k]\!]) \to \mathcal{V}_{nwf}[\![\mathsf{N}_{\ell}]\!]$ of the $f \in \mathcal{V}[\![\gamma]\!]$ with $|f| \in \mathcal{L}_{nwf}[\![|\gamma|]\!]$. Define $\mathcal{V}_{nwf}[\![\Gamma; \Delta]\!]$ and $\mathcal{V}_{nwf}[\![\Gamma; \Delta \vdash E; \gamma]\!]$ in the standard way. (c) We write $p =_{nwf} p'$ when $\mathcal{L}_{nwf} \llbracket \Sigma \vdash p : \gamma \rrbracket |\rho| = \mathcal{L}_{nwf} \llbracket \Sigma \vdash p' : \gamma \rrbracket |\rho|$ for all $|\rho| \in \mathcal{L}_{nwf} \llbracket \Sigma \rrbracket$. We define $\leq_{nwf}, \geq_{nwf}, \ldots$ analogously.

There is still a problem with impredicative-type values. In deriving closed-form upper bounds on recursions, we often need a well-founded upper bound on the value of a variable of an impredicative type. There is no effective way to obtain such bound. We thus do the next best thing: give a canonical such upper bound a name and work with that name.

Definition 28. We add a new combinator, \mathbf{p} , to the second-order polynomials such that $\mathcal{L}_{nwf}[\![\Sigma \vdash (\mathbf{p} p): \gamma]\!] |\rho| = \text{the least } \gamma\text{-well-founded upper bound on } \mathcal{L}_{nwf}[\![\Sigma \vdash p: \gamma]\!] |\rho|$. (See Figure 17 for \mathbf{p} 's typing rule.) For each variable x, we abbreviate $(\mathbf{p} x)$ by \mathbf{p}_x .

The choice **p** makes is analogous to choice of a in the situation where one knows $f \in O(n)$ and picks the least $a \in \omega$ such that $f(n) \leq a \cdot (n+1)$ for all $n \in \omega$. In most uses, **p**_x's are destined to be substituted for by concrete, well-founded terms.

To help work with terms involving impredicative types we introduce:

Definition 29 (Shadowing). Suppose $\Sigma \vdash p:\sigma$. An occurrence of a subterm r of p is shadowed when the occurrence properly appears within another shadowed occurrence or else the occurrence has an enclosing subexpression (t r) where the occurrence of t is of an impredicative type $\sigma \to \tau$ with $tail(\sigma) :\geq tail(\tau)$. A variable x is a shadowed free variable for p when all of x's free occurrences in p are shadowed; otherwise x is an unshadowed free variable for p.

8. SAFE UPPER BOUNDS

The restriction to the \mathcal{V}_{nwf} -semantics solves the problem with impredicative types, but not the problem with flat types. To work towards a solution of this later problem, in this section we introduce the notion of a safe second-order polynomial (Definition 30) and show that any expression (in a simplification of ATR) that does not involve flat-type variables has a safe upper bound. The next section proposes a solution to the flat-type problem: that each flat-type length must have a safe upper bound. Theorem 43, in §10, shows that this proposed solution does indeed work. *Convention:* In this section **b**, γ , σ , and τ range over size types. In writing $p = (x p_1 \dots p_k)$, we mean x is a variable and, when k = 0, p = x.

Definition 30 (Strictness, chariness, and safety). Suppose $\Sigma \vdash p: \gamma$.

(a) We say that p is **b**-strict with respect to Σ when $tail(\gamma) \leq \mathbf{b}$ and every unshadowed free-variable occurrence in p has a type with tail $\leq \mathbf{b}$.

(b) We say that p is **b**-chary with respect to Σ when $\gamma = \mathbf{b}$ and either (i) $p = (x \ q_1 \cdots q_k)$ with each q_i **b**-strict or (ii) $p = p_1 \lor \cdots \lor p_m$, where each p_i satisfies (i). (Note that $\underline{0}$ sneaks in as **b**-chary; take m = 0 in (ii).)

(c) We say that p is γ -safe with respect to Σ if and only if

(i) when $\gamma = \mathsf{T}_{\Box_d}$, then $p =_{\text{nwf}} q \lor r$ where q is γ -strict and r is γ -chary,

(ii) when $\gamma = \mathsf{T}_{\diamond_d}$, then $p =_{\text{nwf}} q + r$ where q is a γ -strict and r is γ -chary r, and

(iii) when $\gamma = \sigma \rightarrow \tau$, then (px) is τ -safe with respect to $\Sigma, x: \sigma$.

With the above notions, we drop the "with respect to Σ " when Σ is clear from context. *Examples:* Recall the bound $p + \bigvee_{j=1}^{n} |y_j|$ of Proposition 1(b). In terms of the sizetypes, the subterm p is T_{\diamond} -strict, the subterm $\bigvee_{j=1}^{n} |y_j|$ is T_{\diamond} -chary, and hence, $p + \bigvee_{j=1}^{n} |y_j|$

$$\mathsf{s}\text{-}I: \quad \frac{\Sigma \vdash s: \mathsf{T}_{\diamond_d}}{\Sigma \vdash (\mathsf{s} \ s): \mathsf{T}_{\diamond_d}} \qquad \qquad \mathsf{R}\text{-}I: \quad \frac{\Sigma_0 \vdash s_0: \mathsf{T}_{\ell'} \to \mathsf{T}_{\ell'}}{\Sigma_0 \cup \Sigma_1 \cup \Sigma_2 \vdash (\mathsf{R} \ s_0 \ s_1 \ s_2): \mathsf{T}_{\ell'}} \quad \begin{pmatrix} \ell' = \\ succ(\ell) \end{pmatrix}$$

Figure 16: The additional typing rules for GR

is T_{\diamond} -safe. Roughly, Proposition 1 implies that each BCL expression has a safe size-bound. Note that if $f: \mathsf{N}_{\diamond} \to \mathsf{N}_{\Box \diamond}$ and $x: \mathsf{N}_{\diamond}$, then |f|(|x|) is $\mathsf{T}_{\Box \diamond}$ -chary, but not $\mathsf{T}_{\Box \diamond}$ -strict.

Strictness and chariness are syntactic notions, whereas safety is a semantic notion because of the use of $=_{nwf}$ in Definition 30(c). Thus:

Lemma 31. If $\Sigma \vdash p$: **b** and p is **b**-strict or **b**-chary, then p is also **b**-safe.

Proof. Since $\underline{0}$ is both **b**-strict and **b**-chary and since $p =_{nwf} p \lor \underline{0} =_{nwf} \underline{0} \lor p =_{nwf} p + \underline{0} =_{nwf} \underline{0} + p$, the lemma follows.

The next lemma notes a key property of safe second-order polynomials.

Lemma 32 (Safe substitution). Fix Σ . Given a γ -safe p_0 , a σ -safe p_1 , and a variable x with $\Sigma(x) = \sigma$, we can effectively find a γ -safe p'_0 such that $p_0[x := p_1] \leq_{\text{nwf}} p'_0$.

Proof. Except for the case when p_1 is a λ -expression, the argument is a straightforward induction. When p_1 is a λ -expression, the substitution can trigger a cascade of other substitutions to deal with. However, as we are working with an applied simply-typed λ -calculus, strong normalization holds [Win93], and hence, these cascades are finite. Consequently, to deal with this case we simply use a stronger induction than before, say on the syntactic structure of p_0 and p_1 and on the length of the longest path of β -reductions to normal form of $p_0[x := p_1]$. This is fairly conventional and left to the reader.¹⁸

Remark 25 informally argued that if e, an ATR expression, does not involve impredicative-type variables, then |e| has a well-founded upper bound. The analogous argument here would be that if e does not involve flat-type variables, then |e| has a safe upper bound. This assertion is true, but not so interesting because most natural crec-expressions have their recursor variable of flat type. To get around this problem we introduce a little formalism, GR (for growth rate) which includes a simple iteration construct that does not depend so heavily on flat-type variables and which captures ATR's growth rate properties *including* ATR's difficulty with flat-type values. We show in Theorem 34 that GR expressions that do not involve flat-type variables have safe upper bounds.

Definition 33. GR's raw terms are given by: $S ::= \mathbf{0}^* | (\lor S S) | (\mathsf{s} S) | (\mathsf{R} S S S) | X | (S S) | (\lambda X.S)$. The typing rules for GR consist of \rightarrow -*I* and \rightarrow -*E* from Figure 2; *Zero-I*, *Const-I*, *Subsumption*, *Shift*, and \lor -*I* from Figure 14; and s -*I* and R -*I* from Figure 16.¹⁹ The intended interpretations of \lor , s , and R are: $(\lor \underline{m} \underline{n}) = \max(\underline{m}, \underline{n}), (\mathsf{s} \underline{m}) = \underline{m} + 1$, and $(\mathsf{R} f \underline{m} \underline{n}) = f^{(m)}(\underline{n})$.

We straightforwardly extend the \mathcal{L}_{nwf} -semantics for second-order polynomials to GR. Note: $\mathcal{L}_{nwf} \llbracket \lambda m, n \cdot (\mathsf{R} \ f \ m \ n) \rrbracket \rho = \lambda m, n \in \omega \cdot n2^m$ when $\rho(f) = \lambda k \in \omega \cdot 2k$. So GR has familiar problems with flat-type values. We note that the GR analogues of Lemmas 19, 20, and 22 all hold. *Terminology:* $\Sigma \vdash s:\sigma$ is *flat-type-variable free* when no variable is explicitly or implicitly assigned a flat type by the judgment.

 $^{^{18}}$ Alternatively, the lemma's proof could be done through a logical relations induction [Win93].

¹⁹Recall from §5 that $succ(\ell)$ = the successor of ℓ in the ordering on labels.

Theorem 34. Given a flat-type-variable free $\Sigma \vdash s: \gamma$, we can effectively find a γ -safe p_s with respect to Σ such that $s \leq_{nwf} p_s$. Moreover, we can choose p_s so that all free variable occurrences are unshadowed.

Proof. Without loss of generality we assume that s is in β -normal form. The argument is a structural induction on the derivation of $\Sigma \vdash s: \gamma$. We consider the cases of the last rule used in the derivation. Let d range over ω .

CASE: Zero-I. Then $s = \underline{0}$ and $\gamma = \mathsf{T}_{\varepsilon}$. So $p_s = \underline{0}$ suffices since $\underline{0}$ is T_{ε} -strict.

CASE: Const-I. Then $s = \underline{k}$ and $\gamma = \mathsf{T}_{\diamond}$. So $p_s = \underline{k}$ suffices since \underline{k} is T_{\diamond} -strict.

CASE: Id-I. Then s = x, a variable. SUBCASE: γ is a base type. Then $p_s = x$ suffices since x is γ -chary. SUBCASE: $\gamma = (\sigma_0, \ldots, \sigma_k) \rightarrow \mathbf{b}$. (Recall the introduction of \mathbf{p}_x in Definition 28.) Let $\Sigma' = \Sigma, x_0: \sigma_0, \ldots, x_k: \sigma_k$ and $p' = (\mathbf{p}_x \ p_0 \ \ldots \ p_k)$ where, for each i, $p_i = x_i$ if $tail(\sigma_i) \leq \mathbf{b}$, and $p_i = \underline{0}_{\sigma_i}$, otherwise. (Note that since s is flat-type-variable free, $tail(\sigma_i) \neq \mathbf{b}$ for each i.) Then p' is **b**-chary with respect to Σ' and $(x \ x_0 \ \ldots \ x_k) \leq_{nwf} p'$. It follows that $p_s = \lambda x_0, \ldots, x_k \cdot p'$ suffices.

CASE: \rightarrow -I. This case follows by the induction hypothesis and clause (iii) in Definition 30(c).

CASE: \rightarrow -E. This case follows by the induction hypothesis and Lemma 32.

CASE: Subsumption. Then by Subsumption we know that $\Sigma \vdash s: \gamma'$ where $\gamma' \leq : \gamma$. Without loss of generality, we assume $\gamma' \leq : \gamma$. By the induction hypothesis there exists p, a γ' -safe size-bound for s with respect to Σ . It follows from Definition 30 that p is γ -strict with respect to Σ . Hence, $p_s = p$ suffices.

CASE: Shift. Recall that if $(\vec{\sigma}) \to \mathbf{b} \propto (\vec{\sigma}') \to \mathbf{b}'$, then, for each *i*, $tail(\sigma_i) \leq \mathbf{b}$ implies $tail(\sigma_i') \leq \mathbf{b}'$ and $tail(\sigma_i) = \mathbf{b}$ implies $tail(\sigma_i') = \mathbf{b}'$. Thus this case follows from Lemma 20 (in both its second-order polynomial and GR versions) and Definition 30.

CASE: s-I. Then $s = (s \ s_1)$ and $\gamma = \mathsf{T}_{\diamond_d}$. So by s-I, we know that $\Sigma \vdash s_1: \mathsf{T}_{\diamond_d}$ and by the induction hypothesis we have that there is a T_{\diamond_d} -strict q and a T_{\diamond_d} -chary r with $s_1 \leq_{\text{nwf}} q + r$. Thus $p_s = (q + 1) + r$ suffices since q + 1 is T_{\diamond_d} -strict.

CASE: $\lor I$. Then $s = (\lor s_0 \ s_1)$. SUBCASE: $\gamma = \mathsf{T}_{\diamond_d}$. So by $\lor I$ we know that $\Sigma_0 \vdash s_0$: T_{\diamond_d} and $\Sigma_1 \vdash s_1$: T_{\diamond_d} , where $\Sigma = \Sigma_0 \cup \Sigma_1$. By the induction hypothesis, there are T_{\diamond_d} -strict q_0 and q_1 and T_{\diamond_d} -chary r_0 and r_1 such that $s_0 \leq_{\mathrm{nwf}} q_0 + r_0$ and $s_1 \leq_{\mathrm{nwf}} q_1 + r_1$. Thus $p_s = (q_0 \lor q_1) + (r_0 \lor r_1)$ suffices since $s \leq_{\mathrm{nwf}} (q_0 + r_0) \lor (q_1 + r_1) \leq_{\mathrm{nwf}} (q_0 \lor q_1) + (r_0 \lor r_1)$ and since $(q_0 \lor q_1)$ is T_{\diamond_d} -strict $(r_0 \lor r_1)$ is T_{\diamond_d} -chary with respect to Σ . SUBCASE: $\gamma = \mathsf{T}_{\Box_d}$. This follows by an easy modification of the above argument.

CASE: R-*I*. Then $s = (R \ s_0 \ s_1 \ s_2)$. SUBCASE: $\gamma = \mathsf{T}_{\diamond d}$. So by R-*I*, we have $\Sigma_0 \vdash s_0: \mathsf{T}_{\diamond d} \to \mathsf{T}_{\diamond d}, \Sigma_1 \vdash s_1: \mathsf{T}_{\Box d}, \text{ and } \Sigma_2 \vdash s_2: \mathsf{T}_{\diamond d}, \text{ where } \Sigma = \Sigma_0 \cup \Sigma_1 \cup \Sigma_2$. Since *s* is flat-type-variable free, we must have $s_1 = \lambda z \cdot s'_1$ where $\Sigma_0, z: \mathsf{T}_{\diamond d} \vdash s'_1: \mathsf{T}_{\diamond d}$. Hence, by the induction hypothesis, there are $\mathsf{T}_{\diamond d}$ -strict q_0 and $q_2, \mathsf{T}_{\diamond d}$ -chary r_0 and r_2 , and $\mathsf{T}_{\Box d}$ -safe p_1 such that $s'_1 \leq_{\mathsf{nwf}} q_0 + r_0, s_1 \leq_{\mathsf{nwf}} p_1, \text{ and } s_2 \leq_{\mathsf{nwf}} q_2 + r_2$. Note that p_1 is also $\mathsf{T}_{\diamond d}$ -strict. Suppose z has no free occurrences in $q_0 + r_0$. Then it follows that $s \leq_{\mathsf{nwf}} (q_0 + r_0) \lor (q_2 + r_2) \leq_{\mathsf{nwf}} (q_0 \lor q_2) + (r_0 \lor r_2); \text{ so } p_s = (q_0 \lor q_2) + (r_0 \lor r_2)$ suffices. Now suppose z does have a free occurrence in $q_0 + r_0$. Since q_0 is $\mathsf{T}_{\diamond d}$ -strict, z cannot occur in q_0 . Since s is flat-type-variable free, it follows that $r_0 =_{\mathsf{nwf}} z \lor r'_0$ where z has no free occurrences in r'_0 and $\mathsf{where } r'_0$ is $\mathsf{T}_{\diamond d}$ -chary. By the inequality $q + (q' + r') \lor r \leq (q + q') + r' \lor r$, it follows that $s \leq_{\mathsf{nwf}} (p_1 * q_0 + q_2) + (r'_0 \lor r_2)$. So, $p_s = (p_1 * q_0 + q_2) + (r'_0 \lor r_2)$ suffices. (Note the parallel to the proof of Proposition 1.) SUBCASE: $\gamma = \mathsf{T}_{\Box d}$. This follows by an easy modification of the above argument.

$$\mathbf{p}\text{-}I: \quad \frac{\Sigma \vdash p:\sigma}{\Sigma \vdash (\mathbf{p} \, p):\sigma} \qquad \qquad \mathbf{q}\text{-}I: \quad \frac{\Sigma \vdash p:\sigma}{\Sigma \vdash (\mathbf{q} \, p):\sigma^{\dagger}} \qquad \qquad \mathbf{r}\text{-}I: \quad \frac{\Sigma \vdash p:\sigma}{\Sigma \vdash (\mathbf{r} \, p):\sigma^{\dagger}}$$

Figure 17: Typing rules for the \mathbf{p} , \mathbf{q} , and \mathbf{r} combinators

9. FLAT TYPES AND WELL-TEMPEREDNESS

To avoid problems like the one of Example 23(b), flat-type values need to be restricted. The GR formalism of the previous section is subject to roughly the same problem as that of Example 23(b), but by Theorem 34 flat-type-variable free GR expressions have safe second-order polynomial bounds. This suggests that a solution to the flat-type values problem is to require all flat-type values to have safe size-bounds. We call this property *well-temperedness*, meaning: all things are in the right proportions.

Definition 35. A $t \in \mathcal{L}_{nwf}[\![\gamma]\!]$ is γ -well-tempered when γ is strict or when γ is flat and there is a closed, γ -safe s with $t \leq \mathcal{L}_{nwf}[\![s]\!]$.

Lemma 36. Suppose $\Sigma \vdash p:\sigma$, $\rho \in \mathcal{L}_{nwf}[\![\Sigma]\!]$, and $\rho(x)$ is $\Sigma(x)$ -well-tempered for each $x \in \text{preimage}(\Sigma)$. Then $\mathcal{L}_{nwf}[\![p]\!] \rho$ is σ -well-tempered.

Lemma 36's proof is an induction on the derivation of $\Sigma \vdash p: \sigma$. Everything is fairly straightforward except that the $\rightarrow -E$ case depends critically on Lemma 32. Lemma 36 indicates that a semantics for the second-order polynomials based on well-temperedness will be well defined.

Definition 37 (The well-tempered semantics).

(a) Inductively define $\mathcal{L}_{wt}[\![\sigma]\!]$ by: $\mathcal{L}_{wt}[\![\mathsf{T}_{\ell}]\!] = \omega$ and, for $\sigma = (\sigma_1, \ldots, \sigma_k) \to \mathsf{T}_{\ell}$, $\mathcal{L}_{wt}[\![\sigma]\!]$ is the restriction to $(\mathcal{L}_{wt}[\![\sigma_1]\!], \ldots, \mathcal{L}_{wt}[\![\sigma_k]\!]) \to \mathcal{L}_{wt}[\![\mathsf{T}_{\ell}]\!]$ of the σ -well-tempered elements of $\mathcal{L}_{nwf}[\![\sigma]\!]$. $\mathcal{L}_{wt}[\![\Sigma]\!]$ and $\mathcal{L}_{wt}[\![\Sigma \vdash p:\sigma]\!]$ are defined in the standard way.

(b) Inductively define $\mathcal{V}_{wt}\llbracket\sigma\rrbracket$ by: $\mathcal{V}_{wt}\llbracket\mathsf{N}_{\ell}\rrbracket = \mathbb{N}$ and, for $\sigma = (\sigma_1, \ldots, \sigma_k) \to \mathsf{N}_{\ell}$, $\mathcal{V}_{wt}\llbracket\sigma\rrbracket$ is the restriction to $(\mathcal{V}_{wt}\llbracket\sigma_1\rrbracket, \ldots, \mathcal{V}_{wt}\llbracket\sigma_k\rrbracket) \to \mathcal{V}_{wt}\llbracket\mathsf{N}_{\ell}\rrbracket$ of the $f \in \mathcal{V}_{nwf}\llbracket\sigma\rrbracket$ with $|f| \in \mathcal{L}_{wt}\llbracket|\sigma|\rrbracket$. $\mathcal{V}_{wt}\llbracket\Gamma; \Delta\rrbracket$ and $\mathcal{V}_{wt}\llbracket\Gamma; \Delta \vdash E; \sigma\rrbracket$ are defined in the standard way.

(c) We write $p =_{\text{wt}} p'$ when $\mathcal{L}_{\text{wt}}[\![\Sigma \vdash p:\sigma]\!] |\rho| = \mathcal{L}_{\text{wt}}[\![\Sigma \vdash p':\sigma]\!] |\rho|$ for all $|\rho| \in \mathcal{L}_{\text{wt}}[\![\Sigma]\!]$. We define $\leq_{\text{wt}}, \geq_{\text{wt}}, \ldots$ analogously.

