

1 **A low-level look at A-normal form**

2 WILLIAM J. BOWMAN, University of British Columbia, Canada

3
4 A-normal form (ANF) is a widely studied intermediate form in which local control and data flow is made
5 explicit in syntax, and a normal form in which many programs with equivalent control-flow graphs have a
6 single normal syntactic representation. However, ANF is difficult to implement effectively and, as we formalize,
7 unsafe for scoped effects such as scoped region-based allocation. The problem, as has often been observed, is
8 the normalization of commuting conversions.

9 We argue that the traditional view of ANF, informed by high-level languages, is wrong. By studying
10 the low-level intensional aspects of ANF, we can derive a normal form in which normalizing commuting
11 conversion is easy, does not require join points, or code duplication, or renormalization after inlining, and is
12 safe for scope. We formalize the connection between ANF and monadic form and their intensional properties,
13 derive an imperative ANF, and design a compiler pipeline from an untyped λ -calculus with scoped regions, to
14 monadic form, to a low-level imperative monadic form in which A-normalization is trivial and safe. We prove
15 that any such compiler preserves, or optimizes, stack and memory behaviour compared to ANF.

16 The main take-away from this work is that, in general, monadic form should be preferred over ANF,
17 and A-normalization should only be done in a low-level imperative intermediate form. This maximizes the
18 advantages of each form, and avoids all the standard problems with ANF.

19 CCS Concepts: • **Theory of computation** → **Type structures**; • **Software and its engineering** → **Formal**
20 **software verification**; *Software performance*.

21 Additional Key Words and Phrases: Compilers, Optimization, A-normal form, CPS, Monadic Form, Intermediate
22 Representation, Normal Form, Normalization

23 **ACM Reference Format:**

24 William J. Bowman. 2024. A low-level look at A-normal form. 1, 1 (April 2024), 24 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

25
26
27 **1 INTRODUCTION**

28 Intermediate representations, forms, and languages are used to simplify program analysis, opti-
29 mization, and compilation, e.g., by (1) explicating abstractions, such as evaluation order or data
30 flow, into explicit representations in syntax; (2) normalizing programs, so that many programs
31 equal under some equivalence class have a single normal representation; or (3) providing practical
32 equational theories for optimization or transformation.

33 A-normal form (ANF) is an intermediate representation widely studied in compilation [3, 8,
34 11, 13, 14, 17]. ANF can be described syntactically as an untyped λ -calculus with **let** where all
35 *computations* must take *values* as their operands, and intermediate computations are explicitly
36 sequenced using **let**. This makes data and local control flow (everything except returning from a
37 call) explicit in the syntax, simplifying analysis, optimization, and compilation.

38 For example, consider the following two equivalent terms, where Listing 2 is the A-normalization
39 of Listing 1.

40
41 Author’s address: William J. Bowman, University of British Columbia, Vancouver, Canada, wjb@williamjbowman.com.

42
43 Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee
44 provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and
45 the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses,
46 contact the owner/author(s).

47 © 2024 Copyright held by the owner/author(s).

48 XXXX-XXXX/2024/4-ART

49 <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

50
51 $(+ (\mathbf{let} (x (f\ 5))\ 0)\ 6)$

52 Listing 1. λ -calculus

53 $(\mathbf{let} (x (f\ 5))\ (+\ 0\ 6))$

54 Listing 2. ANF

55 Code generation from the Listing 2 is simpler than Listing 1, since the computation $(+ 0 6)$ can be
56 compiled to something like `mov y 0; add y 6`. Optimization is simpler; $(+ 0 6)$ is trivially equal to
57 6, whereas optimizing the original term requires some additional control and data flow analysis.
58 ANF explicates many stack frames into syntax, so the ANF abstract machine runs with a smaller
59 stack, in general, than the λ -calculus machine. In this example, the λ -calculus term evaluates the
60 first operand of $+$ to a value, after pushing the frame $(+ \cdot 6)$, while the ANF term never pushes this
61 frame.

62 ANF has several known disadvantages. We discuss these in detail, but in short: ANF is not closed
63 under β -reduction (complicating inlining) and compiling branching expressions into ANF requires
64 care to avoid duplicating expressions or introducing procedure calls.

65 ANF suffers at least one additional disadvantage that is not discussed in the literature: trans-
66 formation into ANF is *unsafe for scope*. In Listing 2, the scope of x is extruded past the addition
67 expression compared to Listing 1. Effects attached to lexical binding can be reordered.

68 For example, we can translate Listing 1 and Listing 2 into a language with the lexically scoped
69 region system of Tofte and Talpin [20], where $(\mathbf{letregion}\ r\ e)$ allocates a new region r for use in e
70 and frees that region when the expression returns, and $(@ r\ e)$ allocates in r the value resulting
71 from evaluating e . All values must be explicitly allocated and passed by reference.