There is still a problem with flat-type values. To give closed-form upper bounds on recursions, we sometimes need to decompose a safe flat-type polynomial into strict and chary parts. (Recall that safety is a semantic, not syntactic, notion.) For flat-type-variable free safe polynomials this is easy. A way of breaking flat-type variables into strict and chary parts would allow us to extend this decomposition to all safe polynomials. We introduce two new combinators to effect such a decomposition. Since there is no canonical way to do this decomposition, we take a different (and trickier) approach from that of Definition 28. Terminology: Let $(\mathbf{b})^{\dagger} = (\mathbf{b})^{\ddagger} = \mathbf{b}$, $((\vec{\sigma}) \to \mathbf{b})^{\dagger} = (\vec{\sigma}') \to \mathbf{b}$, and $((\vec{\sigma}) \to \mathbf{b})^{\ddagger} = (\vec{\sigma}'') \to \mathbf{b}$, where $\vec{\sigma}' =$ the subsequence of σ_i 's in $\vec{\sigma}$ with $tail(\sigma_i) \neq \mathbf{b}$ and $\vec{\sigma}'' =$ the subsequence of σ_i 's in $\vec{\sigma}$ with $tail(\sigma_i) \neq \mathbf{b}$. (Recall: () $\to \mathbf{b} \equiv \mathbf{b}$.)

Definition 38. We add two new combinators, \mathbf{q} and \mathbf{r} , to the second-order polynomials with typing rules given in Figure 17. Suppose $\Sigma = w_1: \tau_1, \ldots, w_n: \tau_n, \quad \Sigma \vdash p: \gamma$, and $\rho \in \mathcal{L}_{wt}[\![\Sigma]\!]$. For γ strict, define $\mathcal{L}_{wt}[\![(\mathbf{q}\,p)]\!] \rho = \underline{0}_{\gamma^{\dagger}}$ and $\mathcal{L}_{wt}[\![(\mathbf{r}\,p)]\!] = \mathcal{L}_{wt}[\![(\mathbf{p}\,p)]\!]$. Suppose $\gamma = (\sigma_0, \ldots, \sigma_k) \to \mathbf{b}$ is flat. Let $\zeta = (\vec{\tau}, \vec{\sigma}) \to \mathbf{b}, (\sigma'_{i_1}, \ldots, \sigma'_{i_m}) \to \mathbf{b} = \gamma^{\dagger}, (\sigma''_{j_1}, \ldots, \sigma''_{j_n}) \to$ $\mathbf{b} = \gamma^{\ddagger}, \ \vec{x} = x_0, \dots, x_k, \ \vec{x}' = x_{i_1}, \dots, x_{i_m}, \ \vec{x}'' = x_{j_1}, \dots, x_{j_n}, \ \text{and} \ \{z_1, \dots, z_u\} = \{w_i \mid \tau_i = \mathbf{b}\} \cup \{x_i \mid \tau_i = \mathbf{b}\} \text{ where the } z_i \text{'s are all distinct. Define } \mathcal{L}_{wt}[\![(\mathbf{q} p)]\!] = \mathcal{L}_{wt}[\![\lambda \vec{x}' \cdot q]\!] \text{ and} \ \mathcal{L}_{wt}[\![(\mathbf{r} p)]\!] = \mathcal{L}_{wt}[\![\lambda \vec{x}'' \cdot r]\!], \ \text{where (i) } q \text{ is } \mathbf{b}\text{-strict with respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(ii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \mathbf{b}\text{-chary with} \ \text{respect to } \Sigma, \vec{x} : \vec{\sigma}, \ \text{(iii) } r \text{ is } \vec{\sigma} \text{ is$

$$\mathcal{L}_{\mathrm{wt}}\llbracket \vdash \lambda \vec{w} \cdot (\mathbf{p} \ p) : \zeta \rrbracket \leq \begin{cases} \mathcal{L}_{\mathrm{wt}}\llbracket \vdash \lambda \vec{w}, \vec{x} \cdot (q + r \lor z_1 \lor \ldots \lor z_u) : \zeta \rrbracket, & \text{for computational } \gamma; \\ \mathcal{L}_{\mathrm{wt}}\llbracket \vdash \lambda \vec{w}, \vec{x} \cdot (q \lor r \lor z_1 \lor \ldots \lor z_u) : \zeta \rrbracket, & \text{for oracular } \gamma. \end{cases}$$

For each variable x, we abbreviate $(\mathbf{q} x)$ by \mathbf{q}_x and $(\mathbf{r} x)$ by \mathbf{r}_x . Also, we take $(\mathbf{q}_x \vec{x}')$ as being **b**-strict and $(\mathbf{r}_x \vec{x}'')$ as being **b**-chary.

Example 39. By the definition of *prn* given in Figure 13 and our proof sketch for Proposition 1, it follows that $|prn| \leq_{\text{wt}} \lambda |e|, |x| \cdot ((|x|+1) * \mathbf{q}_{|e|}(|x|) + \mathbf{r}_{|e|}(|x|))$.

By Definitions 35 and 37, q and r as in Definition 38 must exist. By the axiom of choice, there are functions that pick out particular q and r. **N.B.** The choices of q and r are **arbitrary** subject to satisfying conditions (i), (ii), and (iii) of Definition 38. The semantics for the second-order polynomials is thus parameterized by the functions that pick out the required q's and r's. The choices \mathbf{q} and \mathbf{r} make are analogous to the choices of a and $b \in \omega$ in the situation were one knows that $f \in O(n)$ and picks some arbitrary a and b such that $f(n) \leq a \cdot n + b$ for all n. Such a and b can be used in constructing algebraic upper bounds on expressions involving f. If later we determine concrete a_0 and b_0 such that $f(n) \leq a_0 \cdot n + b_0$ for all n, then said algebraic upper bounds are still valid after the substitution $[a := a_0, b := b_0]$ since the choices of a and b were arbitrary.

Definition 40. Suppose $\Sigma \vdash p: \gamma$, where $\{y_1, \ldots, y_k\} = \{y \mid \Sigma(y) = tail(\gamma)\}$. We say that p is manifestly γ -safe with respect to Σ if and only if the only applications of the \mathbf{p} , \mathbf{q} , and \mathbf{r} combinators are to variables, and:

(a) when $\gamma = \mathsf{T}_{\Box_d}$, then p is of one of the forms: $q, r \lor y_{i_1} \lor \ldots \lor y_{i_n}$, and $q \lor r \lor y_{i_1} \lor \ldots \lor y_{i_n}$, where q is γ -strict, r is γ -chary with no occurrences of any of the y_i 's, and $\{y_{i_1}, \ldots, y_{i_n}\}$ is a (possibly empty) subset of $\{y_1, \ldots, y_k\}$;

(b) when $\gamma = \mathsf{T}_{\diamond_d}$, then p is of one of the forms: $q, r \lor y_{i_1} \lor \ldots \lor y_{i_n}$, and $q + r \lor y_{i_1} \lor \ldots \lor y_{i_n}$, where $q, r, \text{ and } \{y_{i_1}, \ldots, y_{i_n}\}$ are as in (a); and

(c) when $\gamma = (\sigma_0, \ldots, \sigma_m) \to \mathbf{b}$, then the β -normal form of $(p \ \vec{x})$ is manifestly **b**-safe with respect to $\Sigma, x_0; \sigma_0, \ldots, x_m; \sigma_m$.

Lemma 41 (Manifestly safe substitution). Fix Σ . Given a manifestly γ -safe p_0 , a manifestly σ -safe p_1 , and a variable x with $\Sigma(x) = \sigma$, we can effectively find a manifestly γ -safe p'_0 such that $p_0[x := p_1] \leq_{\text{nwf}} p'_0$.

Proof. This is a straightforward adaptation of the proof of Lemma 32.

We now have a reasonable semantics for ATR and the tools to work with this semantics to establish (in Theorem 43) a *safe polynomial boundedness* result for ATR, where:

Definition 42. Suppose $\Gamma; \Delta \vdash e: \sigma$. We say that p is a $|\sigma|$ -safe polynomial size-bound for e with respect to $\Gamma; \Delta$ when p is a $|\sigma|$ -safe second-order polynomial with respect to $|\Gamma; \Delta|$ and $|\mathcal{V}_{wt}[\![e]\!] \rho| \leq \mathcal{L}_{wt}[\![p]\!] |\rho|$ for all $\rho \in \mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$; if in addition p is manifestly $|\sigma|$ -safe with respect to $\Gamma; \Delta$, we say that p is a manifestly $|\sigma|$ -safe polynomial size-bound for e with respect to $\Gamma; \Delta$. (The "with respect to" clause is dropped when it is clear from context.)

N. DANNER AND J. S. ROYER

10. Polynomial size-boundedness

Theorem 43 (Polynomial Boundedness). Given $\Gamma; \Delta \vdash e; \gamma$, we can effectively find p_e , a manifestly $|\gamma|$ -safe polynomial size-bound for e with respect to $\Gamma; \Delta$.

Proof. The argument is a structural induction on the derivation of $\Gamma; \Delta \vdash e; \gamma$. We consider the cases of the last rule used in the derivation. Excluding the **crec** case, everything is fairly straightforward. Fix $\rho \in \mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$. Note that $|\mathcal{V}_{wt}[\![\cdot]\!]\rho|$ is invariant under β - and η -equivalence. So without loss of generality, we assume that e is in β -normal form.

CASES: Int-Id-I and Aff-Id-I. Then e = x, a variable. Subcase: γ is strict. Then $p_e = |x|$ clearly suffices. Subcase: γ is flat. Hence, $level(\gamma) = 1$. Let $(\mathbf{b}_0, \ldots, \mathbf{b}_k) \to \mathbf{b} = \gamma$. Then by Definitions 37 and 38,

$$p_e = \lambda |x_0|, \dots, |x_k| \cdot \left((\mathbf{q}_{|x|} |\overrightarrow{x'|}) \odot (\mathbf{r}_{|x|} |\overrightarrow{x'|}) \lor |y_1| \lor \dots \lor |y_\ell| \right)$$

suffices, where $|\overline{x'}|$ = the subsequence of the $|x_i|$'s with $\mathbf{b}_i \neq \mathbf{b}$ and $\{y_1, \ldots, y_\ell\} = \{y \mid (\Gamma, x_0; \mathbf{b}_0, \ldots, x_k; \mathbf{b}_k; \Delta)(y) = \mathbf{b}\}$, and where $\odot = +$, if γ is computational, and $\odot = \vee$, if γ is oracular.

CASE: Zero-I. Then $e = \epsilon$ and $\gamma = N_{\varepsilon}$. Clearly $p_e = 0$ suffices.

CASE: Const-I. Then e = some constant k and $\gamma = N_{\diamond}$. Clearly $p_e = |k|$ suffices.

CASE: $t_{\mathbf{a}}$ -I. So $\gamma = N_{\diamond_d}$ for some d. Clearly $p_e = \underline{1}$ suffices.

CASE: d-I. Then e = (d e') for some e' and $\gamma = N_{\diamond_d}$ for some d. By the induction hypothesis, there is $p_{e'}$, a manifestly T_{\diamond_d} -safe polynomial size-bound for e' with respect to $\Gamma; \Delta$. Clearly $p_e = p_{e'}$ suffices.

CASE: down-*I*. Then $e = (\text{down } e_0 \ e_1)$ with $\Gamma; \Delta \vdash e_0: \mathbf{b}_0, \ \Gamma; \Delta \vdash e_1: \mathbf{b}_1$, and $\gamma = \mathbf{b}_1$. By the induction hypothesis, there is a p_{e_1} , a manifestly $|\mathbf{b}_1|$ -safe polynomial size-bound for e_1 with respect to $\Gamma; \Delta$. Clearly $p_e = p_{e_1}$ suffices.

CASE: c_a -*I*. Then $e = (c_a e')$ for some e' and $\gamma = N_{\diamond_d}$ for some *d*. By the induction hypothesis, there is $p_{e'}$, a manifestly T_{\diamond_d} -safe polynomial size-bound for e' with respect to $\Gamma; \Delta$. Clearly $p_e = \underline{1} + p_{e'}$ suffices.

CASES: Subsumption and Shift. These follow as in the proof of Theorem 34.

Aside: For the arguments for the \rightarrow -*I* and \rightarrow -*E* cases below, recall from §2.10 that (2.2) and (2.3) provide the definition of length for elements of **TC** of type-level 1 and type-level 2, respectively, and that higher-type lengths are pointwise monotone nondecreasing.

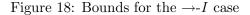
CASE: $\rightarrow -I$. Then $e = \lambda x \cdot e'$ and $\gamma = \sigma \rightarrow \tau$. By our induction hypothesis, there is a $p_{e'}$, a manifestly $|\tau|$ -safe polynomial size bound for e' with respect to $\Gamma, x:\sigma; \Delta$. Let $p_e = \lambda |x| \cdot p_{e'}$. By Definition 30(c), p_e is manifestly $|\gamma|$ -safe with respect to $|\Gamma; \Delta|$. Let vrange over $\mathcal{V}_{wt}[\![\sigma]\!]$. Then, for each $t \in \mathcal{L}_{wt}[\![\sigma]\!]$, we have the chain of bounds of Figure 18. Clearly this p_e suffices.

CASE: $\rightarrow -E$. Then $e = (e_0 \ e_1)$ and for some σ we have that $\Gamma; \Delta \vdash e_0: \sigma \rightarrow \gamma$ and $\Gamma; _\vdash e_1: \sigma$. By the induction hypothesis, there are p_{e_0} and p_{e_1} such that p_{e_0} is a manifestly $(|\sigma| \rightarrow |\tau|)$ -safe polynomial size bound for e_0 and p_{e_1} is a manifestly $|\sigma|$ -safe polynomial size-bound for e_1 . By Lemma 41 we can effectively find a manifestly γ -safe p_e such that $(p_{e_0} \ p_{e_1}) \leq_{\text{wt}} p_e$. Then we have the chain of bounds of Figure 19. Clearly this p_e suffices.

CASE: If-I. Then $e = (\text{if } e_0 \text{ then } e_1 \text{ else } e_2)$. By the induction hypothesis, there are p_{e_1} and p_{e_2} , manifestly $|\gamma|$ -safe polynomial size-bounds for e_1 and e_2 respectively. Clearly $p_e = p_{e_1} \vee p_{e_2}$ suffices.

We have just one case left, but now the real work starts.

$\left \mathcal{V}_{\mathrm{wt}}\llbracket\lambda x . e'\rrbracket\rho\right (t)$	
$= \max\{ \left (\mathcal{V}_{\mathrm{wt}} \llbracket \lambda x \cdot e' \rrbracket \rho)(v) \right \colon v \le t \}$	(by (2.2) and (2.3))
$= \max\{ \left \left(\mathcal{V}_{\mathrm{wt}} \llbracket e' \rrbracket \left(\rho \cup \{ x \mapsto v \} \right) \right \mid v \le t \} \right.$	(by the \mathcal{V}_{wt} -interpretation of λ -terms)
$\leq \max\{\left(\mathcal{L}_{\mathrm{wt}}\llbracket p_{e'} \rrbracket(\rho \cup \{ x \mapsto v \}) \mid v \leq t\}$	(by the choice of $p_{e'}$)
$\leq \left(\mathcal{L}_{\mathrm{wt}}\llbracket p_{e'} \rrbracket \left(\rho \cup \{ x \mapsto t \} \right) \right)$	(by monotonicity)
$= (\mathcal{L}_{\mathrm{wt}}\llbracket \lambda \mathbf{x} . p_{e'} \rrbracket \rho)(t)$	(by the \mathcal{L}_{wt} -interpretation of λ -terms)
$= (\mathcal{L}_{\mathrm{wt}}\llbracket p_e \rrbracket \rho)(t)$	(by the choice of p_e).



$\left \mathcal{V}_{\mathrm{wt}}\llbracket(e_0 \ e_1) ight ceil ho ight $	
$= \left \left(\mathcal{V}_{\mathrm{wt}} \llbracket e_0 \rrbracket \rho \right) \left(\mathcal{V}_{\mathrm{wt}} \llbracket e_1 \rrbracket \rho \right) \right $	(by the \mathcal{V}_{wt} -interpretation of application)
$\leq \left(\mathcal{V}_{\mathrm{wt}}[\![e_0]\!]\rho \right) \left(\mathcal{V}_{\mathrm{wt}}[\![e_1]\!]\rho \right)$	(by (2.2) and (2.3))
$\leq \left(\mathcal{L}_{\mathrm{wt}}\llbracket p_{e_0} \rrbracket \rho \right) \left(\mathcal{L}_{\mathrm{wt}}\llbracket p_{e_1} \rrbracket \rho \right)$	(by monotonicity and the choices of p_{e_0} and p_{e_1})
$= \mathcal{L}_{\mathrm{wt}}\llbracket (p_{e_0} p_{e_1}) \rrbracket \left \rho \right $	(by the \mathcal{L}_{wt} -interpretation of application)
$= \hspace{0.1 cm} \mathcal{L}_{\mathrm{wt}} \llbracket p_{e} \rrbracket \ket{\rho}$	(by the choice of p_e).

Figure 19: Bounds for the \rightarrow -*E* case

CASE: crec-*I*. Then $\gamma = (\mathbf{b}_1, \ldots, \mathbf{b}_k) \to \mathbf{b}_0 \in \mathcal{R}$, so $\mathbf{b}_1 = \mathsf{N}_{\Box_{d_1}}$ for some d_1 , and $e = (\operatorname{crec} a \ (\lambda_r f \cdot A))$ with $a \in \mathbf{0}^*$, $\Gamma; f \colon \gamma \vdash A \colon \gamma$, and TailPos(f, A). (Recall: *TailPos* is defined in Figure 11.) For simplicity we assume $\{\mathbf{b}_1, \ldots, \mathbf{b}_k\} = \{\mathsf{N}_{\varepsilon}, \ldots, \mathsf{N}_{\Box_{d_1}}, \mathsf{N}_{\Box_{d_1}}\} \cup \{\mathbf{b} \upharpoonright \mathsf{N}_{\Box_{d_1}} \lneq : \mathbf{b} \leq : \mathbf{b}_{\max}\}$ for some \mathbf{b}_{\max} . Without loss of generality we suppose:

$$A = \lambda x_1, \dots, x_k \, \mathbf{B}, \tag{10.1}$$

where $\widehat{\Gamma}$; $f: \gamma \vdash B$: \mathbf{b}_0 for $\widehat{\Gamma} = \Gamma, x_1: \mathbf{b}_1, \dots, x_k: \mathbf{b}_k, B$ is in β -normal form, and TailPos(f, B).

Aside: To find p_e for this case, we analyze e's tail recursion and determine size bounds on how large the tail-recursion's arguments can grow. In particular, we show that there is a polynomial bound beyond which the first argument *cannot* grow; hence, by (4.2), this polynomial bounds the depth of e's tail recursion. From this bound on recursion depth and from the size bounds on the tail-recursion arguments, constructing p_e is straightforward. To derive these bounds, we proceed a little informally and work with unfolded versions of e.

Consider the occurrences of f in B. Since we have TailPos(f, B) and $\widehat{\Gamma}; f: \gamma \vdash B: \mathbf{b}_0$, these occurrences must have enclosing expressions of the form $(f \ e_1 \ \dots \ e_k)$, where $\widehat{\Gamma}; _\vdash e_1: \mathbf{b}_1, \dots, \widehat{\Gamma}; _\vdash e_k: \mathbf{b}_k$. For a given such subexpression of B, we know by the induction hypothesis that, for each $i = 1, \dots, k$, there is a p_i , a manifestly \mathbf{b}_i -safe polynomial sizebound for e_i with respect to $\widehat{\Gamma}; _$. Since f occurs but finitely many times in B, we may choose p_1, \dots, p_k so that they bound the size of the corresponding argument expressions for every f-application in B. Without loss of generality, we assume that if $\mathbf{b}_i = \mathbf{b}_j$, then $p_i = p_j$. Using the crec reduction rule (4.2), we expand out one-level of e's crec-recursion and, by using β - and η -reductions, clean things up to obtain

$$e^{(1)} = \lambda \vec{x}$$
 if $|a| \leq |x_1|$ then B else ϵ , where

$$B = \text{the } \beta\eta\text{-normal form of } B[f := (\text{crec } (\mathbf{0} \oplus a) (\lambda_r f \boldsymbol{.} A))]).$$

Clearly, $\mathcal{V}_{wt}[\![e]\!] = \mathcal{V}_{wt}[\![e^{(1)}]\!]$. Let ξ denote the substitution $[|x_1| := p_1, \ldots, |x_k| := p_k]$. From our choices of the p_i 's and \widehat{B} it follows that $(p_1 \xi), \ldots, (p_k \xi)$ bound the size of the corresponding argument expressions for every *f*-application in \widehat{B} . For each *i*, $(p_i \xi)$ can be equivalently expressed in terms of p_i as follows. *Terminology:* An *r* is *strictly* **b**-*chary* when *r* is **b**-chary and contains no occurrences of type-**b** variables.²⁰

(*Note:* In working through the proofs of Lemmas 44 and 45 below, the reader many want to consider the case of: $\gamma = (\mathsf{N}_{\Box_1}, \mathsf{N}_{\Box_0}, \mathsf{N}_{\diamond_1}) \to \mathsf{N}_{\diamond_1}$, f has but one occurrence in $A, p_1(|x_1|, |x_2|, |x_3|) = |g|(|x_2|) \lor |x_1|, p_2(|x_1|, |x_2|, |x_3|) = |x_2|$, and $p_3(|x_1|, |x_2|, |x_3|) = q_3(|x_1|, |x_2|) + |x_3|$, where $g: \mathsf{N}_{\diamond_0} \to \mathsf{N}_{\Box_1}$ and where q_3 is an ordinary polynomial.)

Lemma 44 (The one step lemma). Each p_i can be taken so that:

(a) If $\mathbf{b}_i \leq : \mathsf{N}_{\Box_{d_1}}$, then $p_i \xi =_{\mathrm{wt}} p_i$.

(b) If $\mathsf{N}_{\Box_{d_1}} \leq \mathbf{b}_i = \mathsf{N}_{\Box_d}$, then there is a \mathbf{b}_i -strict q_i and a strictly \mathbf{b}_i -chary r_i such that $p_i \xi =_{\mathrm{wt}} q_i \xi \lor r_i \xi \lor p_i$.

(c) If $\mathsf{N}_{\Box_{d_1}} \leq :\mathbf{b}_i = \mathsf{N}_{\diamond_d}$, then there is a \mathbf{b}_i -strict q_i and a strictly \mathbf{b}_i -chary r_i such that $p_i \xi =_{\mathrm{wt}} q_i \xi + r_i \xi \lor p_i$.

Proof. For each d, let:

$$\{ u_0^d, \dots, u_{b_d}^d \} \stackrel{\text{def}}{=} \{ u \mid \Gamma(u) = \mathsf{N}_{\Box_d} \}. \qquad \{ w_0^d, \dots, w_{b_d'}^d \} \stackrel{\text{def}}{=} \{ u \mid \Gamma(u) = \mathsf{N}_{\diamond_d} \}.$$

$$\{ \overline{u}_0^d, \dots, \overline{u}_{c_d}^d \} \stackrel{\text{def}}{=} \{ x_i \mid \mathbf{b}_i = \mathsf{N}_{\Box_d} \}. \qquad \{ \overline{w}_0^d, \dots, \overline{w}_{c_d'}^d \} \stackrel{\text{def}}{=} \{ x_i \mid \mathbf{b}_i = \mathsf{N}_{\diamond_d} \}.$$

(The \overline{u} 's and \overline{w} 's correspond to the arguments of the recursion while the *u*'s and *w*'s correspond to the other parameters.)

For part (a), we inductively consider the cases of $\mathbf{b}_i = \mathsf{N}_{\varepsilon}, \ \mathsf{N}_{\Box_1}, \ldots, \mathsf{N}_{\Box_d}$ in turn.

CASE: $\mathbf{b}_i = \mathbf{N}_{\varepsilon}$. By the induction hypothesis, we may take p_i to be $q \lor r \lor \hat{t} \lor \bar{t}$, where q is T_{\Box_0} -strict, r is strictly T_{\Box_0} -chary, $\hat{t} = \bigvee_{a=0}^{b_0} |u_a^0|$, and $\bar{t} = \bigvee_{a=0}^{c_0} |\overline{u}_a^0|$. It follows from the size typing rules that the only T_{\Box_0} -strict terms are $=_{\mathrm{wt}} \underline{0}$. So, it suffices to take $p_i = r \lor \hat{t} \lor \bar{t}$. Note that $r = r\xi$ and $\hat{t} = \hat{t}\xi$ since neither r nor \hat{t} have any occurrences of any \overline{u}_a^0 . Also recall that we are assuming that if $\mathbf{b}_i = \mathbf{b}_j$, then $p_i = p_j$. Thus, for each a, $|\overline{u}_a^0|\xi = p_a = p_i$. So, $\overline{t}\xi =_{\mathrm{wt}} p_i = r \lor \hat{t} \lor \overline{t}$. Consequently,

$$p_i \xi =_{\mathrm{wt}} (r \lor \widehat{t} \lor \overline{t}) \xi =_{\mathrm{wt}} r \xi \lor \widehat{t} \xi \lor \overline{t} \xi =_{\mathrm{wt}} r \lor \widehat{t} \lor (r \lor \widehat{t} \lor \overline{t}) =_{\mathrm{wt}} r \lor \widehat{t} \lor \overline{t} =_{\mathrm{wt}} p_i.$$

Hence, our choice of p_i suffices for this case.

CASE: $\mathbf{b}_i = \mathbf{N}_{\Box_1}$. By the induction hypothesis, we can take p_i to be of the form $q \lor r \lor \hat{t} \lor \bar{t}$, where q is T_{\Box_1} -strict, r is strictly T_{\Box_1} -chary, $\hat{t} = \bigvee_{a=0}^{b_1} |u_a^1|$, and $\bar{t} = \bigvee_{a=0}^{c_1} |\overline{u}_a^1|$. We first consider q. Since Γ does not assign any of x_1, \ldots, x_k the type N_{\diamond_0} , the only variables from x_1, \ldots, x_k whose lengths can occur in q are those assigned type N_{ε} . Let $\hat{q} = q \xi$, where for each i' with $\mathbf{b}_{i'} = \mathsf{N}_{\varepsilon}$, we take $p_{i'}$ to satisfy part (a). Hence, it follows that $\hat{q}\xi =_{\mathrm{wt}} \hat{q}$. Also, by the monotonicity of everything in sight, we have that $q \leq_{\mathrm{wt}} \hat{q}$. By the same argument,

²⁰Such an r may contain occurrences of variables of types of the form $(\sigma_0, \ldots, \sigma_k) \to \mathbf{b}$.

for $\hat{r} = r \xi$ we have that $\hat{r} \xi =_{\text{wt}} \hat{r}$ and $r \leq_{\text{wt}} \hat{r}$. So, it suffices to take $p_i = \hat{q} \lor \hat{r} \lor \hat{t} \lor \bar{t}$. Note that $\hat{t} = \hat{t} \xi$ since \hat{t} has no occurrence of any \overline{u}_a^d . Also recall that we are assuming that if $\mathbf{b}_i = \mathbf{b}_j$, then $p_i = p_j$. Thus, for each a, $|\overline{u}_a^1| \xi = p_a = p_i$. So, $\overline{t} \xi =_{\text{wt}} p_i = \hat{q} \lor \hat{r} \lor \hat{t} \lor \bar{t}$. Consequently,

$$p_i \xi =_{\mathrm{wt}} (\widehat{q} \lor \widehat{r} \lor \widehat{t} \lor \overline{t}) \xi =_{\mathrm{wt}} \widehat{q} \xi \lor \widehat{r} \xi \lor \widehat{t} \xi \lor \overline{t} \xi =_{\mathrm{wt}} \widehat{q} \lor \widehat{r} \lor \widehat{t} \lor (\widehat{q} \lor \widehat{r} \lor \widehat{t} \lor \overline{t}) =_{\mathrm{wt}} \widehat{q} \lor \widehat{r} \lor \widehat{t} \lor \overline{t} =_{\mathrm{wt}} p_i.$$

Hence, our choice of p_i suffices for this case.

CASES: $\mathbf{b}_i = \mathsf{N}_{\Box_2}, \ldots, \mathsf{N}_{\Box_{d_1}}$. These cases follow from essentially the same as argument given for the $\mathbf{b}_i = \mathsf{N}_{\Box_1}$ case.

Therefore, part (a) follows.

We henceforth assume that p_i satisfies part (a) for each i with $\mathbf{b}_i \leq : \mathbb{N}_{\Box_{d_i}}$.

For parts (b) and (c), consider the cases of $\mathbf{b}_i = \mathsf{N}_{\diamond_{d_1}}, \, \mathsf{N}_{\square_{d_1+1}}, \dots, \mathbf{b}_{\max}$ in turn.

CASE: $\mathbf{b}_i = \mathsf{N}_{\diamond d_1}$. By the induction hypothesis, we may take p_i to be of the form $q + r \vee \widehat{t} \vee \overline{t}$, where q is $\mathsf{T}_{\diamond d_1}$ -strict, r is strictly $\mathsf{T}_{\diamond d_1}$ -chary, $\widehat{t} = \bigvee_{a=0}^{b'_{d_1}} |w_a^{d_1}|$, and $\overline{t} = \bigvee_{a=0}^{c'_{d_1}} |\overline{w}_a^{d_1}|$. Note that as in the previous cases, $\widehat{t} = \widehat{t}\xi$. Also recall that we are assuming that if $\mathbf{b}_i = \mathbf{b}_j$, then $p_i = p_j$. Thus, for each a, $|\overline{w}_a^{d_1}| \xi = p_a = p_i$. So, $\overline{t}\xi =_{wt} p_i = q + r \vee \widehat{t} \vee \overline{t}$. Consequently,

$$p_i \xi =_{wt} (q + r \lor \widehat{t} \lor \overline{t}) \xi =_{wt} q \xi + r \xi \lor \widehat{t} \xi \lor \overline{t} \xi =_{wt} q \xi + r \xi \lor \widehat{t} \lor \overline{t}$$

$$q \xi + r \xi \lor \widehat{t} \lor (q + r \lor \widehat{t} \lor \overline{t}) =_{wt} q \xi + r \xi \lor (q + r \lor \widehat{t} \lor \overline{t}) =_{wt} q \xi + r \xi \lor p_i.$$

Hence, taking $q_i = q$ and $r_i = r$ suffices for this case.

CASE: $\mathbf{b}_i = \mathsf{N}_{\Box_{d_1+1}}$. By the induction hypothesis, we may take p_i to be of the form $q \lor r \lor \hat{t} \lor \bar{t}$, where q is $\mathsf{T}_{\Box_{d_1+1}}$ -strict, r is strictly $\mathsf{T}_{\Box_{d_1+1}}$ -chary, $\hat{t} = \bigvee_{a=0}^{b'_{d_1+1}} |u_a^{d_1+1}|$, and $\bar{t} = \bigvee_{a=0}^{c'_{d_1+1}} |\overline{u}_a^{d_1+1}|$. By an argument similar to the one for the previous case it follows that taking $q_i = q$ and $r_i = r$ suffices for this case too.

CASES: $\mathbf{b}_i = \mathsf{N}_{\diamond_{d_1+1}}, \ldots, \mathbf{b}_{\max}$. These cases follow from essentially the same as arguments as given for the previous two cases.

Lemma 44

Henceforth we assume that each p_i is as in Lemma 44 and, in the cases where $N_{\Box_{d_1}} \leq \mathbf{b}_i$, q_i and r_i are as in that lemma too. For each $n \in \omega$, define

 $e^{(n)}$ = the β -normal form of the *n*-level unfolding of *e*'s crec-recursion, (10.2)

where β - and η -reductions are used to neaten up things as in the definition of $e^{(1)}$. So, $e^{(0)} = e$ and $e^{(1)} =$ our prior definition of $e^{(1)}$. Let $\xi^{(0)} =$ the empty substitution and $\xi^{(n+1)} = \xi \circ \xi^{(n)} =$ the (n+1)-fold composition ξ . It follows that, with respect to $\widehat{\Gamma}$; , for each *i* and *n*, $(p_i \xi^{(n)})$ is a size bound for *i*-th argument expression of every *f*-application in $e^{(n)}$.

Lemma 45 (The n step lemma). For each i and n:

$$\begin{array}{ll} (a) \ p_i \, \xi^{(n)} &=_{\mathrm{wt}} \ p_i \ when \ \mathbf{b}_i \leq : \mathsf{N}_{\Box_{d_1}}. \\ (b) \ p_i \, \xi^{(n)} &\leq_{\mathrm{wt}} \ (q_i \lor r_i) \, \xi^{(n)} \lor p_i \ when \ \mathsf{N}_{\Box_{d_1}} \leq : \mathbf{b}_i = \mathsf{N}_{\Box_d}. \\ (c) \ p_i \, \xi^{(n)} &\leq_{\mathrm{wt}} \ n \ast (q_i \, \xi^{(n)}) + (r_i \, \xi^{(n)}) \lor p_i \ when \ \mathsf{N}_{\Box_{d_1}} \leq : \mathbf{b}_i = \mathsf{N}_{\diamond_d}. \end{array}$$

Proof. Part (a) follows directly from Lemma 44(a). For parts (b) and (c) we first note that by monotonicity we have that, for all k and i, $(q_i \vee r_i)\xi^{(k)} \leq_{wt} (q_i \vee r_i)\xi^{(k+1)}$. Now, for part (b), it follows immediately from Lemma 44(b) that, for each n and i, we have $p_i \xi^{(n)} =_{wt} \left(\bigvee_{j=0}^n (q_i \vee r_i)\xi^{(j)} \right) \vee p_i$. Hence by the noted monotonicity of $(q_i \vee r_i)\xi^{(\cdot)}$, part (b) follows. For part (c), first fix i such that $\mathbf{b}_i = \mathsf{N}_{\diamond_d}$ with $d \geq d_1$. It follows from an easy induction that for all n, $p_i \xi^{(n)} \leq_{wt} (\sum_{j=1}^n q_i \xi^{(j)}) + (\bigvee_{j=1}^n r_i \xi^{(j)}) \vee p_i$; note the parallel to the argument for the prn-case of Proposition 1. Hence by monotonicity of $q_i\xi^{(\cdot)}$ and $r_i\xi^{(\cdot)}$, part (c) follows.

Lemma 45 By Lemma 45(a) and (4.2) we have

Lemma 46 (Termination). $\mathcal{L}_{wt}[[p_0]] |\rho| \geq the maximum depth of e's crec-recursion.$

For each *i* with $\mathbf{b}_i \leq \mathbf{N}_{\Box_{d_1}}$, let $p'_i = p_i$. For $\sigma = \mathbf{N}_{\diamond_{d_1}}, \ldots, \mathbf{b}_{\max}$ in turn, we inductively define θ_{σ} to be the substitution $[x_j := p'_j + \mathbf{b}_j \leq \sigma]$ and also define, for each *i* with $\mathbf{b}_i = \sigma$:

$$p'_{i} \stackrel{\text{def}}{=} \begin{cases} (r_{i} \theta_{\sigma}) \lor p_{i}, & \text{if } \sigma \text{ is oracular;} \\ p'_{0} \cdot (q_{i} \theta_{\sigma}) + (r_{i} \theta_{\sigma}) \lor p_{i}, & \text{if } \sigma \text{ is computational.} \end{cases}$$

By Lemma 41, for each *i*, we can effectively find a manifestly \mathbf{b}_i -safe p''_i with $p'_i \leq_{wt} p''_i$.

Lemma 47 (Final sizes). For each *i*, p''_i is a manifestly \mathbf{b}_i -safe polynomial size-bound on the *i*-th argument expression in the final step of the crec-recursion in e.

Proof. For each *i* with $\mathbf{b}_i \leq \mathbb{N}_{\Box_{d_1}}$, the conclusion follows from Lemma 45(a). For the $\sigma = \mathsf{N}_{\diamond_{d_1}}$ case, fix an *i* with $\mathbf{b}_i = \mathsf{N}_{\diamond_{d_1}}$. Then the bound for this case follows from Lemmas 45(c) and 46. The $\mathsf{N}_{\Box_{d_1+1}}$ through \mathbf{b}_{\max} cases follow similarly.

Lemma 47

By the induction hypothesis, there exists p_B , a manifestly $|\mathbf{b}_0|$ -safe polynomial sizebound for B (as in (10.1)) with respect to Γ ; $f:\gamma$. By Lemma 41, we can effectively find a manifestly \mathbf{b}_0 -safe \hat{p} such that $p_B [|f| := \underline{0}_{\gamma}, |x_1| := p'_1, \ldots, |x_k| := p'_k] \leq_{wt} \hat{p}$. The effect of the substitution on p_B is to trivialize |f| and replace each $|x_i|$ with the final size bound from Lemma 47. It follows that \hat{p} is a manifestly \mathbf{b}_0 -safe size bound for the value returned by final step of the crec-recursion. Since TailPos(f, A), \hat{p} is also a size bound on the value returned by the entire (tail) recursion. Thus, $p_E = \lambda |x_1|, \ldots, |x_k| \cdot \hat{p}$ suffices for the crec case.

Theorem 43

11. An Abstract machine

Our next major goal is to show that every ATR expression is computable within a second-order polynomial time-bound (Theorem 79). Before formalizing time bounds, we first need to make precise what is being bounded. Below we set out the abstract machine that provides the operational semantics of PCF, BCL, and ATR and, based on this, §11.2 introduces and justifies our notion of the time cost of an expression evaluation.

$$\begin{array}{ccc} \left((B \ e), \hat{\rho}, \kappa \right) & \to & \left(e, \hat{\rho}, \langle \mathsf{op}, B, \kappa \rangle \right) & (a) \\ \left(v, \hat{\rho}, \langle \mathsf{op}, B, \kappa \rangle \right) & \to & \left(\delta_1(B, v), \{\}, \kappa \right) & (b) \end{array}$$

$$((\operatorname{\mathsf{down}} e e'), \hat{\rho}, \kappa) \to (e, \hat{\rho}, \langle \operatorname{\mathsf{dn}}, e', \hat{\rho}, \kappa \rangle) \tag{c}$$

$$(v', \rho', \langle \mathsf{dn}', v, \kappa \rangle) \rightarrow (\delta_2(v, v'), \{\}, \kappa)$$
 (e)

$$(x, \hat{\rho}, \kappa) \rightarrow (v, \hat{\rho}', \kappa)$$
, where $\langle v, \hat{\rho}' \rangle = \hat{\rho}(x)$ (f)

$$\begin{array}{ccc} ((e\ e'),\hat{\rho},\kappa) &\to & (e,\hat{\rho},\langle \arg, e',\hat{\rho},\kappa\rangle) \\ (v,\hat{\rho},\langle \arg, e',\hat{\rho}',\kappa\rangle) &\to & (e'\hat{\rho}',\langle \operatorname{fun},v,\hat{\rho},\kappa\rangle) \\ (h) \end{array}$$

$$\begin{pmatrix} v, \rho', \langle \mathsf{fun}, O, \hat{\rho}, \kappa \rangle \end{pmatrix} \to \begin{pmatrix} O(v'), \{\}, \kappa \end{pmatrix}$$
(j)

$$\left((\text{if } e_{?} \text{ then } e_{t} \text{ else } e_{f}), \hat{\rho}, \kappa \right) \rightarrow \left(e_{?}, \hat{\rho}, \langle \text{test}, e_{t}, e_{f}, \hat{\rho}, \kappa \rangle \right) \tag{k}$$

$$(v_{?}, \hat{\rho}', \langle \mathsf{test}, e_{f}, \hat{\rho}, \kappa \rangle) \quad \to \quad \begin{cases} (e_{t}, \hat{\rho}, \kappa), & \text{if } v_{?} \neq \epsilon; \\ (e_{f}, \hat{\rho}, \kappa), & \text{if } v_{?} = \epsilon. \end{cases}$$
 (l)

$$\left((\mathsf{fix}\,(\lambda x \, \boldsymbol{\cdot}\, e)), \hat{\rho}, \kappa\right) \quad \to \quad \left(e[x := (\mathsf{fix}\,(\lambda x \, \boldsymbol{\cdot}\, e))], \hat{\rho}, \kappa\right) \tag{(m)}$$

$$\begin{array}{rcl} \left((\operatorname{prn} e), \hat{\rho}, \kappa \right) & \to & \left(e', \hat{\rho}, \kappa \right), & \text{where} \\ e' & = & \lambda y \, \textbf{.} (\text{if } y \neq \epsilon \ \text{then} \ (e \ y \ (\operatorname{prn} e \ (d \ y))) \ \text{else} \ (e \ \epsilon \ \epsilon)) \end{array}$$

Figure 20: The CE	K-rewrite rules
-------------------	-----------------

11.1. The CEK machine. The operational semantics for PCF, BCL, and ATR are provided by the abstract machine whose rules are given in Figure 20. The machine is based on Felleisen and Friedman's CEK-machine [FF87] as presented by Felleisen and Flatt [FF06]. States in this machine are triples consisting of: (i) an expression to be reduced or else a value, (ii) an environment, and (iii) a continuation. CEK-environments, closures, and values are defined recursively by:

> CEK-Environments = Variables $\stackrel{\text{finite}}{\rightarrow}$ Closures. Closures = $(\text{Terms} \cup \text{Values}) \times \text{CEK-Environments}.$ Values = Strings \cup Oracles $\cup \lambda$ -Terms.

An *oracle* is just an element of $\bigcup_{k>0} \mathbf{TC}_{(N^k)\to N}$. Note that the result of applying an oracle value $O \in \mathbf{TC}_{(N^{k+1})\to N}$ to a $v \in N$ is the oracle value $O(v) \in \mathbf{TC}_{(N^k)\to N}$, where k > 0. The continuations should be self-explanatory from the rules—and if not, see [FWH01].

The CEK rules use the following variables (plain and decorated) with indicated ranges. B:Basic-Operations (i.e., $c_0, c_1, d, t_0, and t_1$); κ :Continuations; e:Terms; O:Oracles; $\hat{\rho}$:CEK-Environments; v:Values; and x:Variables. Also, $\delta_1(B, v)$ returns the value of the given basic-operation on the given value and $\delta_2(v, v')$ returns down(v, v'). For each expression e and CEK-environment $\hat{\rho}$ with $FV(e) \subseteq \text{preimage}(\hat{\rho})$,

$$\operatorname{eval}_{\operatorname{CEK}}(e,\hat{\rho}) \stackrel{\text{def}}{=} \begin{cases} v, & \text{if } (e,\hat{\rho},\operatorname{halt}) \to^* (v,\hat{\rho}',\operatorname{halt});\\ \text{undefined}, & \text{if there is no such } v \text{ and } \hat{\rho}'. \end{cases}$$

For each ordinary environment $\rho = \{x_1 \mapsto v_1, \ldots, x_k \mapsto v_k\}$, let ρ^* be the corresponding CEK-environment, i.e, $\{x_1 \mapsto (v_1, \{\}), \ldots, x_k \mapsto (v_k, \{\})\}$, and let $\mathsf{eval}_{\mathsf{CEK}}(e, \rho) \stackrel{\text{def}}{=} \mathsf{eval}_{\mathsf{CEK}}(e, \rho^*)$.