72 $(\mathbf{letregion}\ r1$
73 $\ (\ @\ r1\ (+\ (\mathbf{letregion}\ r2$
74 $\ (\mathbf{let}\ (x\ (f\ (@\ r2\ 5)))$
75 $\ (\ @\ r1\ 0)))$
76 $\ (\ @\ r1\ 6))))$

77 Listing 3. λ -calculus with Regions

78 $(\mathbf{letregion}\ r1$
79 $\ (\mathbf{letregion}\ r2$
80 $\ (\mathbf{let}\ (x1\ (@\ r2\ 5))$
81 $\ (\mathbf{let}\ (x\ (f\ x1))$
82 $\ (\mathbf{let}\ (x2\ (@\ r1\ 0))$
83 $\ (\mathbf{let}\ (x3\ (@\ r1\ 6))$
84 $\ (\ @\ r1\ (+\ x2\ x3))))))$

85 Listing 4. ANF with Regions

86 The term in Listing 3 allocates an inner region $r2$, allocating 5 in region $r2$, calling f before
87 allocating a result in the outer region $r1$. $r2$ is freed when the computation of x completes. However,
88 in Listing 4, $r2$ is not freed until the end of the program. We must make the allocation and free
89 operations explicit before we can safely (as in safe for space [19]) target ANF.

90 Typically, ANF is contrasted with continuation-passing style (CPS), another syntactic discipline
91 popular as an intermediate form, which also explicates control and data flow. We ignore CPS until
92 Subsection 6.1, but instead consider a more related alternative: monadic form.

93 Monadic form is the syntactic discipline induced by Moggi's monadic meta-language [15] when
94 treating all non-value expressions as effectful computations (consider, e.g., partiality as the effect).
95 We can pronounce the monadic bind as **let**, and implement monadic return as an untagged inclusion
96 of values into computations. For example, Listing 1 (reproduced in Listing 5 for comparison) could
97 be implemented in monadic form as either Listing 6 or Listing 7, since all ANF terms are also in
98 monadic form. Neither term would be produced by the usual translation of Listing 1, but they're
99 useful examples.

99
100 $(+ (\mathbf{let} (x (f\ 5))$
101 $0))$
102 $6)$

103 Listing 5. λ -calculus
104

$(\mathbf{let} (y (\mathbf{let} (x (f\ 5))$
101 $0))$
102 $(+ y\ 6))$

103 Listing 6. A monadic equivalent
104

$(\mathbf{let} (x (f\ 5))$
101 $(\mathbf{let} (y\ 0)$
102 $(+ y\ 6)))$

103 Listing 7. An ANF equivalent
104

105 Monadic form suffers none of the disadvantages of ANF. However, it is *less normal*. Being less
106 normal, transformations and analyses can be *less obvious* in monadic form compared to ANF. In
107 Listing 6, the computation $(+ y\ 6)$ only takes values as operands, but it's not obvious how to optimize
108 the expression. In Listing 7, we can obviously inline y , since it is bound to a constant.

109 ANF is monadic form with all commuting conversions normalized. Listing 6 and Listing 7 are
110 equal by associativity of bind, one of the commuting conversions. ANF also normalizes commuting
111 conversions for conditionals, such as $(\mathbf{let} (x (\mathbf{if} v\ e_1\ e_2))\ e_3) \equiv (\mathbf{if} v (\mathbf{let} (x\ e_1)\ e_3) (\mathbf{let} (x\ e_2)\ e_3))$. Du-
112 plicating e_3 is undesired, so typically this is abstracted into a continuation, called a *join point*, as in
113 $(\mathbf{let} (x (\mathbf{if} v\ e_1\ e_2))\ e_3) \equiv (\mathbf{let} (j (\lambda (y)\ e_3)) (\mathbf{if} v (\mathbf{let} (x\ e_1)\ (j\ x)) (\mathbf{let} (x\ e_2)\ (j\ x))))$.

114 This difference between ANF and monadic form, normalization of commuting conversions, is at
115 the heart of ANF. It is both why ANF is attractive as an intermediate form, and causes many of the
116 problems of ANF.

117 This is not a novel observation. Kennedy [11] observes that normalizing commuting conversions
118 causes all the well known problems with ANF, before resorting to CPS. Maurer et al. [14] point out
119 these problems as well, observing that while join points avoid duplication, they inhibit optimizations,
120 and create a join point calculus to recover these optimizations in something ANF-like.

121 The problem with these commuting conversions is that, while the programs are extensionally
122 equal even when reasoning about monadic effects, they are intentionally different when we consider
123 details related to efficient execution and compilation. Duplicating code may be extensionally fine, but
124 it's intentionally bad. Commuting conversions are also not equal for otherwise-effectful programs
125 that use monadic form or ANF as intermediate forms but not to express monadic effects. That is
126 why when we add scoped regions, A-normalization causes a new problem.