11.2. The CEK cost model. We assume that the underlying model of computation is along the lines of Kolmogorov and Uspenskii's [KU58] "pointer machines" or Schönhage's *storage modification machines* [Sch80]. A string is represented by a linked list of **0**'s and **1**'s. We take the cost of evaluating an expression e to be the sum of the cost of the steps involved in evaluating e on the CEK machine. We charge unit cost for for CEK-steps that do not involve operations on strings or else carry out operations that work on just the fronts of strings (e.g., c_a , d, and t_a). For steps that involve copying or examining the entirety of arbitrary strings (rules (e), (f), and (j)), our charge involves the sum of the lengths of the strings involved. Specifically:

(j) Oracle application. Applying this rule has $\cot \underline{1} \vee |O(v)|$ when O(v) is of base type and $\underline{1}$ otherwise. (When O(v) is of base type, an application of the oracle pops into memory a string of length |O(v)|. We view the action of entering this string in memory, character-by-character, as observable.)

(e) δ_2 application. Applying this rule has cost $\underline{1} + |v| + |v'|$. (down looks at the entirety of its arguments.)

(f) Environment application. Applying this rule has $\cot \underline{1} \vee |\hat{\rho}(x)|$ when $\hat{\rho}(x)$ is of a base type and $\underline{1}$ otherwise. (Since our CEK machine starts with an arbitrary environment, the environment is essentially another oracle.)

Given this assignments of costs, we introduce:

Definition 48. For each expression e and CEK-environment $\hat{\rho}$,

$$\operatorname{cost}_{\operatorname{CEK}}(e,\hat{\rho}) \stackrel{\text{def}}{=} \begin{cases} s, & \text{if } \operatorname{eval}_{\operatorname{CEK}}(e,\hat{\rho}) \text{ is defined, where } s \text{ is the sum of the costs of the steps in this CEK-computation;} \\ undefined, & \text{otherwise.} \end{cases}$$

and for each ordinary environment ρ , $\operatorname{cost}_{\operatorname{CEK}}(e,\rho) \stackrel{\text{def}}{=} \operatorname{cost}_{\operatorname{CEK}}(e,\rho^*)$.

We note that the standard proof that storage modification machines and Turing machines are polynomially-related models of computation [Sch80] straightforwardly extends to show that, at type-levels 1 and 2, our CEK model of computation and cost is (second-order) polynomially related to Kapron and Cook's oracle Turing machines under their answerlength cost model [KC96].

12. Time bounds

As the next step towards showing polynomial time-boundedness for ATR, the present section sets up a formal framework for working with time bounds. We start by noting the obvious: Run time is **not** an extensional property of programs. That is, \mathcal{V}_{wt} -equivalent

expressions can have quite distinct run time properties. Because of this we introduce \mathcal{T} , a new semantics for ATR that provides upper bounds on the time complexity of expressions.

The setting. Our framework for time complexities uses the following simple setting.

CEK costs. Time costs are assigned to ATR-computations via the CEK cost model.

Worst-case bounds. $\mathcal{T}[\![e]\!]$ will provide a worst-case upper bound on the CEK cost of evaluating e, but not necessarily a tight upper bound.

No free lunch. All evaluations have positive costs. This even applies to "immediately evaluating" expressions (e.g., λ -expressions), since checking whether something "immediate-evaluates" counts as a computation with costs.

Inputs as oracles. We treat each type-level 1 input f as an oracle. In a time-complexity context this means that f is thought of answering any query in one time step, or equivalently, any computation involved in determining the reply to a query happens unobserved off-stage. Thus the cost of a query to f involves only (i) the time to write down a query v, and (ii) the time to read the reply f(v). The times (i) and (ii) are bounded by roughly |v| and |f|(|v|), respectively. Thus our time bounds will ultimately be expressed in terms of the *lengths* of the values of free and input variables.

Currying and time complexity. In common usage, "the time complexity of e" can mean one of two things. When e is of base type, the phrase usually refers to the time required to compute the value of e. We might think of this as *time past*—the time it took to arrive at e's value. When e is of an arrow type and thus describes a procedure, the phrase usually refers to the function that, given the sizes of arguments, returns the maximum time the procedure will take when run on arguments of the specified sizes. We might think of this as *time in possible futures* in which e's value is applied. An expression can have both a past and futures of interest. Consider $(e_0 e_1)$ where e_0 is of type $N_{\varepsilon} \rightarrow N_{\varepsilon} \rightarrow N_{\circ}$ and e_1 is of type N_{ε} . Then $(e_0 e_1)$ has a time complexity in the first sense as it took time to evaluate the expression, and, since $(e_0 e_1)$ is of type $N_{\varepsilon} \rightarrow N_{\diamond}$, it also has a time complexity in the first sense and the potential/futures part of e_0 's time complexity must account for the multiple senses of time complexity just attributed to $(e_0 e_1)$. Type-level-2 expressions add further twists to the story. Our treatment of time complexity takes into account these extended senses.

Costs and potentials. In the following the time complexity of an expression e always has two components: a *cost* and a *potential*. A cost is always a positive (tally) integer and is intended to be an upper bound on the time it takes to evaluate e. The form of a potential depends on the type of e. Suppose e is of a base (i.e., string) type. Then e's potential is intended to be an upper bound on the length of its value, an element of ω . The length of e's value describes the potential of e in the sense that when e's value is used, its length is the only facet of the value that plays a role in determining time complexities. Now suppose eis of type, say, $N_{\varepsilon} \to N_{\diamond}$. Then e's potential will be an $f_e \in (\omega \to \omega \times \omega)$ that maps a $p \in \omega$ (the length/potential of the value of an argument of e) to a $(c_r, p_r) \in \omega \times \omega$ where c_r is the cost of applying the value of e to something of length p and p_r is the length/potential of the result. Note that (c_r, p_r) is a time complexity for something of base type. Generalizing from this, our motto will be: The potential of a type- $(\sigma \rightarrow \tau)$ thing is a map from potentials of type- σ things to time complexities of type- τ things.²¹

Our first task in making good on this motto is to situate time complexities in a suitable semantic model.²²

A model for time complexities. The *time types* are the result of the following translations $(\|\cdot\| \text{ and } \langle\!\langle \cdot \rangle\!\rangle)$ of ATR types:

 $\|\sigma\| \stackrel{\text{def}}{=} \mathsf{T} \times \langle\!\langle \sigma \rangle\!\rangle. \qquad \langle\!\langle \mathsf{N}_\ell \rangle\!\rangle \stackrel{\text{def}}{=} \mathsf{T}_\ell. \qquad \langle\!\langle \sigma \to \tau \rangle\!\rangle \stackrel{\text{def}}{=} \langle\!\langle \sigma \rangle\!\rangle \to \|\tau\|\,.$

So, $\|N_{\ell_1} \to N_{\ell_2} \to N_{\ell_0}\| = T \times (T_{\ell_1} \to T \times (T_{\ell_2} \to T \times T_{\ell_0}))$ and $\|(N_{\ell_1} \to N_{\ell_2}) \to N_{\ell_0}\| = T \times ((T_{\ell_1} \to T \times T_{\ell_2}) \to T \times T_{\ell_0})$. The time types are thus a subset of the simple product types over $\{T, T_{\epsilon}, T_{\diamond}, T_{\Box\diamond}, \ldots\}$. The intent is that T is the type of costs, the T_{ℓ} 's help describe lengths, $\|\gamma\|$ is the type of complexity bounds of type- γ objects, and $\langle\!\langle \gamma \rangle\!\rangle$ is the type of potentials of type- γ objects. (Note: $\langle\!\langle \sigma \to \tau \rangle\!\rangle$'s definition parallels the motto.)

Our proof of polynomial time-boundedness for ATR (Theorem 79) needs to intertwine the size estimates implicit in potentials and the size bounds of Theorem 43. The semantics for the time types thus needs to be an extension of the \mathcal{L}_{wt} -semantics. To define this extension we use a combinator, Pot, defined in Definition 60 below. For the moment it is enough to know that, for each ATR-type σ and $p \in \mathcal{L}_{wt}[\![\langle \sigma \rangle \rangle]\!]$, $\mathsf{Pot}(p) \in \mathcal{L}_{wt}[\![\sigma|]\!]$ is a canonical projection of p to a type- $|\sigma|$ size bound. Following the definition of Pot, Lemma 61 notes that all of the notions introduced between here and there mesh properly.

Definition 49 (\mathcal{L}_{wt} extended to the time types). Suppose σ and τ are ATR types. Then $\mathcal{L}_{wt}[\![\|\sigma\|]\!] \stackrel{\text{def}}{=} \omega \times \mathcal{L}_{wt}[\![\langle\!\langle \sigma \rangle\!\rangle]\!]$ and $\mathcal{L}_{wt}[\![\langle\!\langle \sigma \rangle\!\rangle]\!]$ is inductively defined by $\mathcal{L}_{wt}[\![\langle\!\langle n_{\ell} \rangle\!\rangle]\!] \stackrel{\text{def}}{=} \omega$ and $\mathcal{L}_{wt}[\![\langle\!\langle \sigma \rangle\!\rangle]\!] \stackrel{\text{def}}{=}$ the set of all monotone Kleene-Kreisel functionals $f: \mathcal{L}_{wt}[\![\langle\!\langle \sigma \rangle\!\rangle]\!] \to \mathcal{L}_{wt}[\![\|\tau\|]\!]$ such that: (i) $\mathsf{Pot}(f) \in \mathcal{L}_{wt}[\![|\sigma \to \tau|]\!]$ and (ii) $\mathsf{Pot}(f(p_1)) = \mathsf{Pot}(f(p_2))$ whenever $\mathsf{Pot}(p_1) = \mathsf{Pot}(p_2)$.

Condition (i) above restricts $\mathcal{L}_{wt}[\![\langle\!\langle \sigma \to \tau \rangle\!\rangle]\!]$ so that the projection Pot acts as advertised. Condition (ii) restricts each $f \in \mathcal{L}_{wt}[\![\langle\!\langle \sigma \to \tau \rangle\!\rangle]\!]$ so that the size information in f(p) depends only on the size information in p.

We can now define the \mathcal{T} (time-complexity) and \mathcal{P} (potential) interpretations of the ATR types. (The \mathcal{P} -interpretation is a notational convenience.)

Definition 50. Suppose σ is an ATR-type. Then $\mathcal{T}[\![\sigma]\!] \stackrel{\text{def}}{=} \mathcal{L}_{\text{wt}}[\![\|\sigma\|]\!]$ and $\mathcal{P}[\![\sigma]\!] \stackrel{\text{def}}{=} \mathcal{L}_{\text{wt}}[\![\langle\!\langle \sigma \rangle\!\rangle]\!]$.

The \mathcal{T} -interpretation of constants and oracles. The following two definitions introduce a translation from the \mathcal{V}_{wt} model into the \mathcal{T} model. We use this translation to assign time complexities to program inputs: string constants and oracles.

Definition 51. Let $||a|| \stackrel{\text{def}}{=} (\underline{1} \vee |a|, \langle\langle a \rangle\rangle)$ and $\langle\langle a \rangle\rangle \stackrel{\text{def}}{=} |a|$ for each $a \in \mathcal{V}_{\text{wt}}[\![\mathsf{N}_{\ell}]\!]$.

²¹In a more general setting (e.g., call-by-name), a ($\sigma \rightarrow \tau$) potential is a map from σ -time-complexities to τ -time-complexities, as an operator may be applied to an unevaluated operand.

 $^{^{22}}$ **N.B.** The time-complexity cost/potential distinction appears in prior work [San90, Shu85, VS03]. Remark 82 below discusses this prior work and how it relates to ours.

By Lemma 61(a) below, $||a|| \in \mathcal{T}[N_{\ell}]$. We view ||a|| as the time complexity of the string/integer constant a. The interpretation of the cost component of ||a|| is that the cost of evaluating the constant a is the cost of writing down a character by character. (When $a = \epsilon$, we still charge <u>1</u>.)

Definition 52. Let $||f|| \stackrel{\text{def}}{=} (\underline{1}, \langle\!\langle f \rangle\!\rangle)$ and $\langle\!\langle f \rangle\!\rangle \stackrel{\text{def}}{=} \lambda p \in \mathcal{P}[\![\sigma]\!] \cdot \max\{||(fv)|| : \langle\!\langle v \rangle\!\rangle \le p\}$ for each $f \in \mathcal{V}_{\text{wt}}[\![\sigma \to \tau]\!]$.

By Lemma 61(a) below, $||f|| \in \mathcal{T}[\![\sigma \to \tau]\!]$. We view ||f|| as the time complexity of f as an oracle: the only time costs associated with applying f are those involved in setting up applications of f and reading off the results. Recall that under call-by-value, a λ -expression immediately evaluates to itself. The function-symbol f will be treated analogously to a λ -term. Hence, the cost component of ||f|| is <u>1</u>. The definition of $\langle\langle f \rangle\rangle$ parallels both our informal discussion of the notion of the potential of a type-level 1 function and the definition of the length of functions of type levels 1 and 2 in §2.10. One can show that when f is a type-level 2, $\langle\langle f \rangle\rangle$ is total. (The argument is similar to the proof of the totality of the type-level 2 notion of length defined by (2.3) in §2.10.)

Definition 51 and the type-level 1 part of Definition 52 describe the time complexities of possible ATR inputs. The following lemma unpacks the definition of $\langle \langle f \rangle \rangle$ for f of type-level 1. The proof is a straightforward induction and hence omitted.

Lemma 53. For $f \in \mathcal{V}_{wt}[\![(\mathsf{N}_{\ell_1},\ldots,\mathsf{N}_{\ell_k}) \to \mathsf{N}_{\ell_0}]\!]$, $\langle\!\langle f \rangle\!\rangle = q_1$ where $q_i = \lambda p_i \in \omega \cdot (\underline{1}, q_{i+1})$ (for $1 \leq i < k$) and $q_k = \lambda p_k \in \omega \cdot (\underline{1} \vee |f|(p_1,\ldots,p_k), |f|(p_1,\ldots,p_k))$.

T-Applications.

Definition 54.

(a) Suppose $t_0 \in \mathcal{T}[\![\sigma \to \tau]\!]$ and $t_1 \in \mathcal{T}[\![\sigma]\!]$, where $t_0 = (c_0, p_0), t_1 = (c_1, p_1)$, and $(c_r, p_r) = p_0(p_1)$. Then $t_0 \star t_1 \stackrel{\text{def}}{=} (c_0 + c_1 + c_r + \underline{3}, p_r)$.

(b) Suppose $t_0 \in \mathcal{T}[\![(\sigma_1, \ldots, \sigma_k) \to \tau]\!], t_1 \in \mathcal{T}[\![\sigma_1]\!], \ldots, t_k \in \mathcal{T}[\![\sigma_k]\!]$. Then $t_0 \star \vec{t} \stackrel{\text{def}}{=} t_0 \star t_1 \star \ldots \star t_k$. (The \star operation left associates.)

By Lemma 61(b) below, $t_0 \star t_1 \in \mathcal{T}[\![\tau]\!]$ when $t_0 \in \mathcal{T}[\![\sigma \to \tau]\!]$ and $t_1 \in \mathcal{T}[\![\sigma]\!]$. Suppose that t_0 (respectively, t_1) is the time complexity of a type- $(\sigma \to \tau)$ expression e_0 (respectively, type- σ expression e_1). Then $t_0 \star t_1$ is intended to be the time complexity of $(e_0 e_1)$. The cost component of $t_0 \star t_1$ is: (the cost of evaluating e_0) + (the cost of evaluating e_1) + (the cost of applying e_0 's value to e_1 's value) + <u>3</u>, where the <u>3</u> is the CEK-overhead of an application. The potential component is simply the potential of the result of the application. The next lemma works out of the effect of the \star operation for type-level 1 oracles.

Lemma 55. Suppose $f \in \mathcal{V}_{wt}[\![(\mathsf{N}_{\ell_1}, \dots, \mathsf{N}_{\ell_k}) \to \mathsf{N}_{\ell_0}]\!], v_1 \in \mathcal{V}_{wt}[\![\mathsf{N}_{\ell_1}]\!], \dots, v_k \in \mathcal{V}_{wt}[\![\mathsf{N}_{\ell_k}]\!].$ Then

$$|f|| \star \overrightarrow{\|v\|} = \left(\left(\sum_{i=1}^{k} (\underline{1} \lor |v_i|) \right) + \underline{1} \lor |f|(\overrightarrow{|v|}) + \underline{5k-1}, |f|(\overrightarrow{|v|}) \right),$$
(12.1)

where $\overrightarrow{\|v\|}$ abbreviates $\|v_1\|, \ldots, \|v_k\|$ and $\overrightarrow{|v|}$ abbreviates $|v_1|, \ldots, |v_k|$.

The proof is a straightforward calculation. Equation (12.1) can be interpreted as giving an upper bound on the time complexity of applying an oracle f to arguments v_1, \ldots, v_k . Let us consider the cost component of the k = 1 and k = 2 cases of (12.1) in more detail.

<u>1</u>	=	the cost of evaluating f
$\underline{1} \vee v_1 $	=	the cost of evaluating v_1 , i.e., the cost of writing down the value v_1
$\underline{1} \vee f (v_1)$	=	the cost of applying f to v_1 , i.e., the cost of writing down $f(v_1)$'s value
<u>3</u>	=	the overhead of the application

Figure 21: Break down of the cost component of $||f|| \star ||v_1||$

For k = 1, the right-hand side of (12.1) simplifies to: $((\underline{1} \lor |v_1|) + \underline{1} \lor |f|(|v_1|) + \underline{4}, |f|(|v_1|))$. Its cost component is broken down in Figure 21. For k = 2, the right-hand side of (12.1) simplifies to: $((\underline{1} \lor |v_1|) + (\underline{1} \lor |v_2|) + \underline{1} \lor |f|(|v_1|, |v_2|) + \underline{9}, |f|(|v_1|, |v_2|))$. We leave it to the reader to break down its cost component.

 \mathcal{T} -Environments. As a companion to \mathcal{T} -application we shall define an analogue of currying in \mathcal{T} . First, we introduce \mathcal{T} -environments. Recall that in a call-by-value language, variables name *values* [Plo75], i.e., the end result of a (terminating) evaluation. Thus, a value does not need to be evaluated again, at least no more than an input value does. Hence, if a \mathcal{T} -environment maps a variable to a type- γ time complexity (c, p), then c should be: $\underline{1} \vee p$, when γ is a base type, and $\underline{1}$, when γ is an arrow type.

Definition 56. Suppose σ and τ vary over ATR types and Γ ; Δ is an ATR is type context.

- (a) $\|\Gamma; \Delta\| \stackrel{\text{def}}{=} \{ x \mapsto \|\sigma\| \colon (\Gamma; \Delta)(x) = \sigma \}.$
- (b) For $p \in \mathcal{P}[[N_{\ell}]]$, $\mathsf{val}(p) \stackrel{\text{def}}{=} (\underline{1} \lor p, p)$.
- (c) For $p \in \mathcal{P}\llbracket \sigma \to \tau \rrbracket$, $\mathsf{val}(p) \stackrel{\text{def}}{=} (\underline{1}, p)$.
- (d) $\mathcal{T}_{\text{val}}\llbracket \sigma \rrbracket \stackrel{\text{def}}{=} \{ \operatorname{val}(p) \mid p \in \mathcal{P}\llbracket \sigma \rrbracket \}.$

(e) $\mathcal{T}\llbracket\Gamma; \Delta\rrbracket$ is the set of all finite maps of the form $\{x_1 \mapsto t_1, \ldots, x_k \mapsto t_k\}$, where $\{x_1, \ldots, x_k\} = \operatorname{preimage}(\Gamma; \Delta)$, and, for $i = 1, \ldots, k$, $t_i \in \mathcal{T}_{\operatorname{val}}\llbracket(\Gamma; \Delta)(x_i))\rrbracket$.

(f) For each $\rho \in \mathcal{V}_{wt}[\Gamma; \Delta]$, define $\|\rho\| \in \mathcal{T}[\Gamma; \Delta]$ by $\|\rho\|(x) = \|\rho(x)\|$. Such as $\|\rho\|$ is called an *oracle environment*.

Convention: We use ρ as a variable over \mathcal{T} -environments. **N.B.** Not every ρ of interest is an oracle environment.

 \mathcal{T} -currying. Here then is our time-complexity analogue to currying. Recall that $\mathcal{T}[\![\Gamma; \Delta \vdash e: \tau]\!]$ will be (when we get around to defining it) a function from $\mathcal{T}[\![\Gamma; \Delta]\!]$ to $\mathcal{T}[\![\tau]\!]$.

Definition 57. Suppose (i) $\Gamma; \Delta$ is a ATR type context with $(\Gamma; \Delta)(x_i) = \sigma_i$, for $i = 1, \ldots, k$; (ii) $\Gamma'; \Delta'$ is the result of removing $x_1: \sigma_1$ from $\Gamma; \Delta$; and (iii) X is a function from $\mathcal{T}\llbracket\Gamma; \Delta\rrbracket$ to $\mathcal{T}\llbracket\tau\rrbracket$. Then $\Lambda_{\star}(x_1, X)$ is the function from $\mathcal{T}\llbracket\Gamma'; \Delta'\rrbracket$ to $\mathcal{T}\llbracket\sigma_1 \to \tau\rrbracket$ given by:

$$\Lambda_{\star}(x_1, X) \, \varrho' \stackrel{\text{def}}{=} \left(\underline{1}, \, \lambda p \in \mathcal{P}[\![\sigma_1]\!] \, (X \, (\varrho' \cup \{ x_1 \mapsto \mathsf{val}(p) \})) \right), \tag{12.2}$$

where $\varrho' \in \mathcal{T}\llbracket\Gamma'; \Delta'\rrbracket$. Also, $\Lambda_{\star}(x_1, x_2, \dots, x_k, X) \stackrel{\text{def}}{=} \Lambda_{\star}(x_1, \Lambda_{\star}(x_2, \dots, x_k, X))$ when k > 1.

Note the complementary roles of Λ_{\star} and \star : Λ_{\star} shifts the past (the cost) into the future (the potential) and \star shifts part of the future (the potential) into the past (the cost). This being complexity theory, there are carrying charges on all this shifting. This is illustrated in the next lemma that shows how Λ_{\star} and \star interact. First, we introduce:

Definition 58. $dally(d, (c, p)) \stackrel{\text{def}}{=} (c + d, p)$ for $d \in \omega$ and (c, p), a time complexity.

Lemma 59 (Almost the η -law). Suppose Γ , Δ , X, \vec{x} , $\vec{\sigma}$, and τ are as in Definition 57. Let $\Gamma'; \Delta'$ be the result of removing $x_1: \sigma_1, \ldots, x_k: \sigma_k$ from $\Gamma; \Delta$. Let $\varrho \in \mathcal{T}[\![\Gamma; \Delta]\!]$ and let ϱ' be the restriction of ϱ to preimage($\Gamma'; \Delta'$). Then

$$\left(\Lambda_{\star}(x_1,\ldots,x_k,X)\,\varrho'\right)\star\varrho(x_1)\star\ldots\star\varrho(x_k) = dally(\underline{5\cdot k+4} + \sum_{i=1}^k c_i, X\,\varrho), \quad (12.3)$$

where $(c_1,p_1) = \varrho(x_1),\ldots,(c_k,p_k) = \varrho(x_k).$

The lemma's proof is another straightforward calculation.

Projections. The next definition introduces a way of recovering more conventional bounds from time complexities. Note, by Definitions 51 and 52, and Lemmas 53 and 55, when v is a string constant or a type-1 oracle, the value of ||v|| is a function of the value of ||v||. So, by an abuse of notation, we treat ||v|| as a function of |v| for such v.

Definition 60. Suppose σ and $(\sigma_1, \ldots, \sigma_k) \to \mathsf{N}_{\ell}$ are ATR types.

(a) For each $t \in \mathcal{T}[\sigma]$, let $cost(t) \stackrel{\text{def}}{=} \pi_1(t)$ and $pot(t) \stackrel{\text{def}}{=} \pi_2(t)$. (So, t = (cost(t), pot(t)).)

(b) For each $t \in \mathcal{T}[\![N_\ell]\!]$, let $\mathsf{Cost}(t) = \mathit{cost}(t)$ and $\mathsf{Pot}(t) = \mathit{pot}(t)$ and, for each $t \in \mathcal{T}[\![(\sigma_1, \ldots, \sigma_k) \to \mathsf{N}_\ell]\!]$, let:

 $\mathsf{Cost}(t) \stackrel{\text{def}}{=} \lambda |\overrightarrow{v}| \cdot cost(t \star ||\overrightarrow{v}||). \qquad \mathsf{Pot}(t) \stackrel{\text{def}}{=} \lambda |\overrightarrow{v}| \cdot pot(t \star ||\overrightarrow{v}||).$

where $|\overrightarrow{v}|$ abbreviates $|v_1| \in \mathcal{L}_{wt}[\sigma_1], \ldots, |v_k| \in \mathcal{L}_{wt}[\sigma_k]$ and $||\overrightarrow{v}||$ abbreviates $||v_1||, \ldots, ||v_k||$. (So, $t \star ||\overrightarrow{v}|| = (\mathsf{Cost}(t)(|\overrightarrow{v}|), \mathsf{Pot}(t)(|\overrightarrow{v}|))$.) We call $\mathsf{Cost}(t)$ and $\mathsf{Pot}(t)$, respectively, the base cost and base potential of t.

(c) For each $p \in \mathcal{P}[\![\sigma]\!]$, let $\mathsf{Pot}(p) \stackrel{\text{def}}{=} \mathsf{Pot}((\underline{1}, p))$.

Suppose t is the time complexity of e of type $(\vec{\sigma}) \to \mathsf{N}_{\ell}$. Then both $\mathsf{Cost}(t)$ and $\mathsf{Pot}(t)$ are functions of the sizes of possible arguments of e. The intent is that $\mathsf{Cost}(t)(|\vec{v}|)$ is an upper bound on the time cost of first evaluating e and then applying its value to arguments of the specified sizes and that $\mathsf{Pot}(t)$ is an upper bound on the length of e's value.

With Pot's definition in hand, we make good on the promise to check that the notions defined between Definitions 49 and 60 make sense.

Lemma 61. Suppose σ and $\sigma \rightarrow \tau$ are ATR types.

(a) For each $v \in \mathcal{V}_{wt}[\sigma]$, $||v|| \in \mathcal{T}[\sigma]$ and $\mathsf{Pot}(v) = |v|$.

(b) For each $t_0 \in \mathcal{T}\llbracket \sigma \to \tau \rrbracket$ and $t_1 \in \mathcal{T}\llbracket \sigma \rrbracket$, $t_0 \star t_1 \in \mathcal{T}\llbracket \tau \rrbracket$.

(c) Λ_{\star} is well-defined in the sense that the left-hand side of (12.2) is in $\mathcal{T}[\![\sigma_1 \to \tau]\!]$ as asserted in Definition 61.

All three parts follow straightforwardly from the definitions.

Time-complexity polynomials. To complete the basic time-complexity framework, we define an extension of the second-order polynomials for the simple product types over T, T_{ε} , T_{\circ} ,... under the \mathcal{L} -semantics. The restriction of these to the time types under the \mathcal{L}_{wt} -semantics are the *time-complexity polynomials*. First we extend the grammar for raw

$$\begin{split} \mathcal{T}\llbracket k \rrbracket \varrho & \stackrel{\text{def}}{=} \|k\| & \mathcal{T}\llbracket (\mathbf{c_a} \ e_0) \rrbracket \varrho & \stackrel{\text{def}}{=} (c_0 + \underline{2}, \ p_0 + \underline{1}). \\ \mathcal{T}\llbracket (\mathbf{t_a} \ e_0) \rrbracket \varrho & \stackrel{\text{def}}{=} (c_0 + \underline{2}, \underline{1}). & \mathcal{T}\llbracket (\mathbf{d} \ e_0) \rrbracket \varrho & \stackrel{\text{def}}{=} (c_0 + \underline{2}, \ (p_0 - 1) \lor \underline{0}). \\ \mathcal{T}\llbracket v \rrbracket \varrho & \stackrel{\text{def}}{=} \varrho(v). & \mathcal{T}\llbracket (\operatorname{down} \ e_0 \ e_1) \rrbracket \varrho & \stackrel{\text{def}}{=} (c_0 + c_1 + p_0 + p_1 + \underline{3}, \ \min(p_0, p_1)) \\ \mathcal{T}\llbracket (\lambda x \cdot e_0) \rrbracket \varrho & \stackrel{\text{def}}{=} \Lambda_{\star}(x, \mathcal{T}\llbracket e_0 \rrbracket) \varrho. & \mathcal{T}\llbracket (e_0 \ e_1) \rrbracket \varrho & \stackrel{\text{def}}{=} (\mathcal{T}\llbracket e_0 \rrbracket \varrho) \star (\mathcal{T}\llbracket e_1 \rrbracket \varrho). \end{split}$$

 $\mathcal{T}\llbracket(\mathsf{if}\ e_0\ \mathsf{then}\ e_1\ \mathsf{else}\ e_2)\rrbracket\varrho\ \stackrel{\mathrm{def}}{=}\ (c_0+\underline{2},\underline{0})+(c_1,p_1)\vee(c_2,p_2).$

Above: k is a string constant, $\varrho \in \mathcal{T}[[\Gamma; \Delta]]$, and $(c_i, p_i) = \mathcal{T}[[\Gamma; \Delta \vdash e_i; \gamma_i]] \varrho$ for i = 0, 1, 2.

Figure 22: The
$$\mathcal{T}$$
-interpretation of ATR⁻.

expressions to include: $P ::= (P, P) | \pi_1(P) | \pi_2(P)$. Then we add the following new typing rules for second-order polynomials:

$$\frac{\Sigma \vdash p: \sigma_1 \times \sigma_2}{\Sigma \vdash \pi_i(p): \sigma_i} \qquad \frac{\Sigma_1 \vdash p_1: \sigma_1 \quad \Sigma_2 \vdash p_2: \sigma_2}{\Sigma_1 \cup \Sigma_2 \vdash (p_1, p_2): \sigma_1 \times \sigma_2} \qquad \frac{\Sigma_1 \vdash p_1: \sigma \quad \Sigma_2 \vdash p_2: \sigma}{\Sigma_1 \cup \Sigma_2 \vdash p_1 \odot p_2: \sigma}$$

where σ , σ_1 , and σ_2 simple product types over T , T_{ε} , T_{\diamond} ,... and \odot stands for any of *, +, or \lor . Next we extend the arithmetic operations to all types by recursively defining, for each γ and each $u, v \in \mathcal{L}[\![\gamma]\!]$:

$$u \odot v \stackrel{\text{def}}{=} \begin{cases} \text{the standard thing,} & \text{if } \gamma = \mathsf{T}; \\ (\pi_1(u) \odot \pi_1(v), \pi_2(u) \odot \pi_2(v)), & \text{if } \gamma = \sigma \times \tau; \\ \lambda z \in \mathcal{L}\llbracket \sigma \rrbracket (u(z) \odot v(z)), & \text{if } \gamma = \sigma \to \tau. \end{cases}$$
(12.4)

Finally, the \mathcal{L} -interpretation of the polynomials is just the standard definition.

Remark 62. Note that q_1 of Lemma 53 and the right-hand sides of (12.1) and (12.3) are well-typed, time-complexity polynomials. Also note that by Definition 54(a), if q_1 and q_2 are time-complexity polynomials with $\|\Gamma; \Delta\| \vdash q_1: \|\sigma \to \tau\|$ and $\|\Gamma; \Delta\| \vdash q_2: \|\sigma\|$, then $q_1 \star q_2$ is a time-complexity polynomial with $\|\Gamma; \Delta\| \vdash q_1 \star q_2: \|\tau\|$.

13. The time-complexity interpretation of ATR⁻

Here we establish a polynomial time-boundedness result for ATR^- , the subsystem of ATR obtained by dropping the crec construct. Definition 63 introduces the \mathcal{T} -interpretation of ATR^- and the proof of Theorem 67 shows that ATR^- -expressions have time complexities that are polynomial bounded and well-behaved in other ways. All of this turns out to be pleasantly straightforward. The hard work comes in the following two sections: §14 establishes a key time-complexity decomposition property concerning the affine types and §15 uses this decomposition to define the \mathcal{T} -interpretation of crec expressions and to prove a polynomial boundedness theorem for ATR time complexities.

Convention: Through out this section suppose that γ , σ , and τ are ATR types and Γ ; Δ is an ATR type context.

Definition 63. Figure 22 provides the \mathcal{T} -interpretation for each ATR⁻ construct.

We note that our \mathcal{T} -interpretation of ATR^- is well-defined in the sense that $\mathcal{T}\llbracket\Gamma; \Delta \vdash e: \gamma \rrbracket \varrho \in \mathcal{T}\llbracket\gamma \rrbracket$ for each ATR^- judgment $\Gamma; \Delta \vdash e: \gamma$ and $\varrho \in \mathcal{T}\llbracket\Gamma; \Delta \rrbracket$. (This follows from Lemma 61 and some straightforward calculations.) Here is a simple application of Definition 63. Let $\mathbf{g} = (\lambda y \cdot (\mathbf{c_0} \ (\mathbf{c_0} \ y))): \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond}$ and $\mathsf{A} = (\lambda f \cdot (\lambda x \cdot (f \ x))): (\mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond}) \to \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond}$. We write $\mathcal{T}\llbrackete\rrbracket$ for $\mathcal{T}\llbrackete\rrbracket$ to cut some clutter. The reader may check that:

$$\begin{aligned} \mathcal{T}\llbracket \mathbf{g} \rrbracket &= (\underline{1}, \lambda p_y \in \mathcal{P}\llbracket \mathsf{N}_{\varepsilon} \rrbracket \cdot (\underline{1} \lor p_y + \underline{4}, p_y + \underline{2})). \\ \mathcal{T}\llbracket \mathsf{A} \rrbracket &= (\underline{1}, \lambda p_f \in \mathcal{P}\llbracket \mathsf{N}_{\varepsilon} \to \mathsf{N}_{\diamond} \rrbracket \cdot (\underline{1}, \lambda p_x \in \mathcal{P}\llbracket \mathsf{N}_{\varepsilon} \rrbracket \cdot \mathsf{val}(p_f) \star \mathsf{val}(p_x))). \\ \mathcal{T}\llbracket (\mathsf{A} \mathsf{g}) \rrbracket &= (\mathcal{T}\llbracket \mathsf{A} \rrbracket) \star (\mathcal{T}\llbracket \mathsf{g} \rrbracket) &= dally(7, \mathcal{T}\llbracket \mathsf{g} \rrbracket). \end{aligned}$$

There are three key things to establish about the time complexities assigned by \mathcal{T} , that they are: not too big, not too small, and well-behaved. "Not too big" means that the time complexities are polynomially-bounded in the sense of Definition 64 below. "Not too small" means that $\operatorname{cost}_{\operatorname{CEK}}(e,\rho) \leq \operatorname{cost}(\mathcal{T}[\![e]\!] \|\rho\|)$ and $|\mathcal{V}_{\mathrm{wt}}[\![e]\!] \rho| \leq \operatorname{Pot}(\mathcal{T}[\![e]\!] \|\rho\|)$. This "not too small" property (soundness) is introduced in Definition 65. Finally, "well-behaved" means that the \mathcal{T} -assigned time complexities are monotone (Definition 66) which requires that $\mathcal{T}[\![e]\!] \rho \leq \mathcal{T}[\![e]\!] \rho'$ when $\rho \leq \rho'$ (see Definition 66(a)) and that when $\mathcal{T}[\![e]\!] \rho$ is a function, it is pointwise, monotone nondecreasing. Monotonicity plays an important role in dealing with crec. Theorem 67 establishes that the \mathcal{T} -interpretation of ATR⁻ satisfies each of these properties. Let \mathcal{F} range over programming formalisms (e.g., ATR⁻ or ATR) in the following.

Definition 64 (Polynomial time-boundedness). A \mathcal{T} -interpretation of \mathcal{F} is polynomial timebounded when, given $\Gamma; \Delta \vdash_{\mathcal{F}} e; \gamma$, we can effectively find a time-complexity polynomial p_e with $|\Gamma; \Delta| \vdash p_e: ||\gamma||$ such that $\mathcal{T}[\![e]\!] ||\rho|| \leq \mathcal{L}_{wt}[\![p_e]\!] |\rho||$ for each $\rho \in \mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$.

Definition 65 (Soundness). A \mathcal{T} -interpretation of \mathcal{F} is sound when, for each $\Gamma; \Delta \vdash_{\mathcal{F}} e: \gamma$ and each $\rho \in \mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$, we have $\operatorname{cost}_{CEK}(e, \rho) \leq \operatorname{cost}(\mathcal{T}[\![e]\!] \|\rho\|)$ and $|\mathcal{V}_{wt}[\![e]\!] \rho| \leq \operatorname{Pot}(\mathcal{T}[\![e]\!] \|\rho\|)$.

Definition 66 (Monotonicity).

(a) For $\rho, \rho' \in \mathcal{T}[\Gamma; \Delta]$, we write $\rho \leq \rho'$ when $\rho(x) \leq \rho'(x)$ for each $x \in \text{preimage}(\Gamma; \Delta)$.

(b) We say that a \mathcal{T} -interpretation of \mathcal{F} is *monotone* when, for each $\Gamma; \Delta \vdash_{\mathcal{F}} e: \gamma$: (i) $\mathcal{T}\llbracket e \rrbracket$ is a pointwise, monotone nondecreasing function from $\mathcal{T}\llbracket \Gamma; \Delta \rrbracket$ to $\mathcal{T}\llbracket \gamma \rrbracket$, and (ii) if $\gamma = (\sigma_0, \ldots, \sigma_k) \to \mathbf{b}$, then the function from $\mathcal{T}\llbracket \Gamma; \Delta \rrbracket \times \mathcal{T}\llbracket \sigma_0 \rrbracket \times \cdots \times \mathcal{T}\llbracket \sigma_k \rrbracket$ to $\mathcal{T}\llbracket \mathbf{b} \rrbracket$ given by $(\varrho, v_0, \ldots, v_k) \mapsto ((\mathcal{T}\llbracket e \rrbracket \varrho) v_0 \ldots v_k)$ is pointwise, monotone nondecreasing.

Theorem 67. The \mathcal{T} -interpretation of ATR^- is (a) polynomial time-bounded, (b) monotone, and (c) sound.

The proofs of parts (a) and (b) are straightforward standard structural inductions, but the argument for (c) is a logical-relations arguments [Win93]. Before proving the above we first introduce a few useful time-complexity polynomials.

Definition 68. N.B. The following definitions are purely syntactic. Suppose $v \in \mathbb{N}$. Let $||v|| \stackrel{\text{def}}{=} (\underline{1} \lor |v|, |v|)$. For each \mathbf{b} , let $||x||_{\mathbf{b}} \stackrel{\text{def}}{=} (\underline{1} \lor |x|, |x|)$. For each $\gamma = (\mathbf{b}_1, \ldots, \mathbf{b}_k) \to \mathbf{b}_0$, let $||x||_{\gamma} \stackrel{\text{def}}{=} (\underline{1}, q_1)$ where $q_i = \lambda p_i \cdot (\underline{1}, q_{i+1})$, for each i with $1 \le i < k$, and $q_k = \lambda p_k \cdot (\underline{1} \lor |x|(p_1, \ldots, p_k), |x|(p_1, \ldots, p_k))$. (Recall Lemma 53.)

Note that if $\Gamma; \Delta \vdash x; \gamma$ where x is a variable, then $|\Gamma; \Delta| \vdash ||x|| : ||\gamma||$.

Proof of Theorem 67(a): Polynomial time-boundedness. Fix an ATR⁻-judgment $\Gamma; \Delta \vdash e; \gamma$. Let ρ range over $\mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$. We have to effectively construct a t.c. polynomial q_e as required by Definition 64. The argument is yet another a structural induction on the derivation of $\Gamma; \Delta \vdash e; \gamma$. We consider the cases of the last rule used in the derivation.

CASES: Zero-I and Const-I. Then $e = v \in \mathbb{N}$ and γ is a base type. Let $q_e = ||v||$. By Definition 63, $\mathcal{T}[v] ||\rho|| = (\underline{1} \vee |v|, |v|) = \mathcal{L}_{wt}[\![q_e]\!] |\rho|$ and thus q_e suffices.

CASES: Int-Id-I and Aff-Id-I. Then e = x, a variable. Then by Definition 63, $\mathcal{T}[x] \|\rho\| = \|\rho\|(x) = \|\rho(x)\|$. Let $q_e = \|x\|_{\gamma}$. Subcase: γ is a base type. By Definition 68(b), $q_e = (\underline{1} \lor |x|, |x|)$. So by Definition 51, $\|\rho(x)\| = \mathcal{L}_{wt}[q_e] |\rho|$ and thus q_e suffices. Subcase: $\gamma = (\mathbf{b}_1, \ldots, \mathbf{b}_k) \to \mathbf{b}_0$. By Definition 68(c), $q_e = (\underline{1}, q_1)$, where q_1, \ldots, q_k are as in that definition. By Lemma 53, $\|\rho(x)\| = \mathcal{L}_{wt}[q_e] |\rho|$ and thus q_e suffices.

CASE: $\mathbf{c_a}$ -I, where $\mathbf{a} \in \{\mathbf{0}, \mathbf{1}\}$. Then $e = (\mathbf{c_a} \ e_0)$ for some e_0 and $\gamma = \mathsf{N}_{\diamond_d}$ for some d. Let $(c_0, p_0) = \mathcal{T}\llbracket e_0 \rrbracket \|\rho\|$. By Definition 63, $\mathcal{T}\llbracket e \rrbracket \|\rho\| = (c_0 + 2, p_0 + 1)$. By the induction hypothesis, we can construct q_{e_0} with $|\Gamma; \Delta| \vdash q_{e_0}: \mathsf{N}_{\diamond_d}$ such that $\mathcal{T}\llbracket e_0 \rrbracket \|\rho\| \leq \mathcal{L}_{\mathrm{wt}}\llbracket q_{e_0} \rrbracket |\rho|$. Thus, $q_e = q_{e_0} + (2, 1)$ suffices.

CASES: t_0 -*I*, t_1 -*I*, down-*I*, d-*I*, \rightarrow -*e*, and *If*-*I*. These follow by arguments analogous to the proof for the c_a -*I* case.

CASES: Subsumption and Shifting. There is nothing to prove here.