127 In this paper, we formalize a novel observation: join points are unnecessary and commuting
128 conversions are not a problem. Instead, commuting conversions are a specification, not a literal im-
129 plementation technique. By fully understanding the intentional aspects of commuting conversions,
130 we can gain the benefits of ANF with none of the drawbacks. We make these intentional aspects
131 formal using abstract machines. To contrast to Maurer et al. [14], our motto is: *work in monadic form,*
132 *but think in abstract machines*. If we think of $(\mathbf{let} (x (\mathbf{if} v\ e_1\ e_2))\ e_3)$ as its computation in an abstract
133 machine, we might render it as $(\mathbf{begin} (\mathbf{set!} x (\mathbf{if} v\ e_1\ e_2))\ e_3)$, which by (a low-level interpretation
134 of) associativity is the same as $(\mathbf{begin} (\mathbf{if} v (\mathbf{set!} x\ e_1) (\mathbf{set!} x\ e_2))\ e_3)$. No duplication, no join points
135 needed. In fact, this idea is understood by some compiler writers; the high-performance Chez
136 Scheme compiler uses a similar transformation (Subsection 6.4). Our paper performs a rational
137 reconstruction and explanation of what some compiler writers already do.

138 Concretely, our contributions are:
139
140

- 141 (1) A novel formalization of monadic form as a subset of A-normal form, which is important
142 for clarifying the distinction between ANF and monadic form, and identifying the source of
143 problems with ANF (Section 2).
- 144 (2) A novel analysis of the abstract machine of high-level ANF and monadic form, with which
145 we formalize how ANF optimizes stack usage (Section 3).
146
147

$$\begin{aligned}
v &::= \iota \mid x \mid (\lambda (x) e) \\
e &::= v \mid (op \vec{e}) \mid (e e) \mid (\mathbf{let} (x e) e) \mid (\mathbf{if0} e e e) \\
E &::= \cdot \mid (\mathbf{let} (x E) e) \mid (\mathbf{if0} E e e) \mid (E e) \mid (v E) \mid (op \vec{v} E \vec{e}) \\
O &::= v \mid op
\end{aligned}$$

$$\begin{aligned}
E[(\mathbf{let} (x e_1) e_2)] &\longrightarrow_A (\mathbf{let} (x e_1) E[e_2]) & A_1 \\
&\text{where } E \neq \cdot \\
E[(\mathbf{if0} v e_1 e_2)] &\longrightarrow_A (\mathbf{if0} v E[e_1] E[e_2]) & A_2 \\
&\text{where } E \neq \cdot \\
E[(O \vec{v})] &\longrightarrow_A (\mathbf{let} (x' (O \vec{v})) E[x']) & A_3 \\
&\text{where } E \neq \cdot, E \neq E'[(\mathbf{let} (x \cdot) e)], \text{ fresh } x'
\end{aligned}$$

Fig. 1. A -normalization for λ -calculus

- (3) A novel formalization of a counterexample to the correctness of A -normal form with respect to scoped effects (*unsafe for scope*), which is important for understanding limitations of ANF in compilation (Subsection 3.2).
- (4) A novel formalization of imperative variants of monadic form and ANF derived from the machine semantics and admissible equations, and an algorithm for normalizing commuting conversions without join points or duplication. These are important as normal forms for compilation, optimization, and analysis as they enable the advantages of ANF but suffers none of the known disadvantages (Section 4).
- (5) A novel analysis of the abstract machine for imperative monadic form, with which we prove that using monadic form followed by imperative A -normalization has the same or better performance characteristics as ANF. In particular, we show that any such compiler (1) avoids code duplication and join point introduced by commuting conversion in high-level ANF; (2) preserves the stack behaviour of ANF; and (3) preserves the memory usage of scoped regions, unlike the ANF compiler (Section 4).
- (6) A model compiler designed to use monadic form as a high-level intermediate language, and imperative ANF as a low-level IL. We argue that monadic form followed by imperative A -normalization simplifies compiler implementation compared to ANF (Section 5).

All machines, reduction systems, languages, compilers, examples, and counterexample are implemented in a PLT Redex [6, 12], available in the anonymous supplementary materials.

2 A-NORMAL AND MONADIC FORM, FORMALLY

2.1 A-normal Form

A -normal form (ANF) is often called “administrative normal form” or sometimes “administrative form”, but it is important and useful to think about ANF not as a vague form related to administrative reductions, but formally and precisely as a normal form with respect to a set of reductions, as it was originally formalized [8].

Formally, A -normal form was introduced as the form normal with respect to the set of reductions $A = \{A_1, A_2, A_3\}$, defined in Figure 1. A term e represents an arbitrary λ -calculus expression. An operator op is some n -ary primitive operator, such as addition, and must appear in operator position. An ι is some ground value, such as a natural number, and a value v is any syntactic value. We use the slight abuse of notation $(O \vec{e})$ to mean an application—of either a function or an operator, and restrict the function position to a value as necessary. The A -reductions use the call-by-value

(Values) $V ::= \iota \mid x \mid (\lambda (x) M)$
 (Computations) $N ::= V \mid (V V) \mid (op \vec{V})$
 (Configuration) $M ::= N \mid (\mathbf{let} (x N) M) \mid (\mathbf{if0} V M M)$

Fig. 2. A-normal form

evaluation contexts E to identify a non-value in evaluation position, and lift and bind it explicitly. This way, the data flow and local control flow is made explicit in the syntax.

The A -reductions require some side conditions about the evaluation context to ensure non-circular rewrites, and therefore to guarantee termination. We assume uniqueness of names and consider terms up to α -equivalence, as is standard.