CASE: $\rightarrow -E$. Then $e = (e_0 \ e_1)$ for some e_0 and e_1 with $\Gamma; \Delta \vdash e_0: \tau \rightarrow \gamma$ and $\Gamma; _ \vdash e_1: \tau$. By the induction hypothesis, we can construct q_0 and q_1 , bounding time-complexity polynomials for $\mathcal{T}\llbracket e_0 \rrbracket$ and $\mathcal{T}\llbracket e_1 \rrbracket$, respectively. Let $q_e = q_0 \star q_1$. By Remark 62, q_e is a time-complexity polynomial and it follows from the monotonicity of \star that $\mathcal{T}\llbracket (e_0 \ e_1) \rrbracket \parallel \rho \parallel = (\mathcal{T}\llbracket e_0 \rrbracket \parallel \rho \parallel) \star (\mathcal{T}\llbracket e_1 \rrbracket \parallel \rho \parallel) \leq (q_0 \ \parallel \rho \parallel) \star (q_1 \ \parallel \rho \parallel) = q_e \ \parallel \rho \parallel$. Thus, q_e suffices.

CASE: $\rightarrow -I$. Then $\gamma = \sigma \rightarrow \tau$ and $e = (\lambda x \cdot e_0)$ for some e_0 with $\Gamma, x: \sigma; \Delta \vdash e_0; \tau$. By Definitions 57 and 63 we thus have $\mathcal{T}\llbracket e \rrbracket = \Lambda_*(x, \mathcal{T}\llbracket e_0 \rrbracket)$. By the induction hypothesis, we can construct q_{e_0} with $|\Gamma, x: \sigma; \Delta| \vdash q_{e_0}: ||\tau||$ with $\mathcal{T}\llbracket e_0 \rrbracket ||\rho'|| \leq \mathcal{L}_{wt}\llbracket q_{e_0} \rrbracket ||\rho'|$ for each $\rho' \in \mathcal{V}_{wt}\llbracket \Gamma, x: \sigma; \Delta \rrbracket$. Subcase: σ is a base type. So, $\langle \langle \sigma \rangle \rangle = \mathsf{T} \times \sigma$. Let $q_e = (\underline{1}, \lambda | x | \cdot q_{e_0})$. A straightforward argument shows that q_e suffices for the polynomial bound. Subcase: $\sigma = (\sigma_1, \ldots, \sigma_k) \rightarrow \mathbf{b}$. Let p' be the expression $\lambda ||p|| \cdot \pi_2((\underline{1}, p) \star ||p||)$, where $|p|| = |y_1|, \ldots, |y_k|$ and $||p|| = (\underline{1} \lor |y_1|, |y_1|), \ldots, (\underline{1} \lor |y_k|, |y_k|)$. (See Definition 60(b).) Let p'' be the expansion of p' in which p is treated as being of type $\langle \langle \sigma \rangle \rangle$ and the \mathcal{T} -applications are expanded out per Definition 54. It follows that p'' is a time complexity polynomial with $|\Gamma|, p: \langle \langle \sigma \rangle \rangle; |\Delta| \vdash p'': |\sigma|$. Let $q_e = (\underline{1}, \lambda p \cdot q_{e_0}[|x| := p''])$. Again, a straightforward argument shows that q_e suffices for the polynomial of p' in which p is polynomial bound.

Theorem 67(a)

Proof of Theorem 67(b): Monotonicity. This argument follows along the lines of the proof of part (a) and is left to the reader. \Box

Theorem 67(b)

For the proof of soundness, we shall first define a logical relation $\sqsubseteq_{\gamma}^{\text{tc}}$ between CEKclosures and time-complexities. Roughly, $e\hat{\rho} \sqsubseteq_{\gamma}^{\text{tc}}(c,p)$ says that the time complexity (c,p)bounds the cost of evaluating the closure $e\hat{\rho}$. Conventions on CEK-closures: CEK-closures are written $e\hat{\rho}$. (We always assume $FV(e) \subseteq \text{preimage}(\hat{\rho})$.) A CEK-closure $e\hat{\rho}$ is called a value when e is a CEK-value. $e\hat{\rho} \downarrow v\hat{\rho}'$ means that starting from $(e, \hat{\rho}, \mathsf{halt})$, the CEKmachine eventually ends up with $(v, \hat{\rho}', \mathsf{halt})$, where $v\hat{\rho}'$ is a value. Below, v ranges over CEK-values and p and q range over potentials.

Definition 69.

(a) For each ATR-type γ we define a relation $\sqsubseteq_{\gamma}^{\text{tc}}$ between type- γ CEK-closures and time complexities and a second relation $\sqsubseteq_{\gamma}^{\text{pot}}$ between type- γ CEK-closures and potentials as follows.

- $e\hat{\rho} \sqsubseteq_{\gamma}^{\mathrm{tc}}(c,p) \equiv_{\mathrm{def}} \mathrm{cost}_{\mathrm{CEK}}(e,\hat{\rho}) \leq c \& v\hat{\rho}' \sqsubseteq_{\gamma}^{\mathrm{pot}} p$, where $e\hat{\rho} \downarrow v\hat{\rho}'$.
- $v\hat{\rho} \sqsubseteq_{\mathbf{b}}^{\text{pot}} p \equiv_{\text{def}} |v\hat{\rho}| \le p.$
- $(\lambda x \cdot e)\hat{\rho} \sqsubseteq_{\sigma \to \tau}^{\text{pot}} p \equiv_{\text{def}} \text{ for all } v\hat{\rho}' \text{ and all } q \text{ with } v\hat{\rho}' \sqsubseteq_{\sigma}^{\text{pot}} q, \ e(\hat{\rho}[x \mapsto v\hat{\rho}']) \sqsubseteq_{\tau}^{\text{tc}} p(q).$ $O\hat{\rho} \sqsubseteq_{\gamma \to \tau}^{\text{pot}} p \equiv_{\text{def}} \text{ for all } v\hat{\rho}' \text{ and all } q \text{ with } v\hat{\rho}' \sqsubseteq_{\gamma}^{\text{pot}} q, \ O(v\hat{\rho}')\{\} \sqsubseteq_{\tau}^{\text{tc}} p(q).$

(b) Suppose $\rho \in \mathcal{T}\llbracket\Gamma; \Delta\rrbracket$. We write $\hat{\rho} \sqsubseteq \rho$ when, for each $x \in \text{preimage}(\hat{\rho}), x\hat{\rho} \sqsubseteq_{\gamma}^{\text{tc}} \rho(x),$ where $\gamma = (\Gamma; \Delta)(x)$.

(c) Suppose $\Gamma; \Delta \vdash e: \gamma$ and $X: \mathcal{T}\llbracket\Gamma; \Delta\rrbracket \to \mathcal{T}\llbracket\gamma\rrbracket$. We write $e \sqsubseteq_{\gamma}^{\mathrm{tc}} X$ when, for all CEK-environments $\hat{\rho}$ and $\rho \in \mathcal{T}\llbracket\Gamma; \Delta\rrbracket$ with $\hat{\rho} \sqsubseteq \rho$, $e\hat{\rho} \sqsubseteq_{\gamma}^{\mathrm{tc}} X \rho$.

Lemma 70.

(a) Suppose x is a variable and $v\hat{\rho} \sqsubseteq_{\gamma}^{\text{pot}} q$. Then $x(\hat{\rho}' \cup \{x \mapsto v\hat{\rho}\}) \sqsubseteq_{\gamma}^{\text{tc}} \mathsf{val}(q)$.

(b) Suppose $\Gamma; \Delta \vdash e: \gamma$. Then $e \sqsubseteq_{\gamma}^{tc} \mathcal{T}\llbracket e \rrbracket$.

(c) Suppose $e\hat{\rho}$ is a type- γ CEK-closure and t and t' are type- γ time complexities with $e\hat{\rho} \sqsubseteq_{\gamma}^{\mathrm{tc}} t \text{ and } t \leq t'.$ Then $e\hat{\rho} \sqsubseteq_{\gamma}^{\mathrm{tc}} t'.$

Proof. Part (a). Since $v\hat{\rho} \sqsubseteq_{\gamma}^{\text{pot}} q = pot(\mathsf{val}(q))$, we just need to show that $\operatorname{cost}_{\operatorname{CEK}}(x(\hat{\rho}' \cup x))$ $\{x \mapsto v\hat{\rho}\}\} \leq cost(val(q))$. If γ is a base type, then $|v| \leq q$, hence $cost_{CEK}(x(\hat{\rho}' \cup \{x \mapsto v\hat{\rho}\}))$ $v\hat{\rho}$)) = $\underline{1} \lor |v| \le \underline{1} \lor q = cost(val(q)).$

Part (b). The argument is a structural induction on the derivation of $\Gamma; \Delta \vdash e; \gamma$. We consider the cases of the last rule used in the derivation. Fix a CEK-environment $\hat{\rho}$ and a $\varrho \in \mathcal{T}\llbracket\Gamma; \Delta\rrbracket$ with $\hat{\rho} \sqsubseteq \varrho$.

CASE: Zero-I and Const-I. Then e = v, a string constant. So, $\mathcal{T}[e] \rho = (1 \lor |v|, |v|)$, $\operatorname{cost}_{\operatorname{CEK}}(e, \hat{\rho}) = \underline{1} \leq \underline{1} \vee |v|$, and $|e \rho| = |v|$. Hence, e is as required.

CASE: Int-Id-I and Aff-Id-I. Then e = x, a variable. Since $\hat{\rho} \sqsubseteq \varrho$, we have $x\hat{\rho} \sqsubseteq_{\gamma}^{\text{tc}}$ $\rho(x) = \mathcal{T}[x] \rho$. Hence, *e* is as required.

CASE: $c_{\mathbf{a}}$ -*I*, where $\mathbf{a} \in \{0, 1\}$. Then $e = (c_{\mathbf{a}} e_0)$ where $\Gamma; \Delta \vdash e_0; \gamma$ and γ is a base type. Let $(c_0, p_0) = \mathcal{T}\llbracket e_0 \rrbracket \varrho$ and suppose $e_0 \hat{\rho} \downarrow v \hat{\rho}'$. By the induction hypothesis applied to e_0 , we know $\operatorname{cost}_{\operatorname{CEK}}(e_0, \hat{\rho}) \leq c_0$ and $|v\hat{\rho}'| \leq p_0$. By inspection of the CEK machine and the definition of $\operatorname{cost}_{\operatorname{CEK}}$, $\operatorname{cost}_{\operatorname{CEK}}(\mathsf{c}_{\mathbf{a}} e_0, \hat{\rho}) = \operatorname{cost}_{\operatorname{CEK}}(e_0, \hat{\rho}) + \underline{2} \leq c_0 + \underline{2}$. It also follows that $(\mathbf{c}_{\mathbf{a}} \ e_0)\hat{\rho} \downarrow (\mathbf{a} \oplus v)\hat{\rho}'$ and $|(\mathbf{a} \oplus v)\hat{\rho}'| = |v\hat{\rho}'| + \underline{1} \leq p_0 + \underline{1}$. By Definition 63, $\mathcal{T}\llbracket e \rrbracket \parallel \rho \parallel = (c_0 + \underline{2}, p_0 + \underline{1}).$ Hence, e is as required.

CASES: t_0-I , t_1-I , down-I, d-I, $\rightarrow -E$, and If-I. These follow by arguments analogous to the proof for the c_a -*I* case.

CASES: Subsumption and Shifting. There is nothing to prove here.

CASE: $\rightarrow I$. Then $\gamma = \sigma \rightarrow \tau$ and $e = (\lambda x \cdot e_0)$ for some e_0 with $\Gamma, x: \sigma; \Delta \vdash e_0: \tau$. So, by Definition 63, $cost(\mathcal{T}[\lambda x \cdot e_0]] \varrho) = \underline{1}$ and $\lambda x \cdot e_0 \varrho$ is itself a value. Since $cost_{CEK}(\lambda x \cdot e_0, \hat{\rho}) = \underline{1}$ 1, all that is left to show is that $(\lambda x \cdot e_0)\hat{\rho} \sqsubseteq_{\sigma \to \tau}^{\text{pot}} \operatorname{pot}(\mathcal{T}\llbracket\lambda x \cdot e_0] \varrho)$. Let $p = \operatorname{pot}(\mathcal{T}\llbracket\lambda x \cdot e_0] \varrho) =$ $pot(\Lambda_{\star}(x, \mathcal{T}\llbracket e_0 \rrbracket) \varrho) = \lambda p' \in \mathcal{P}\llbracket \sigma \rrbracket (\mathcal{T}\llbracket e_0 \rrbracket (\varrho' \cup \{x \mapsto \mathsf{val}(p')\})), \text{ let } v\hat{\rho} \text{ be an arbitrary type-}\sigma$ value and let q be an arbitrary potential with $v\hat{\rho} \sqsubseteq_{\sigma}^{\text{pot}} q$. Then establishing $(\lambda x \cdot e_0)\hat{\rho} \sqsubseteq_{\sigma \to \tau}^{\text{pot}}$ $\mathcal{T}[\lambda x \cdot e_0] \rho$ is equivalent to showing $e_0(\hat{\rho}[x \mapsto v\hat{\rho}']) \sqsubseteq_{\tau}^{\mathrm{tc}} p(q)$ By part (a), $x(\hat{\rho} \cup \{x \mapsto v\rho\} \sqsubseteq_{\sigma}^{\mathrm{tc}})$ $\mathsf{val}(q)$. Hence, $\hat{\rho} \cup \{x \mapsto v\rho\} \sqsubseteq \varrho \cup \{x \mapsto \mathsf{val}(q)\}$. Thus, by the induction hypothesis on e_0 , $e_0(\hat{\rho}[x \mapsto v\hat{\rho}']) \sqsubseteq_{\tau}^{\text{tc}} \mathcal{T}\llbracket e_0 \rrbracket (\varrho' \cup \{x \mapsto \mathsf{val}(q)\}) = p(q).$ Hence, e is as required.

CASE: $\rightarrow -E$. Then $e = (e_0 \ e_1)$ for some e_0 and e_1 with $\Gamma; \Delta \vdash e_0: \sigma \rightarrow \gamma$ and $\Gamma; _\vdash e_1: \sigma$. Suppose $e_0\hat{\rho} \downarrow v_0\hat{\rho}_0$, $e_1\hat{\rho} \downarrow v_1\hat{\rho}_1$, $(e_0 \ e_1)\hat{\rho} \downarrow v_r\hat{\rho}_r$, $(c_0, p_0) = \mathcal{T}[\![e_0]\!]\varrho$, $(c_1, p_1) = \mathcal{T}[\![e_1]\!]\varrho$, and $(c_r, p_r) = p_0(p_1)$. By the induction hypothesis on e_0 and e_1 :

(a)
$$\operatorname{cost}_{\operatorname{CEK}}(e_0, \hat{\rho}) \le c_0.$$
 (b) $v_0 \hat{\rho}_0 \sqsubseteq_{\sigma}^{\operatorname{pot}} p_0.$ (13.1)

(a)
$$\operatorname{cost}_{\operatorname{CEK}}(e_1, \hat{\rho}) \le c_1.$$
 (b) $v_1 \hat{\rho}_1 \sqsubseteq_{\sigma}^{\operatorname{pot}} p_1.$ (13.2)

There are two subcases to consider based on the form of v_0 . Subcase: $v_0 = \lambda x \cdot e'_0$ for some $\Gamma, x: \sigma; \Delta \vdash e'_0; \gamma$. Then (13.1b) means that, for all type- τ values $v\hat{\rho}''$ and all q with $v\hat{\rho}'' \sqsubseteq_{\sigma}^{\text{pot}} q$, we have $e'_0(\hat{\rho}' \cup \{x \mapsto v\hat{\rho}''\}) \sqsubseteq_{\gamma}^{\text{tc}} p_0(q)$. So by (13.2b), $e'_0\hat{\rho}'_0 \sqsubseteq_{\gamma}^{\text{tc}} (c_r, p_r)$, where $\hat{\rho}'_0 = \hat{\rho}' \cup \{x \mapsto v_1\hat{\rho}_1\}$. Now

$$\begin{aligned}
\cot_{CEK}((e_0 \ e_1), \hat{\rho}) &= \ \cot_{CEK}(e_0, \hat{\rho}) + \cot_{CEK}(e_1, \hat{\rho}) + \cot_{CEK}(e'_0, \hat{\rho}'_0) + \underline{3} \\
& (by \ Figure \ 20 \ \& \ Definition \ 48) \\
&\leq \ c_0 + c_1 + c_r + \underline{3} \\
&= \ \cot(\mathcal{T}[(e_0 \ e_1)]] \varrho) \quad (by \ Definition \ 63).
\end{aligned}$$

Note that $e'_0 \hat{\rho}'_0 \downarrow v_r \hat{\rho}_r$. So by $e'_0 \hat{\rho}'_0 \sqsubseteq_{\gamma}^{\text{tc}} (c_r, p_r)$, $v_r \hat{\rho}_r \sqsubseteq_{\gamma}^{\text{pot}} p_r = pot(\mathcal{T}[\![(e_0 e_1)]\!] \varrho)$. Hence, in this subcase e is as required. Subcase: v_0 is an oracle. The argument here is a repeat, mutatis mutandis, of the proof of previous subcase.

Part (c). The argument follows along the lines of the proof of (b).

Lemma 70

Proof of Theorem 67(c): Soundness. This follows straightforwardly from Lemma 70(b) and Definition 60. \Box

Scholium 71. The \mathcal{T} -interpretation of ATR⁻ (and later, ATR) sits in-between the actual costs of evaluating expressions on our CEK machine and the sought-after polynomial time-bounds on these costs. Why is working with \mathcal{T} -interpretations preferable to working directly with executions of CEK machines and their costs? Part of the reason is that \mathcal{T} -interpretations have built-in to them the cost-potential aspects expressions. One would somehow have to replicate these in working directly with CEK-computations. Another part of the reason is that \mathcal{T} -interpretations collapse the many possible paths of a CEK-computation into a single time-complexity. The \mathcal{T} -interpretation of if-then-else is chiefly responsible for these collapses. Scholium 80 notes that these collapses are a source of some trouble in dealing with crec-expressions.

14. An Affine decomposition of time complexities

When analyzing the time complexity of a program, one often needs to decompose its time complexity into pieces that may have little to do with the program's apparent syntactic structure. Theorem 74 below is a general time-complexity decomposition result for ATR expressions. The ATR typing rules for affinely restricted variables are critical in ensuring this time-complexity decomposition. The decomposition is used in the next section to obtain the recurrences for the analysis of the time complexity of **crec** expressions. Note that the theorem presupposes that that $\mathcal{T}[\cdot]$ is defined on **crec** expressions. However, since no affinely restricted variable can occur free in a well-typed **crec** expression and since the application of the theorem will be within a structural induction, this presupposition does not add any difficulties.

Remark 72. In fact, the time-complexity of a crec expression e will be defined in terms of time-complexities of expressions built up from subexpressions of e using term constructors other than crec. Thus a completely standard structural induction for establishing soundness does not quite work. A fully formal proof would have first established results such as "if $e_0 \sqsubseteq_{\sigma \to \tau}^{\text{tc}} X_0$ and $e_1 \sqsubseteq_{\sigma}^{\text{tc}} X_1$, then $(e_0 e_1) \sqsubseteq_{\tau}^{\text{tc}} X_0 \star X_1$ " where the X_i 's are general mappings from \mathcal{T} -environments to time complexities. These lemmas would then be used to carry out the induction steps of a structural induction which, in all but the crec case, would just quote the relevant lemma. Rather than impose this additional level of detail on the reader, we have opted for a less formal approach here and will assume that if we inductively have soundness for a subterm e, then we also have it for terms built up from e without crec.

To help in the statement and proof of the Affine Decomposition Theorem, we introduce the following definitions and conventions.

Definition 73.

(a) $(c_1, p_1) \uplus (c_2, p_2) \stackrel{\text{def}}{=} (c_1 + c_2, p_1 \lor p_2)$, where $(c_1, p_1), (c_2, p_2) \in \mathcal{T}[\![\gamma]\!]$. (Clearly, $(c_1, p_1) \uplus (c_2, p_{h2}) \in \mathcal{T}[\![\gamma]\!]$.)

(b) For each ATR-type γ , define ϵ_{γ} inductively by: $\epsilon_{N_{\ell}} = \epsilon$ and $\epsilon_{\sigma \to \tau} = \lambda x \cdot \epsilon_{\tau}$. (Clearly, $\vdash \epsilon_{\gamma} : \gamma \text{ and } |\mathcal{V}_{wt}[\![\epsilon_{\gamma}]\!] \{\} = \underline{0}_{|\gamma|}$.)

(c) Given $f: (\sigma_1, \ldots, \sigma_k) \to \mathsf{N}_{\ell}$, an expression of the form $(f \ e_1 \ \ldots \ e_k)$ is called a *full* application of f.

Conventions on factoring out environments: Suppose \odot is a binary operation on time complexities. We often write $\mathcal{T}\llbracket e_0 \rrbracket \odot \mathcal{T}\llbracket e_1 \rrbracket$ for $\varrho \mapsto (\mathcal{T}\llbracket e_0 \rrbracket \varrho) \odot (\mathcal{T}\llbracket e_1 \rrbracket) \varrho)$. For example: $(\mathcal{T}\llbracket e_0 \rrbracket \uplus \mathcal{T}\llbracket e_1 \rrbracket) \varrho = (\mathcal{T}\llbracket e_0 \rrbracket \varrho) \uplus (\mathcal{T}\llbracket e_1 \rrbracket) \varrho)$ and $(\mathcal{T}\llbracket e_0 \rrbracket \star \mathcal{T}\llbracket e_1 \rrbracket) \varrho = (\mathcal{T}\llbracket e_0 \rrbracket \varrho) \star (\mathcal{T}\llbracket e_1 \rrbracket) \varrho)$. We extend this convention to *n*-ary operations. For example: $\mathsf{val}(\mathcal{T}\llbracket e_1 \rrbracket) \varrho = \mathsf{val}(\mathcal{T}\llbracket e_1 \rrbracket) \varrho$ and $(\mathcal{T}\llbracket e_0 \rrbracket \star \ldots \star \mathcal{T}\llbracket e_k \rrbracket) \varrho = (\mathcal{T}\llbracket e_0 \rrbracket \varrho) \star \ldots \star (\mathcal{T}\llbracket e_k \rrbracket) \varrho)$. We also generalize this last equality as follows. Suppose X is a map from $\mathcal{T}\llbracket \Gamma; \Delta \rrbracket$ to $\mathcal{T}\llbracket \sigma_i \rrbracket$. Then $X \star \vec{Y}$ denotes the map $\mathcal{T}\llbracket \Gamma; \Delta \rrbracket$ to $\mathcal{T}\llbracket \mathsf{N}_\ell \rrbracket$ given by: $(X \star \vec{Y}) \varrho = (X \varrho) \star (Y_1 \varrho) \star \ldots \star (Y_k \varrho)$.

Theorem 74 (Affine decomposition). Suppose Γ ; $f: \gamma \vdash e: \mathsf{N}_{\ell_0}$, where $\gamma = (\mathsf{N}_{\ell_1}, \ldots, \mathsf{N}_{\ell_k})$ $\rightarrow \mathsf{N}_{\ell_0} \in \mathcal{R}$ and TailPos(f, e). Let ζ denote the substitution $[f:=\epsilon_{\gamma}]$. Then

$$\mathcal{T}\llbracket e \rrbracket \leq \mathcal{T}\llbracket e \zeta \rrbracket \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}), \tag{14.1}$$

where $(f e_1^1 \dots e_k^1), \dots, (f e_1^m \dots e_k^m)$ are the full applications of f occurring in e and $t_j = \bigvee_{i=1}^m \operatorname{val}(\mathcal{T}[\![e_j^i]\!])$ for $j = 1, \dots, k$.

By Lemma 11 we know that there is at most one use of an affinely restricted variable in an expression. In terms of costs, one can thus interpret (14.1) as saying that the cost of evaluating e can be bounded by the sum of: (i) the cost of evaluating $e\zeta$, which includes the all of the costs of e except for the possible application of the value of f to the values of its arguments, and (ii) $cost((\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho)$, which clearly bounds the cost of any such fapplication. In terms of potentials, one can interpret (14.1) as saying that the size of the value of e is bounded by the maximum of (i) the size of the value of $e\zeta$, which covers all the cases where f is not applied, and (ii) $pot((\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho)$, which covers all the cases where f is applied.

If (14.1) solely concerned CEK costs, the above remarks would almost constitute a proof. However, (14.1) is about \mathcal{T} -interpretations of expressions and $\mathcal{T}[\![e]\!]$ is an approximation to

$$\begin{aligned} \mathcal{T}\llbracket A \rrbracket \varrho &= (cost(\mathcal{T}\llbracket A_0 \rrbracket \varrho) + \underline{2}, \underline{0}) \ \uplus \ \bigvee_{i=1}^2 \mathcal{T}\llbracket A_i \rrbracket \varrho \\ &\leq (cost(\mathcal{T}\llbracket A_0 \rrbracket \varrho) + \underline{2}, \underline{0}) \ \uplus \ \bigvee_{i=1}^2 \left(\mathcal{T}\llbracket A_i \zeta \rrbracket \varrho \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho \right) \\ &\leq (cost(\mathcal{T}\llbracket A_0 \rrbracket \varrho) + \underline{2}, \underline{0}) \ \uplus \ \left(\bigvee_{i=1}^2 \mathcal{T}\llbracket A_i \zeta \rrbracket \varrho \right) \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho \\ &\leq \left((cost(\mathcal{T}\llbracket A_0 \zeta \rrbracket \varrho) + \underline{2}, \underline{0}) \ \uplus \ \left(\bigvee_{i=1}^2 \mathcal{T}\llbracket A_i \zeta \rrbracket \varrho \right) \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho \\ &\leq \left((cost(\mathcal{T}\llbracket A_0 \zeta \rrbracket \varrho) + \underline{2}, \underline{0}) \ \uplus \ \bigvee_{i=1}^2 \mathcal{T}\llbracket A_i \zeta \rrbracket \varrho \right) \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho \\ &= \mathcal{T}\llbracket A \zeta \rrbracket \varrho \ \uplus \ (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho. \end{aligned}$$

the true time complexities involved in evaluating e. The theorem asserts that our \mathcal{T} -interpretation of ATR is verisimilar enough to capture this property of time complexities. This later requires a little work.

Proof of Theorem 74. Fix $\rho \in \mathcal{T}[[\Gamma; f: \gamma]]$. Without loss of generality, we assume there are no bound occurrences of f in e. We argue by structural induction that for each A, a subterm of e with $\Gamma; f: \gamma \vdash A: \mathbb{N}_{\ell_0}$, we have

$$\mathcal{T}\llbracket A \rrbracket \varrho \leq \mathcal{T}\llbracket A \zeta \rrbracket \varrho \ \uplus \ \left(\mathcal{T}\llbracket f \rrbracket \star \vec{t} \right) \varrho, \tag{14.2}$$

where the \vec{t} 's are as in the lemma's statement. It follows from TailPos(f, e) that the following three cases are the only ones to consider.

Case 1: f fails to occur in A. Then (14.2) follows immediately.

Case 2: $A = (f e_1 \ldots e_k)$, where $\Gamma; _\vdash e_1: \mathsf{N}_{\ell_1}, \ldots, \Gamma; _\vdash e_k: \mathsf{N}_{\ell_k}$. By the monotonicity of $\mathcal{T}\llbracket f \rrbracket$ and the \mathcal{T} -interpretation of application from Figure 22, it follows that (14.2) holds for A.

Case 3: $A = (\text{if } A_0 \text{ then } A_1 \text{ else } A_2)$ where f occurs in A_1 or A_2 or both. By Definitions 63 and 74(c), $\mathcal{T}\llbracket A \rrbracket \varrho = (cost(\mathcal{T}\llbracket A_0 \rrbracket \varrho) + \underline{2}, \underline{0}) \uplus \bigvee_{i=1}^2 \mathcal{T}\llbracket A_i \rrbracket \varrho$. Note: $A_0 = A_0 \zeta$ since f cannot appear in A_0 . By the induction hypothesis applied to A_1 and A_2 , $\mathcal{T}\llbracket A_i \rrbracket \varrho \leq \mathcal{T}\llbracket A_i \rrbracket \varrho$ $\Downarrow (\mathcal{T}\llbracket f \rrbracket \star \vec{t}) \varrho$ for i = 1, 2. Thus we have the chain of bounds of Figure 23.

Scholium 75. As demonstrated in [DR07], handling forms of recursion beyond tail recursion requires notions of decomposition more sophisticated than (14.1). Moreover, if explicit —o-types were added to ATR, then the decomposition also becomes more involved than (14.1).

For the analysis of **crec** expressions we need the following corollary to Theorem 74. We leave its proof to the reader who should be mindful of Remark 72 above.

Corollary 76. Suppose $\Gamma; f: \gamma \vdash A: \gamma$, where $\gamma = (\mathsf{N}_{\ell_1}, \ldots, \mathsf{N}_{\ell_k}) \to \mathsf{N}_{\ell_0} \in \mathcal{R}$, $A = \lambda u_1, \ldots, u_k$. B, TailPos(f, A), $\Gamma(x_1) = \mathsf{N}_{\ell_1}, \ldots, \Gamma(x_k) = \mathsf{N}_{\ell_k}$, and ζ is as before. Then $\mathcal{T}\llbracket(A \ \vec{x})\rrbracket \leq \mathcal{T}\llbracket(A \ \vec{x}) \zeta \rrbracket \uplus (\mathcal{T}\llbracket f \rrbracket \star \vec{t})$, where $(f \ e_1^1 \ \ldots \ e_k^1), \ldots, (f \ e_1^m \ \ldots \ e_k^m)$ are the full applications of f occurring in B and $t_j = (\bigvee_{i=1}^m \mathsf{val}(\mathcal{T}\llbracket e_i^i \rrbracket))[\vec{u} := \vec{x}]$ for $j = 1, \ldots, k$.

15. The time-complexity interpretation of ATR

We are now in a position to consider the time complexity properties of crec expressions. Remark 77 below motivates the \mathcal{T} -interpretation of crec expressions given in Definition 78. The remark's analysis will be reused in establishing soundness and polynomial time-boundedness for ATR.

Remark 77. Suppose Γ ; $f: \gamma \vdash A: \gamma$, where $\gamma = (\mathbf{\vec{b}}) \to \mathbf{b}_0 \in \mathcal{R}$ and TailPos(f, A). For each $a \in \mathbb{N}$, let $e_a = (\operatorname{crec} a \ (\lambda_r f \cdot A))$. Thus, Γ ; $_\vdash e_a: \gamma$. Our goal is to express $\mathcal{T}\llbracket e_a \rrbracket$ in terms of $\mathcal{T}\llbracket e_{\mathbf{0}\oplus a} \rrbracket$ so as to later extract recurrences, the solutions of which will provide a closed form polynomial time-bound for e_a . So suppose in the following that $\mathcal{T}\llbracket e_{\mathbf{0}\oplus a} \rrbracket$ has a settled value and that \mathcal{T} -soundness holds for all proper subterms of e_a and their expansions below. In a CEK evaluation of e_a , in one step e_a is rewritten to $\lambda \vec{x} \cdot B_a$, where

$$B_a = (\text{if } |a| \le |x_1| \text{ then } C_a \text{ else } \epsilon) \text{ and } C_a = (A \vec{x})[f := e_{\mathbf{0} \oplus a}].$$

Let $\overline{\Gamma} = \Gamma, \vec{x}: \vec{\mathbf{b}}$. So, $\overline{\Gamma}; f: \gamma \vdash B_a: \mathbf{b}_0$ and $\overline{\Gamma}; f: \gamma \vdash C_a: \mathbf{b}_0$. Fix a CEK-environment $\hat{\rho}$ and a $\rho \in \mathcal{T}[[\overline{\Gamma}; f: \gamma]]$ with $\hat{\rho} \sqsubseteq \rho$. From Figure 20 and Definition 48 it follows that

$$\operatorname{cost}_{\operatorname{CEK}}(B_a, \hat{\rho}) \leq \underline{2} \cdot |\hat{\rho}(x_1)| + \underline{2} \cdot |a| + \underline{5} + \begin{cases} \operatorname{cost}_{\operatorname{CEK}}(C_a, \hat{\rho}), & \text{if } |a| \leq |\hat{\rho}(x_1)|; \\ \underline{1}, & \text{otherwise.} \end{cases}$$
(15.1)

By our \mathcal{T} -soundness assumptions,

$$C_a \sqsubseteq_{\mathbf{b}_0}^{\mathrm{tc}} \mathcal{T}\llbracket C_a \rrbracket.$$
(15.2)

Let ζ be the substitution $[f := \epsilon_{\gamma}]$. By Corollary 76 applied to $(A \ \vec{x})$: $\mathcal{T}\llbracket(A \ \vec{x})\rrbracket \leq \mathcal{T}\llbracket(A \ \vec{x}) \zeta\rrbracket \ \uplus \ (\mathcal{T}\llbracket f\rrbracket \star \vec{t})$, where t_1, \ldots, t_k are as in Theorem 74. Let ξ be the substitution $[f := e_{\mathbf{0} \oplus a}]$. Since f has no occurrence in \vec{t} , we have that $\mathcal{T}\llbracket(A \ \vec{x}) \xi\rrbracket \leq \mathcal{T}\llbracket(A \ \vec{x}) \zeta \xi\rrbracket \ \uplus (\mathcal{T}\llbracket f \ t) \star \vec{t})$ which can be restated as:

$$\mathcal{T}\llbracket C_a \rrbracket \leq \mathcal{T}\llbracket (A \ \vec{x}) \zeta \rrbracket \ \uplus \ (\mathcal{T}\llbracket e_{\mathbf{0} \oplus a} \rrbracket \star \vec{t}).$$
(15.3)

Since $\hat{\rho} \sqsubseteq \varrho$, $|\hat{\rho}(x_1)| \le pot(\mathcal{T}[\![x_1]\!]\varrho)$. So, by Lemma 70(c), (15.1), (15.2), and (15.3), $B_a \hat{\rho} \sqsubseteq_{\mathbf{b}_0}^{\mathrm{tc}} X_a \varrho$, where $X_a: \mathcal{T}[\![\overline{\Gamma}; f: \gamma]\!] \to \mathcal{T}[\![\mathbf{b}_0]\!]$ is given by

$$X_{a} \varrho' = \begin{cases} dally(\underline{c}, \mathcal{T}\llbracket(A \ \vec{x}) \zeta \rrbracket \varrho') \ \uplus \ (\mathcal{T}\llbracket e_{\mathbf{0} \oplus a} \rrbracket \star \vec{t}) \varrho', & \text{if } |a| \le p_{1}; \\ (\underline{c+1}, \underline{0}), & \text{otherwise}; \end{cases}$$

where $p_1 = pot(\mathcal{T}\llbracket x_1 \rrbracket \varrho'), c = 2 \cdot p_1 + 2 \cdot |a| + 5$, and \vec{t} is as before.

By the analysis for the $\rightarrow I$ case in Theorem 67's proof, $(\lambda \vec{x} \cdot B_a) \sqsubseteq_{\gamma}^{\text{tc}} \Lambda_{\star}(\vec{x}, X_a)$. As $\operatorname{cost}_{\operatorname{CEK}}(\lambda \vec{x} \cdot B_a, \underline{)} = 1$, we have that $e_a \sqsubseteq_{\gamma}^{\text{tc}} Y_a$, where $Y_a \stackrel{\text{def}}{=} dally(1, \Lambda_{\star}(\vec{x}, X_a))$.

Definition 78 (The \mathcal{T} -interpretation of ATR). $\mathcal{T}[[\Gamma;] \vdash (\operatorname{crec} a \ (\lambda_r f \cdot A)): \gamma]] \stackrel{\text{def}}{=} Y_a$, where Y_a is as above. Figure 22 provides the the \mathcal{T} -interpretations for the other ATR constructs.

The well-definedness of $\mathcal{T}[[\Gamma;] \vdash (\operatorname{crec} a (\lambda_r f \cdot A)): \gamma]]$ is part of:

Theorem 79. The \mathcal{T} -interpretation of ATR is (a) polynomial time-bounded, (b) monotone, and (c) sound, as well as (d) well-defined.

Proof sketch. All the parts are shown simultaneously by a structural induction on the derivation of $\Gamma; \Delta \vdash e; \gamma$. Along with parts (a)–(d) we also show:

Claim: For all $\rho \in \mathcal{T}[\Gamma; \Delta]$, $\mathsf{Pot}(\mathcal{T}[e]] \rho) \leq \mathcal{T}[\widetilde{p_e}] \rho$, where p_e is the polynomial sizebound for e from Theorem 43 and $\widetilde{p_e}$ is the result of replacing each occurrence of each variable |x| in p_e with $\mathsf{Pot}(x)$. (E.g., if $p = \lambda |z| \cdot (2 * |g|(|z|) + 1)$, then $\widetilde{p_e} = \lambda |z| \cdot (2 * \mathsf{Pot}(g)(\mathsf{Pot}(|z|)) + 1)$.)

Intuitively, \tilde{p}_e is the version of p_e that is over base potentials (Definition 60(b)) instead of lengths and the Claim says that the upper bound on size that is implicit in our \mathcal{T} interpretation, is at least as good as the size bounds of Theorem 43. Through the Claim we are able to make use, in a time-complexity context, of the polynomial bound on the depth of **crec**-recursions from the proof of Theorem 43.

Here, then, is the induction.

For each case, except the crec one, parts (a), (b), and (c) are as in the proof of Theorem 67; part (d) is evident, and the Claim follows from an inspection of the bounds assigned in the proof of Theorem 43 and Definition 63. We thus consider the case of $e = (\operatorname{crec} a(\lambda_r f \cdot A))$ where $\Gamma; f: \gamma \vdash A: \gamma, \gamma = (\sigma_1, \ldots, \sigma_k) \to \mathbf{b}_0 \in \mathcal{R}$, and TailPos(f, A). Without loss of generality, we assume a is a tally string <u>n</u>. So, $\mathbf{0} \oplus a = \underline{n+1}$.

We first import the notation from Remark 77. So, $e = e_{\underline{n}}$, where $e_{\underline{n}}$ is as in Remark 77 with $a = \underline{n}$. Also let $\overrightarrow{\mathcal{T}[x]}$ denote $\mathcal{T}[x_1], \ldots, \mathcal{T}[x_k], \varrho \in \mathcal{T}[\overline{\Gamma};]$, and let m range over $\{\underline{n, n+1, \ldots}\}$. Then, by Remark 77, Definition 78, and Lemma 59 we have: $(\mathcal{T}[e_{\underline{m}}] \star \overrightarrow{\mathcal{T}[x]}) \varrho = (Y_{\underline{m}} \varrho) \star \overrightarrow{\varrho(x)} = (dally(1, \Lambda_{\star}(\vec{x}, X_{\underline{m}}))\varrho) \star \overrightarrow{\varrho(x)} = \mathcal{T}[r_0] \varrho \ \uplus \ X_{\underline{m}} \varrho$, where $r_0 = (\underline{5k+4} + cost(\varrho(x_1)) + \cdots + cost(\varrho(x_k)), \underline{0})$. Let $r_{1,m} = r_0 \ \uplus (\underline{2} \cdot pot(\varrho(x_1)) + \underline{2} \cdot \underline{m} + \underline{6}, \underline{0})$ and $r_{2,m} = r_0 \ \uplus (\underline{2} \cdot pot(\varrho(x_1)) + \underline{2} \cdot \underline{m} + \underline{5}, \underline{0})$. Then, by the definition of $X_{\underline{n}}$ in Remark 77,

$$(\mathcal{T}\llbracket e_{\underline{m}} \rrbracket \star \overrightarrow{\mathcal{T}\llbracket x} \rrbracket) \varrho = \begin{cases} \mathcal{T}\llbracket r_{1,m} \rrbracket \varrho, & \text{if } pot(\varrho(x_1)) \leq \underline{n}; \\ \mathcal{T}\llbracket r_{2,m} \rrbracket \varrho \ \uplus \ \mathcal{T}\llbracket (A \vec{x}) \zeta \rrbracket \varrho \ \uplus \ (\mathcal{T}\llbracket e_{\underline{m+1}} \rrbracket \star \vec{t}) \varrho, \\ & \text{otherwise.} \end{cases}$$
(15.4)

Now let us import some notation from the proof of Theorem 43: Let p_1, \ldots, p_k be the manifestly safe polynomials that bound the sizes of the arguments of f in A and let p'_1, \ldots, p'_k be the polynomials that bound the *final* sizes of said arguments.

Part (d): Well-definedness. Let $\varrho_n = \varrho$. Combine the m = n and m = n + 1 versions of (15.4) to express $(\mathcal{T}\llbracket e_n \rrbracket \star \mathcal{T}\llbracket x \rrbracket) \varrho_n$ in terms of $\mathcal{T}\llbracket e_{n+2} \rrbracket$ and $\varrho_{n+1} =$ the update to ϱ_n produced by the application $(\mathcal{T}\llbracket e_{n+1} \rrbracket \star \vec{t}) \varrho_n$. It follows from the Claim that $pot(\varrho_{n+1}(x_1)) \leq \mathcal{T}\llbracket \tilde{p'_1} \rrbracket \varrho_n$ in terms of of $\mathcal{T}\llbracket e_{\underline{m+1}} \rrbracket$ and $\varrho_m = the update$ to ϱ_{m-1} produced by the application $(\mathcal{T}\llbracket e_{\underline{m+1}} \rrbracket \star \vec{t}) \varrho_n$. It follows from the Claim that $pot(\varrho_{n+1}(x_1)) \leq \mathcal{T}\llbracket \tilde{p'_1} \rrbracket \varrho_n$ in terms of of $\mathcal{T}\llbracket e_{\underline{m+1}} \rrbracket$ and $\varrho_m =$ the update to ϱ_{m-1} produced by the application $(\mathcal{T}\llbracket e_{\underline{m}} \rrbracket \star \vec{t}) \varrho_{m-1}$. The Claim still tells us that $pot(\varrho_{n+1}(x_1)) \leq \mathcal{T}\llbracket \tilde{p'_1} \rrbracket \varrho$. Hence, the otherwise clause of (15.4) can hold only finitely many m. Thus, it follows that, $\mathcal{T}\llbracket e_n \rrbracket$ is defined and total.

Part (b): Monotonicity. Note that the terms $\mathcal{L}_{wt}[\![r_{1,m}]\!]$ and $\mathcal{L}_{wt}[\![r_{2,m}]\!]$ clearly satisfy monotonicity. It follows from the induction hypothesis that the terms $\mathcal{T}[\![(A\vec{x})\zeta]\!]$ and t_1, \ldots, t_k also satisfy monotonicity. It follows from (15.4) that if, for a particular m, the $\mathcal{T}[\![e_{\underline{m}+1}]\!]$ term satisfies monotonicity, then so does $\mathcal{T}[\![e_{\underline{m}}]\!]$. Hence, by the finiteness of the expansion it follows that $\mathcal{T}[\![e_n]\!]$ satisfies monotonicity.