If we *normalize* a λ -calculus expression by reducing the transitive compatible closure of the set A , we reach a *normal* form with respect to A , *i.e.*, A -normal form, described syntactically by the grammar in Figure 2. The non-terminal N represents computations, while M represents program configurations that sequence computations.

As an example of A -normalization, consider the term in Listing 8, which A -normalizes to the term in Listing 9 by A_3 and then A_1 . All computations are explicitly sequenced using **let**, so all operands are values.

(+ (+ 2 2))	(let (x1 (+ 2 2))
(let (x 1)	(let (x 1)
(f x)))	(let (x2 (f x))
	(+ x1 x2))))

Listing 8. λ -calculus exampleListing 9. A -normalization

Unfortunately, ANF has some drawbacks. ANF is not closed under β -reduction, complicating its calculus. The term (**let** (x ((λ (x') M) V)) x) ought to be β -equivalent to (**let** (x M[x' := V]) x), where x' is substituted by v in M . But this expression is invalid since M cannot appear on the right-hand side of **let**. The ANF β -equivalence must renormalize all commuting conversions. This is a drawback as β -equivalence models inlining optimizations, and renormalization is inconvenient and expensive for a compiler. The A_2 rule causes exponential code duplication by duplicating the continuation E . Consider Listing 11. The term *LARGE* is duplicated 2^3 times while A -normalizing Listing 10— 2^n where n is equal to the occurrences of **if** for which *LARGE* is in the evaluation context. Compilers using ANF avoid this using join points, but this canonical solution causes more problems; we discuss this in Section 5.

(let (x (if0 (if0 (if0 0 0 1)	(if0 0 (if0 0
0	(if0 0 (let (x 0) <i>LARGE</i>) (let (x 1) <i>LARGE</i>))
1)	(if0 1 (let (x 0) <i>LARGE</i>) (let (x 1) <i>LARGE</i>)))
0	(if0 1
1))	(if0 0 (let (x 0) <i>LARGE</i>) (let (x 1) <i>LARGE</i>))
<i>LARGE</i>)	(if0 1 (let (x 0) <i>LARGE</i>) (let (x 1) <i>LARGE</i>)))

Listing 10. Nested Branching

Listing 11. ANF Exponential Duplication

There are further benefits and problems with ANF, but these only become obvious when we consider machines for executing in ANF. We return to these in Section 3.

$$\begin{array}{l}
\text{(Value)} \quad U ::= \iota \mid x \mid (\lambda (x) C) \\
\text{(Computations)} \quad C ::= U \mid (U U) \mid (op \vec{U}) \mid (\mathbf{let} (x C) C) \mid (\mathbf{if0} U C C)
\end{array}$$

Fig. 3. Monadic Form Syntax

2.2 Monadic Form

Because ANF is not merely a syntactic description, but the normal form of a set of reductions, we can tweak those reductions to study alternatives normal forms. In fact, monadic form can be defined as a subset of the A -reductions. We can work backwards from its syntax to the reductions for which monadic form is normal.

We describe monadic form with the grammar in [Figure 3](#), where we syntactically separate values U and effectful computations C . If we consider the effect as partiality, then we are explicitly forcing terms to value before calling each operation. In this definition, **let** corresponds to monadic bind, and return is implicit by the untagged inclusion of U in C .

This form is preserved under the monad laws and commuting conversions.

$$(\mathbf{let} (x U) E[x]) \equiv E[U] \quad \text{(Left Identity)}$$

$$(\mathbf{let} (x C) x) \equiv C \quad \text{(Right Identity)}$$

$$(\mathbf{let} (y (\mathbf{let} (x C) C_1)) C_2) \equiv (\mathbf{let} (x C) (\mathbf{let} (y C_1) C_2)) \quad \text{(Associativity)}$$

$$(\mathbf{let} (x (\mathbf{if0} U C_1 C_2)) C) \equiv (\mathbf{if0} U (\mathbf{let} (x C_1) C) (\mathbf{let} (x C_2) C)) \quad \text{(Commute)}$$

Both sides of the equations are in monadic form.

ANF is a normalization of monadic form: it normalizes the two commuting conversions, [Equation Associativity](#) and [Equation Commute](#). The left-hand side of each of two rules is *not* in ANF, while the right-hand side is.

There is one further commuting conversion common in compilation literature that is inexpressible in our formalization of monadic form.

$$(\mathbf{if0} (\mathbf{if0} e_1 e_2) e_4 e_5) \equiv (\mathbf{if0} e (\mathbf{if0} e_1 e_4 e_5) (\mathbf{if0} e_2 e_4 e_5)) \quad \text{(Case-of-Case)}$$

Neither side is valid in our monadic form, since **if** must branch on a value, but **if** is a computation. Both monadic form and ANF normalize this equation; we return to it in [Subsection 6.2](#), as it's important for optimization.

We can derive monadic form as a normal form as follows. First, we modify the definition of evaluation contexts. ANF was defined in a call-by-value¹ setting and, following CPS, concerned with internalizing the evaluation order. Monadic form is different: it (non-strictly) expresses composing effectful computations, ordering *only* the effects. This is reflected directly in the monad laws: to reach monadic form, we should not force terms to either side of the laws, and both computations and values are allowed on the right-hand side of **let**.

We define the non-strict monadic evaluation contexts E^b , and the B -reductions for them, in [Figure 4](#). We omit **let**, but otherwise E^b is the same as E . The **let** evaluation context in ANF is responsible for normalizing commuting conversions, which each have a non-value in the right-hand side of a **let**. Recall [Listing 10](#) has a conditional expression in the strict evaluation context $(\mathbf{let} (x \cdot) e)$. We define the set $B = \{B_1, B_2, B_3\}$ of monadic reductions. These force a non-value in evaluation position to a value, just as like A -reductions, but using monadic evaluation contexts. B_1 and B_3 are essentially unchanged, but we drop one now-unnecessary termination side condition

¹ANF doesn't require strict evaluation; e.g., let could be non-strict in the operational semantics for terms in ANF. But ANF is normal with respect to A , which is defined using strict evaluation contexts.

$$\begin{aligned}
E^b & ::= \cdot \mid (\mathbf{if0} \ E^b \ e \ e) \mid (E^b \ e) \mid (\nu \ E^b) \mid (op \ \vec{v} \ E^b \ \vec{e}) \\
E^b[(\mathbf{let} \ (x \ e_1) \ e_2)] & \longrightarrow (\mathbf{let} \ (x \ e_1) \ E^b[e_2]) & B_1 \quad \text{where } E^b \neq \cdot \\
E^b[(\mathbf{if0} \ \nu \ e_1 \ e_2)] & \longrightarrow (\mathbf{let} \ (x' \ (\mathbf{if0} \ \nu \ e_1 \ e_2)) \ E^b[x']) & B_2 \quad \text{where } E^b \neq \cdot, \text{ fresh } x' \\
E^b[(O \ \vec{v})] & \longrightarrow (\mathbf{let} \ (x' \ (O \ \vec{v})) \ E^b[x']) & B_3 \quad \text{where } E \neq \cdot, \text{ fresh } x'
\end{aligned}$$

Fig. 4. B-normalizations for λ -calculus

restricting E^b . These two rules explicitly force operands to values via bind. B_2 is an optimization of A_2 to avoid code duplication. B_2 produces a monadic form expression, but not an expression in A-normal form. While we could use A_2 in B (they produce equal terms by [Equation Commute](#)), B_2 avoids the problems of A_2 and is more faithful to monadic form since it *avoids* unnecessarily normalizing the equation.

Some B-normal forms are A-normal, but not all, and all A-normal forms are B-normal. Consider our examples from earlier. [Listing 8](#) has the same B-normal form and A-normal form, given in [Listing 9](#). The second example, [Listing 10](#), has a different B-normal form. Its A-normal form ([Listing 11](#)) included exponential code duplication, but its B-normal form ([Listing 13](#)) does not. Instead, the B-normal form normalizes [Equation Case-of-Case](#) by introducing an intermediate **let**, which is required by monadic form, but does not normalize [Equation Associativity](#) or [Equation Commute](#).

(**let** (x (**if0** (**if0** (**if0** 0 0 1) 0 1) 0 1) 0 1))
LARGE)

Listing 12. Nested Branching λ -calculus

(**let** (x (**let** (x₁ (**if0** 0 0 1))
(**let** (x₂ (**if0** x₁ 0 1))
(**if0** x₂ 0 1))))
LARGE)

Listing 13. B-normalized

(**if0** 0 (**let** (x₁ 0) (**if0** x₁ (**let** (x₂ 0) (**if0** x₂ (**let** (x 0) LARGE) (**let** (x 1) LARGE))))
(**let** (x₂ 1) (**if0** x₂ (**let** (x 0) LARGE) (**let** (x 1) LARGE))))))
(**let** (x₁ 1) (**if0** x₁ (**let** (x₂ 0) (**if0** x₂ (**let** (x 0) LARGE) (**let** (x 1) LARGE))))
(**let** (x₂ 1) (**if0** x₂ (**let** (x 0) LARGE) (**let** (x 1) LARGE))))))

Listing 14. B- then A-normalized

B-normal form avoids code duplication by binding all computations, including **let** and **if**, rather than only binding values and primitive operations. This interpretation makes sense monadically, but violates the ANF discipline, which seeks to make the order of evaluation of the inner **lets** syntactically explicit by lifting them.

The specification of ANF and monadic form as a reduction system is useful for formalization. It enables separating A-normal form into its component pieces. As we will see, the reduction systems also serve as specifications so that we can prove certain properties of *any* compiler that targets the normal forms.

Now that we know, formally, what these normal forms are, let us turn to their intensional properties in the form of their machine semantics.