Part (c) and the Claim. By arguments along the lines of the one just given for monotonicity, one can establish soundness and the Claim for e_n . Part (a): Polynomial time-boundedness. Recall from Definition 64, the definition of polynomial time-boundedness, the key inequality to be shown is $\mathcal{T}\llbracket e \rrbracket \|\rho\| \leq \mathcal{L}_{wt}\llbracket p_e \rrbracket |\rho|$ for each $\rho \in \mathcal{V}_{wt}\llbracket \Gamma; \Delta \rrbracket$. So if ρ is an ATR-environment and x is a variable with a string or oracle value, then $\mathcal{T}\llbracket \mathsf{Pot}(x) \rrbracket \|\rho\| = \mathcal{L}_{wt}\llbracket |x| \rrbracket |\rho|$. Thus, for each e', $\mathcal{T}\llbracket \widetilde{p_{e'}} \rrbracket \|\rho\| = \mathcal{L}_{wt}\llbracket p_{e'} \rrbracket |\rho|$.

Now, it follows from the induction hypothesis that there is an \overline{r} , a polynomial timebound for $(A \vec{x})\zeta$ relative to Γ ; . Let $r'_1 = r_{1,p'_1}$ and $r'_2 = r_{2,p'_1}$. Let ξ and ξ' respectively denote the substitutions $[|x_1| := p_1, \ldots, |x_k| := p_k]$ and $[|x_1| := p'_1, \ldots, |x_k| := p'_k]$, where $p_1, \ldots, p_k, p'_1, \ldots, p'_k$ are the polynomials from Theorem 43 introduced before. Note that $||x_i|| \xi = (\underline{1} \lor p_i, p_i)$. By the Claim, for each $j = 1, \ldots, k$, $\mathcal{T}[\![t_j]\!] ||\overline{\rho}|| \leq \mathcal{L}_{wt}[\![(\underline{1} \lor p_j, p_j)]\!] |\overline{\rho}| = \mathcal{L}_{wt}[\![(||x_j|| \xi]\!] |\overline{\rho}|]$. Hence, assuming $pot(\mathcal{L}_{wt}[\![|x_1|]\!] |\overline{\rho}|) > \underline{n}$,

$$\begin{aligned} (\mathcal{T}\llbracket e_{\underline{n+1}} \rrbracket \star \vec{t}) \| \overline{\rho} \| &\leq (\mathcal{T}\llbracket e_{\underline{n+1}} \rrbracket \star (\overline{\|x\|} \xi)) \| \overline{\rho} \| & \text{(by monotonicity)} \\ &\leq \mathcal{L}_{\text{wt}} \llbracket (r'_2 \uplus \overline{r}) \xi \rrbracket | \overline{\rho} | \uplus (\mathcal{L}\llbracket e_{\underline{n+2}} \rrbracket \star (\vec{t} \xi)) \| \overline{\rho} \| & \text{(by (15.4))} \\ &\leq \mathcal{L}_{\text{wt}} \llbracket (r_{2,m} \uplus \overline{r}) \xi \rrbracket | \overline{\rho} | \uplus (\mathcal{L}\llbracket e_{\underline{n+2}} \rrbracket \star (\vec{t} \xi)) \| \overline{\rho} \| & \text{(by monotonicity).} \end{aligned}$$

Clearly, we can repeat the above expansion $(p_1 - n)$ -many times (i.e., until termination), collect terms, and produce the desired polynomial bound. Here is the algebra. Let $s = r'_1 \xi' \uplus$ $\bigcup_{m=0}^{p'_1 - n} (r'_2 \uplus \overline{r}) \xi^{(m)}, s_1 = cost(r'_1 \xi') + p'_1 \cdot (cost(r'_2 \xi') + cost(\overline{r}\xi')), \text{ and } s_2 = pot((r'_1 \lor r'_2 \lor \overline{r})\xi').$ Then $(\mathcal{T}\llbracket e_{\underline{n}} \rrbracket \star \lVert x \rVert) \parallel \overline{\rho} \parallel \leq \mathcal{L}_{wt} \llbracket s \rrbracket \mid \overline{\rho} \mid \leq \mathcal{L}_{wt} \llbracket (s_1, s_2) \rrbracket \mid \overline{\rho} \mid$ by a straightforward argument. Thus, $\lambda \mid \overline{x} \mid \cdot (s_1, s_2)$ suffices as the polynomial time bound for $e_{\underline{n}}$.

Scholium 80. Note that we resorted to reasoning directly about CEK-costs to obtain (15.1). This is because if we had used Definition 63's \mathcal{T} -interpretation of if-then-else, then we would have be left without a base case in our recursive unfoldings of crec-expressions.

We note that as a consequence of parts (a) and (c) of Theorem 79 we have:

Corollary 81. For each $\Gamma; \Delta \vdash e; \gamma$, there is a second-order polynomial q_e with $|\Gamma; \Delta| \vdash q_e: |\gamma|$ such that $\operatorname{cost}_{CEK}(e, \rho) \leq \mathcal{L}_{wt}[\![q_e]\!]|\rho|$ for each $\rho \in \mathcal{V}_{wt}[\![\Gamma; \Delta]\!]$.

Remark 82 (Related work). The time-complexity cost/potential distinction appears in prior work. A version of this distinction can be found in Sands' Ph.D. thesis [San90]. Shultis [Shu85] sketched how to use the distinction in order to give time-complexity semantics for reasoning about the run-time programs that involve higher types. Van Stone [VS03] gives a much more detailed and sophisticated semantics for a variant of PCF using the cost/potential distinction. Very roughly, Shultis and Van Stone were focused on giving static analyses to extract time-bounds for functional programs that compute first-order functions. The time-complexity semantics of this paper was developed independently of Shultis' and Van Stone's work. We also note that Benzinger's work [Ben01, Ben04] on automatically inferring the complexity of Nuprl programs made extensive use of higher-type recurrence equations.

16. Complexity-theoretic completeness

Our final result on ATR is that each type-1 and type-2 BFF is ATR computable. Conventions: In this section, let $\sigma = (\sigma_1, \ldots, \sigma_k) \to \mathbb{N}$ range over simple types over \mathbb{N} of levels 1 or 2, and let $\gamma, \gamma_0, \gamma_1, \ldots$ range over ATR types. Recall from §2.14 that $f \in \mathcal{V}[\![\sigma]\!]$ is basic feasible when there is a closed type- σ , PCF-expression e_f and a second-order polynomial function q_f such that $\mathcal{V}[\![e_f]\!] = f$ and, for all $v_i \in \mathcal{V}[\![\sigma_1]\!], \ldots, v_k \in \mathcal{V}[\![\sigma_k]\!]$, CEK-time $(e_f, v_1, \ldots, v_k) \leq q_f(|v_1|, \ldots, |v_k|)$. Let BFF_{σ} = the class of all type- σ BFFs. **Definition 83.** We say that each base type is *unhindered* and that $(\gamma_1, \ldots, \gamma_k) \to \mathsf{N}_\ell$ is *unhindered* when $(\gamma_1, \ldots, \gamma_k) \to \mathsf{N}_\ell$ is strict, predicative and each γ_i unhindered.

Note that $\mathcal{V}_{wt}[\![\gamma]\!] = \mathcal{V}[\![shape(\gamma)]\!]$ if and only if γ is unhindered.

Theorem 84. BFF_{σ} = { \mathcal{V}_{wt} [$\vdash e: \gamma$] : $\sigma = shape(\gamma) \& \gamma$ is unhindered } for each σ .

Proof. Fix σ and let $\mathcal{U}_{\sigma} = \{ \mathcal{V}_{wt} \llbracket \vdash e : \gamma \rrbracket \mid \sigma = shape(\gamma) \& \gamma \text{ is unhindered } \}.$

Claim 1: $\mathcal{U}_{\sigma} \subseteq BFF_{\sigma}$. Proof: It is straightforward to express a crec-recursion with PCF's fix-construct with only polynomially-much over head on the cost of the simulation. Hence, the claim follows from Theorem 79.

Claim 2: $\text{BFF}_{\sigma} \subseteq \mathcal{U}_{\sigma}$. Proof: Kapron and Cook [KC96] showed that the type-2 basic feasible functionals are characterized by the functions computable in second-order polynomial time-bounded oracle Turing machines (OTMs). Proposition 18 from [IKR01] shows how to simulate any second-order polynomial time-bounded oracle Turing machine using that paper's ITLP_2 programming formalism. That simulation is easily adapted to ATR. Hence, the claim follows.

Note: The proof's two claims are constructive in that: (i) given a closed ATR-expression e of unhindered type, one can construct an equivalent PCF expression e' and a second-order polynomial p_e that bounds the run time of e', and (ii) given an OTM **M** and a second-order polynomial p that bounds the run time of **M**, one can construct an ATR-expression that computes the same function as **M**.

Claim 2 can be extended beyond unhindered types as follows. For each ATR arrowtype $\gamma = (\gamma_1, \ldots, \gamma_k) \rightarrow N_{\ell}$, and each type-*shape*(γ) OTM **M**, we say that **M** computes a BFF $_{\gamma}$ -function when there is a type- $|\gamma|$ polynomial p such that the run time of **M** on (\vec{v}) is bounded by $p(|v_1|, \ldots, |v_k|)$. The proof of Claim 2 lifts to show: for all ATR arrow-types γ , each BFF $_{\gamma}$ -function is ATR computable.

17. Conclusions

ATR is a small functional language, based on PCF, which has the property that each ATR program has a second-order polynomial time-bound. The ATR-computable functions include the basic feasible functionals at type-levels 1 and 2. However, the ATR-computable functions contain other functions, such as prn, that are *not* basic feasible in the original sense of Cook and Urquhart [CU93]. ATR is able to express such functions thanks to its type system and supporting semantics that work together to control growth rates and time complexities. Without some such controls feasible recursion schemes, such as prn, cannot be first-class objects of a programming language.

The ATR type-system and semantics were crafted so that ATR's complexity properties could be established through adaptations of standard tools for the analysis of conventional programming languages (e.g., intuitionistic and affine types, denotational semantics for ATR and its time complexity, and an abstract machine that provides both an operational semantics for ATR and a basis for the time-complexity semantics). As ATR is based on PCF (a theoretical first-cousin of both ML and Haskell), our results suggest that one might be able to craft "feasible sublanguages" of ML and Haskell that are both theoretically well-supported and tolerable for programmers.

ATR and its semantic and analytic frameworks are certainly not the final word on any issue. Here we discuss several possible extensions of our work.

MORE GENERAL RECURSIONS. In [DR07] we consider an expansion of ATR that allows a fairly wide range of affine (one-use) recursions. In particular, the expanded ATR can fairly naturally express the classic insertion- and selection-sort algorithms. Handling this larger set of recursions requires some nontrivial extensions of our framework for analyzing time-complexities.

Dealing with nonlinear recursions (e.g., the standard quicksort algorithm) is trickier to handle because there must be independent clocks on each branch of the recursion that together guarantee certain global upper bounds.

RECURSIONS WITH TYPE-LEVEL 1 PARAMETERS. Another possible extension of ATR would be to allow type-level 1 parameters in crec-recursions so that, for example, one could give a continuation-passing-style definition of *prn*. Because type-1 parameters in recursions act to recursively define functions, these parameters must be affinely restricted just like principle recursor variables of crec-expressions. Consequently, such an extension must also include explicit —o-types to restrict these parameters. However, along with the —o-types come (explicitly or implicitly) tensor-products and these cause problems in analyzing crecrecursions (e.g., one is forced account for all the possible interactions of the affine parameters in the course of a recursion and so the naïve "polynomial" time-bounds are exponential in size).

LAZY EVALUATION. For a lazy (e.g., call-by-need) version of ATR, one would need to: (i) construct an abstract machine for this lazy-ATR, (ii) modify the \mathcal{T} -semantics a bit to accommodate the lazy constructs; and (iii) rework the \mathcal{T} -interpretation of ATR which would then have to be shown monotone, sound, and constructively polynomial time-bounded. (Since the well-tempered semantics is extensional, it requires very few changes for a lazy-ATR.) If our lazy-ATR allowed infinite strings, then the \mathcal{V}_{wt} -semantics would also have to be modified. Note that Sands [San90] and Van Stone [VS03] both consider lazy evaluation in their work.

LISTS AND STREAMS. There are multiple senses of the "size" of a list. For example, the run-time of *reverse* should depend on just a list's length, whereas the run-time of a search depends on both the list's length and the sizes of the list's elements. Any useful extension of ATR that includes lists needs to account for these multiple senses of size in the type system and the well-tempered and time-complexity semantics. If lists are combined with laziness, then we also have the problem of handling infinite lists. However, ATR and its semantics already handle one flavor of infinite object, i.e., type-level 1 inputs, so handling a second flavor of infinite object many not be too hard.

TYPE CHECKING, TYPE INFERENCE, TIME-BOUND INFERENCE. We have not studied the problem of ATR type checking. But since ATR is just an applied simply typed lambda calculus with subtyping, standard type-checking tools should suffice. Type inference is a much more interesting problem. We suspect that a useful type inference algorithm could be based on Frederiksen and Jones' [FJ04] work on applying size-change analysis to detect whether programs run in polynomial time. Another interesting problem would be to start with a well-typed ATR program and then extract reasonably tight size and time bounds (as opposed to the not-so-tight bounds given by Theorem 79).

BEYOND TYPE-LEVEL 2. There are semantic and complexity-theoretic issues to be resolved in order to extend the semantics of ATR to type-levels 3 and above. The key problem is that our definition of the length of a type-2 function (2.3) does not generalize to typelevel 3. This is because for $\Psi \in \mathbf{MC}_{((N \to N) \to N) \to N}$ and $G \in \mathbf{MC}_{(N \to N) \to N}$, we can have $\sup\{|\Psi(F)| \mid |F| \leq |G|\} = \infty$, even when G is 0–1 valued. To fix this problem one can introduce a different notion of length that incorporates information about a function's modulus of continuity. It appears that ATR and the \mathcal{V}_{wt} - and \mathcal{T} -semantics extend to this new setting. However, it also appears that this new notion of length gives us a new notion of higher-type feasibility that goes beyond the BFFs. Sorting out what is going on here should be the source of other adventures.

References

- [Bar96] A. Barber, Dual intuitionistic linear logic, Tech. report, LFCS, Univ of Edinburgh, 1996.
- [BC92] S. Bellantoni and S. Cook, A new recursion-theoretic characterization of the polytime functions, Computational Complexity 2 (1992), 97–110.
- [Ben01] R. Benzinger, Automated complexity analysis of Nuprl extracted programs, Journal of Functional Programming 11 (2001), 3–31.
- [Ben04] _____, Automated higher-order complexity analysis, Theoretical Computer Science **318** (2004), 79–103.
- [BNS00] S. Bellantoni, K.-H. Niggl, and H. Schwichtenberg, Characterising polytime through higher type recursion, Annals of Pure and Applied Logic 104 (2000), 17–30.
- [BP97] A. Barber and G. Plotkin, *Dual intuitionistic linear logic*, Tech. report, LFCS, Univ of Edinburgh, 1997.
- [CK90] S. Cook and B. Kapron, Characterizations of the basic feasible functions of finite type, Feasible Mathematics: A Mathematical Sciences Institute Workshop (S. Buss and P. Scott, eds.), Birkhäuser, 1990, pp. 71–95.
- [Cob65] A. Cobham, The intrinsic computational difficulty of functions, Proceedings of the International Conference on Logic, Methodology and Philosophy (Y. Bar Hillel, ed.), North-Holland, 1965, pp. 24– 30.
- [CU93] S. Cook and A. Urquhart, Functional interpretations of feasibly constructive arithmetic, Annals of Pure and Applied Logic 63 (1993), 103–200.
- [DR06] N. Danner and J. Royer, Adventures in time and space, 33th ACM Symposium on Principles of Programming Languages (S. Peyton Jones, ed.), ACM Press, 2006, pp. 168–179.
- [DR07] _____, Time-complexity semantics for feasible affine recursions, Computation and Logic in the Real World: Third Conference of Computability in Europe, CiE 2007 (S.B. Cooper, B. Löwe, and A. Sorbi, eds.), Lecture Notes in Computer Science, vol. 4497, Springer-Verlag, 2007, to appear.
- [FF87] M. Felleisen and D. Friedman, Control operators, the SECD-machine, and the lambda calculus, Formal Descriptions of Programming Concepts III, 1987, pp. 193–217.
- [FF06] M. Felleisen and M. Flatt, Programming languages and lambda calculi, unpublished manuscript, 2006.
- [FJ04] C. Frederiksen and N. Jones, Recognition of polynomial-time programs, Tech. Report TOPPS/D-501, DIKU, University of Copenhagen, 2004.
- [FWH01] D. Friedman, M. Wand, and C. Haynes, Essentials of programming langauges, second ed., MIT Press, 2001.
- [Gol01] O. Goldreich, Foundations of cryptography, Vol. I: Basic tools, Cambridge University Press, 2001.
- [Gur90] D. J. Gurr, Semantic frameworks for complexity, Ph.D. thesis, University of Edinburgh, 1990.
- [Hof00] M. Hofmann, Programming languages capturing complexity classes, SIGACT News 31 (2000), 31– 42.
- [Hof02] _____, The strength of non-size increasing computation, 29th ACM Symposium on Principles of Programming Languages (J. Michell, ed.), ACM Press, 2002, pp. 260–269.
- [Hof03] _____, Linear types and non-size increasing polynomial time computation, Information and Computation 183 (2003), 57–85.
- [IKR01] R. Irwin, B. Kapron, and J. Royer, On characterizations of the basic feasible functionals, Part I, Journal of Functional Programming 11 (2001), 117–153.

- [IKR02] _____, On characterizations of the basic feasible functionals, Part II, unpublished manuscript, 2002.
- [Kap91] B. Kapron, Feasible computation in higher types, Ph.D. thesis, Department of Computer Science, University of Toronto, 1991.
- [KC96] B. Kapron and S. Cook, A new characterization of type 2 feasibility, SIAM Journal on Computing 25 (1996), 117–132.
- [KU58] A.N. Kolmogorov and V.A. Uspenskii, On the definition of an algorithm, Uspekhi Mat. Nauk 13 (1958), 2–28.
- [Lei94] D. Leivant, A foundational delineation of poly-time, Information and Computation 110 (1994), 391–420.
- [Lei95] _____, Ramified recurrence and computational complexity I: Word recurrence and poly-time, Feasible Mathematics II (P. Clote and J. Remmel, eds.), Birkhäuser, 1995, pp. 320–343.
- [Lei03] _____, Feasible functionals and intersection of ramified types, Proceedings of the Second Workshop on Intersection Types and Related Systems, Electronic Notes in Theoretical Computer Science, vol. 70, Elsevier Science Publishers, 2003, pp. 1–14.
- [LM93] D. Leivant and J.-Y. Marion, Lambda calculus characterizations of polytime, Fundamentæ Informaticæ 19 (1993), 167–184.
- [Lon04] J. Longley, On the ubiquity of certain total type structures (Extended abstract), Proceedings of the Workshop on Domains VI (M. Escardó and A. Jung, eds.), Electronic Notes in Theoretical Computer Science, vol. 73, Elsevier Science Publishers, 2004, pp. 87–109.
- [Lon05] _____, Notions of computability at higher types I, Logic Colloquium 2000 (R. Cori, A. Razborov, S. Torcevic, and C. Wood, eds.), Lecture Notes in Logic, vol. 19, A. K. Peters, 2005.
- [Mar72] S. Marchenkov, The computable enumerations of families of general recursive functions, Algebra and Logic 11 (1972), 326–336.
- [Meh74] K. Mehlhorn, Polynomial and abstract subrecursive classes, Proceedings of the Sixth Annual ACM Symposuium on the Theory of Computing, 1974, pp. 96–109.
- [Meh76] _____, Polynomial and abstract subrecursive classes, Journal of Computer and System Science 12 (1976), 147–178.
- [Nor99] D. Normann, The continuous functionals, Handbook of Computability Theory (E. R. Griffor, ed.), North-Holland, 1999, pp. 251–275.
- [O'H03] P. O'Hearn, On bunched typing, Journal of Functional Programming 13 (2003), 747–796.
- [Pie02] B. Pierce, Types and programming languages, MIT Press, 2002.
- [Plo75] G. Plotkin, Call-by-name, call-by-value and the λ -calculus, Theoretical Computer Science 1 (1975), 125–159.
- [Plo77] _____, *LCF considered as a programming language*, Theoretical Computer Science **5** (1977), 223–255.
- [RC94] J. Royer and J. Case, Subrecursive programming systems: Complexity & succinctness, Birkhäuser, 1994.
- [Rey72] J. Reynolds, Definitional interpreters for higher-order programming languages, Proceedings of the ACM National Conference, 1972, pp. 717–740.
- [Rey93] J. Reynolds, The discoveries of continuations, Lisp and Symbolic Computation 6 (1993), 233–247.
- [Rey98] J. Reynolds, Definitional interpreters for higher-order programming languages, Higher-Order and Symbolic Computation 11 (1998), 363–397, reprint of [Rey72].
- [Roy87] J. Royer, A connotational theory of program structure, Lecture Notes in Computer Science, vol. 273, Springer-Verlag, 1987.
- [San90] D. Sands, Calculi for time analysis of functional programs, Ph.D. thesis, University of London, 1990.
- [Sch80] A. Schönhage, Storage modification machines, SIAM Journal on Computing 8 (1980), 490–508.
- [Sch96] H. Schwichtenberg, Density and choice for total continuous functionals, Kreiseliana (P. Odifreddi, ed.), A.K. Peters, 1996, pp. 335–362.
- [Shu85] J. Shultis, On the complexity of higher-order programs, Tech. Report CU-CS-288-85, University of Colorado, Boulder, 1985.
- [VS03] K. Van Stone, A denotational approach to measuring complexity in functional programs, Ph.D. thesis, School of Computer Science, Carnegie Mellon University, 2003.
- [Win93] G. Winskel, Formal semantics, MIT Press, 1993.