3 MACHINE SEMANTICS OF ANF

In this section, we introduce abstract machines for ANF and monadic form. We formalize two important facts: (1) ANF is an optimization, but (2) ANF is unsafe for scope. We show that A-normalization optimizes stack usage—there are fewer frames in use after A-normalization compared

344	Frame	$F ::= (\cdot e) \mid (v \cdot) \mid (\mathbf{let} (x \cdot) e) \mid (\mathbf{if0} \cdot e e) \mid (op \vec{v} \cdot \vec{e})$	
345	Kontinuation (as stack of frames)	$K ::= \mathbf{mt} \mid F :: K$	
346	$e; K \rightarrow_\lambda e; K$		
347			
348	$((\lambda (x) e) v); K$	$\rightarrow_\lambda e[x := v]; K$	β
349	$(\mathbf{let} (x v) e); K$	$\rightarrow_\lambda e[x := v]; K$	ζ
350	$(\mathbf{if0} 0 e_1 e_2); K$	$\rightarrow_\lambda e_1; K$	IFZ
351	$(\mathbf{if0} v e_1 e_2); K$	$\rightarrow_\lambda e_2; K$	IFNZ where $v \neq 0$
352	$(op \vec{v}); K$	$\rightarrow_\lambda \delta[[op \vec{v}]]; K$	PRIMOP
353	$v; (F :: K)$	$\rightarrow_\lambda F[v]; K$	POP
354	$F[e]; K$	$\rightarrow_\lambda e; (F :: K)$	PUSH where $e \neq v$
355			

Fig. 5. λ -calculus CK Machine

to monadic form, in general. However, we then introduce a region calculus with scoped regions, and show that A -normalization of the region calculus results in more regions and longer region lifetimes than in monadic form, in general. We call this *unsafe for scope*, in the sense of unsafe for space [19]: A -normalization is unsafe for scope-based effects.

3.1 ANF vs λ Machines

We start with a CK machine [5–7] for arbitrary λ -calculus expressions, defined in Figure 5. An expression e represents the code pointer. Intuitively K is the rest of the computation, *i.e.*, the kontinuation, but we represent it as a stack of frames F to better study the effect of ANF on the control stack.

The machine transitions are completely standard. β -reduction performs capture-avoiding substitution as a metafunction over terms, replacing a variable by a value. Primitive operators evaluate by some denotation defined by $\delta[[_]]$. Non-redexes push, and values pop, until the machine terminates with a value and an empty kontinuation \mathbf{mt} . The rule PUSH is deceptively complex; it parses a term into a frame F and subterm e , an operation most real machine do not support directly.

Our example from Listing 8 evaluates in the CK machine are follows.

	$(+ (+ 2 2) (\mathbf{let} (x 1) (f x))); \mathbf{mt}$		
375	$\rightarrow_\lambda (+ 2 2);$	$((+ \cdot (\mathbf{let} (x 1) (f x)))::\mathbf{mt})$	PUSH
376	$\rightarrow_\lambda 4;$	$((+ \cdot (\mathbf{let} (x 1) (f x)))::\mathbf{mt})$	PRIMOP
377	$\rightarrow_\lambda (+ 4 (\mathbf{let} (x 1) (f x)));$	\mathbf{mt}	POP
378	$\rightarrow_\lambda (\mathbf{let} (x 1) (f x));$	$((+ 4 \cdot)::\mathbf{mt})$	PUSH
379	$\rightarrow_\lambda (f 1);$	$(+ 4 \cdot)::\mathbf{mt}$	ζ
380			

After, A -normalization (Listing 9), (1) all computations are either on the right-hand side of a \mathbf{let} , so we can combine the rules PRIMOP and ζ , or in tail position (the injection from N into M)² so can be evaluated in place without a PUSH; (2) all operands are values, so evaluating an operator need never push a frame. Both (1) and (2) imply the only place a frame can be pushed is in a non-tail function call.

Therefore, we can define the *optimized* and *simpler* machine for A - and B -normal forms, *i.e.* monadic forms, in Figure 6. We slightly abuse notation to define two different machines \rightarrow_{anf} for A -normal forms, and \rightarrow_{bnf} for B -normal forms. Since \rightarrow_{anf} is a subset of \rightarrow_{bnf} (except for differences in metavariables), we present only one of the \rightarrow_{bnf} rules formally. The machine is optimized in

²This makes tail calls a semantically distinct concept in this machine, and not an optimization.

$C; K \rightarrow_{nf} C; K$		
$(\mathbf{let} (x' ((\lambda (x) M) V)) M'); K$	$\rightarrow_{nf} M[x := V]; (\mathbf{let} (x' \cdot) M'):: K$	CALL
$V; (F :: K)$	$\rightarrow_{nf} F[V]; K$	RETURN
$((\lambda (x) M) V); K$	$\rightarrow_{nf} M[x := V]; K$	TAIL-CALL
$(\mathbf{let} (x V) M); K$	$\rightarrow_{nf} M[x := V]; K$	MOVE
$(\mathbf{if0} 0 M_1 M_2); K$	$\rightarrow_{nf} M_1; K$	IFZ
$(\mathbf{if0} V M_1 M_2); K$	$\rightarrow_{nf} M_2; K$	IFNZ where $V \neq 0$
$(\mathbf{let} (x (op \vec{V})) M); K$	$\rightarrow_{nf} M[x := \delta[op \vec{V}]]; K$	PRIMOP
$(op \vec{V}); K$	$\rightarrow_{nf} \delta[op \vec{V}]; K$	TAIL-PRIMOP
$(\mathbf{let} (x C_1) C_2); K$	$\rightarrow_{bnf} C_1; (\mathbf{let} (x \cdot) C_2):: K$	BIND where $C_1 \neq N$

Fig. 6. A- and B-normal forms CK Machines

the sense that it requires strictly fewer transitions to evaluate ANF terms compared to the λ CK machine, and simpler in the sense that the machine never needs to decompose a term into frame and expression but looks only at the top-level computation, which maps well to a real machine. In fact, the machine is so straightforward, one might imagine that it could be interpreted as specifying a low-level register-transfer language; we do that in Section 4, after studying the high-level ANF machine in more detail. The BIND rule evaluates monadic terms that are not A-normal. For non-A-normal terms, the RETURN rule is both the return instruction for calls, but also the monadic return.

Now, Listing 9 evaluates as follows.

$(\mathbf{let} (x_1 (+ 2 2)) (\mathbf{let} (x 1) (\mathbf{let} (x_2 (f x)) (+ x_1 x_2))));$	\mathbf{mt}	
$\rightarrow_{anf} (\mathbf{let} (x 1) (\mathbf{let} (x_2 (f x)) (+ 4 x_2)));$	\mathbf{mt}	PRIMOP
$\rightarrow_{anf} (\mathbf{let} (x_2 (f 1)) (+ 4 x_2));$	\mathbf{mt}	MOVE

The intermediate push and pop transitions are eliminated.

Monadic form does not remove all intermediate push and pop transitions, in general. It does reduce some stack usage; recall that the above example is also B-normal, so B-normalization also eliminates all of its stack usage. However, monadic form still supports binding non-trivial computations, such as $(\mathbf{let} (x (\mathbf{if0} 0 0 1)) C)$ (as in Listing 13), and so the BIND rule is needed and pushes a frame. Only the optimization (2) applies in a B-normal machine, but not optimization (1).

We formalize this in Theorem 3.1. We define the metafunction MAX-STACK on machine traces \mathcal{D} as follows, and prove that A-normalizing any B-normal form optimizes stack usage.

$\text{MAX-STACK}[_]: (\mathcal{D} : (C; K \rightarrow_{nf}^* C; K)) \rightarrow \mathbb{N}$
$\text{MAX-STACK}[C; K \rightarrow_{nf}^0 C; K] = 0$
$\text{MAX-STACK}[C; K \rightarrow_{nf} C'; K' \rightarrow_{nf}^* C''; K''] = \text{MAX}[\text{LEN}[K], \text{MAX-STACK}[C'; K' \rightarrow_{nf}^* C''; K'']]$

THEOREM 3.1 (A-NORMALIZATION OPTIMIZES THE STACK (KONTINUATION)). *If $e \rightarrow_A^* M$ where $\mathcal{D}_B : (e; \mathbf{mt} \rightarrow_{bnf}^* v; \mathbf{mt})$ and $\mathcal{D}_A : (M; \mathbf{mt} \rightarrow_{anf}^* v; \mathbf{mt})$, then trace \mathcal{D}_A uses no more frames than \mathcal{D}_B . That is, $\text{MAX-STACK}[\mathcal{D}_B] \geq \text{MAX-STACK}[\mathcal{D}_A]$. Furthermore, there exists a program C for which a trace \mathcal{D}_A uses strictly fewer frames.*

PROOF. The proof is straightforward by induction on the trace \mathcal{D}_B . When \mathcal{D}_B takes a BIND step, C must take an A-reduction, and the A-reduction on the C will either not increase or will decrease the $\text{MAX-STACK}[_]$. \square

442 $r \in \text{Regions}$ $o \in \text{Offset}$
 443 $e ::= \dots \mid (\mathbf{letregion} \ r \ e) \mid (@ \ r \ e)$ $a ::= (r . o)$
 444 $E ::= \dots \mid (@ \ r \ E)$ $S ::= \emptyset \mid (S, a \mapsto v)$
 445 $F ::= \dots \mid (\mathbf{free} \ r) \mid (@ \ r \cdot)$

Fig. 7. λ^r Syntax

450 This is more than a theoretical effect. The CertiCoq compiler has had two backends, a CPS and an
 451 ANF backend, and ANF outperforms CPS [16–18]. Paraskevopoulou [16] attributes this to additional
 452 heap allocation in CPS, noting that “sophisticated techniques to reduce heap allocation” could be
 453 used in the CPS backend. CPS can result in additional heap allocation by allocating a continuation
 454 for local control. Each frame in the CK machine corresponds to a continuation, and could cause
 455 heap allocation of a closure in CPS. The ANF backend, by contrast, avoids this allocation, since
 456 these frames are made explicit syntactically and never allocate. This is no more than we should
 457 expect; A -normalization was partially derived from eliminating administrative redexes introduced
 458 by CPS translation. This suggests that merely using ANF simplifies control-related optimization
 459 compared to CPS, at least in the context of a fully mechanically verified compiler.

3.2 A Region Calculus and Machine

462 Unfortunately, A -normalization is an unsafe optimization in some settings. In general, A -normalization
 463 is *unsafe for scope*: lexical scope is extended, changing variable extent, lifetimes, and other effects related
 464 to lexical scope. Extensionally, such as in a pure calculus when reasoning up to α -equivalence,
 465 this is irrelevant. However, it matters when reasoning about intensional properties or effectful
 466 calculi. We consider a version of the scoped regions of Tofte and Talpin [20].

467 We first extend our λ -calculus with lexically scoped regions and region allocation; the syntax is
 468 given in Figure 7. The expression $(\mathbf{letregion} \ r \ e)$ runs e with a new region r allocated in the store
 469 S^3 . After e returns, the region r is freed. In this calculus, all values must be explicitly allocated in
 470 a region in the store. The form $(@ \ r \ e)$ evaluates to a reference, by evaluating e to a value that is
 471 stored in region r . All primitive operations and value must be explicitly allocated; the result of a
 472 function must also be allocated before it returns. We also extend evaluation contexts, for A - and
 473 B -normalization of this syntax, and add a new frames for the machine semantics. The new frame
 474 $(\mathbf{free} \ r)$ is unusual; it is not a frame of an evaluation context, but instead directs the machine to
 475 free region r .

476 In Figure 8, we give a CSK machine [5–7] for λ^r . The machine includes a straightforward
 477 extension of the transitions from the λ CK machine in Figure 5. The machine is unrealistic in that
 478 all values are passed by references, whereas a realistic machine might avoid boxing word-sized
 479 data. However, this detail is irrelevant for our purposes: some ANF terms bind non-word sized data,
 480 so *any* extension of a region lifetime is unsafe for space.

481 The machine is essentially similar to the λ CK machine, with two main differences: the represen-
 482 tation of run-time values, and the new transitions for regions. The new region-related transitions
 483 are given below the line. Run-time values are addresses, a pair $(r . o)$ of a region and an offset into
 484 that region, so all computations dereference an address a in the store S , written $S(a)$. Syntactic
 485 values must be allocated in a region, indicated by the frame, in the \mathbf{ALLOC} transition. $(\mathbf{letregion} \ r \ e)$

487 ³The allocation pattern forms a (data) stack (of regions), as explained by Tofte and Talpin [20]. The stack of regions is
 488 irrelevant to us, while the control stack (kontinuation) K is important, so we refer to allocation as occurring in the heap,
 489 and use “stack” to refer to the control stack.

$C; S; K \rightarrow_{\lambda^r} C; S; K$			
$(a_1 a_2); S; K$	\rightarrow_{λ^r}	$e[x := a_2]; S; K$	CALL where $(\lambda(x) e) = S(a_1)$
$(\mathbf{let} (x a) e); S; K$	\rightarrow_{λ^r}	$e[x := a]; S; K$	MOVE
$(\mathbf{if0} a e_1 e_2); S; K$	\rightarrow_{λ^r}	$e_1; S; K$	IFZ where $0 = S(a)$
$(\mathbf{if0} a e_1 e_2); S; K$	\rightarrow_{λ^r}	$e_2; S; K$	IFNZ where $0 \neq S(a)$
$(op \vec{a}); S; K$	\rightarrow_{λ^r}	$\llbracket op S(\vec{a}) \rrbracket; S; K$	PRIMOP
$F[e]; S; K$	\rightarrow_{λ^r}	$e; S; (F :: K)$	PUSH where $e \neq a$
$a; S; (F :: K)$	\rightarrow_{λ^r}	$F[S(a)]; S; K$	POP
<hr/>			
$(\mathbf{letregion} r e); S; K$	\rightarrow_{λ^r}	$e; S; ((\mathbf{free} r) :: K)$	RALLOC
$a; S; ((\mathbf{free} r) :: K)$	\rightarrow_{λ^r}	$a; \mathbf{FREE}\llbracket S, r \rrbracket; K$	RFREE
$v; S; ((@ r \cdot) :: K)$	\rightarrow_{λ^r}	$(r \cdot o); (S[(r.o) := v]; K)$	ALLOC fresh o

Fig. 8. λ^r CSK Machine

(implicitly) allocates a new region, and adds a **free** frame to the stack. The metafunction $\mathbf{FREE}\llbracket S, r \rrbracket$ deallocates all addresses with the region r in S .

Consider the following example term's evaluation in this machine. This example is an λ^r equivalent of $(* (* 1 2) (* 3 4))$, with region sizes and lifetimes minimized. r_0 is an initial region in which the final result is allocated.

$(\mathbf{letregion} r_2$		\emptyset	mt
$(@ r_0 (* (\mathbf{letregion} r_1 (@ r_2 (* (@ r_1 1) (@ r_1 2))))$			
$(\mathbf{letregion} r_3 (@ r_2 (* (@ r_3 3) (@ r_3 4))))$			
$\rightarrow_{\lambda^r} (@ r_0 (* (\mathbf{letregion} r_1 (@ r_2 (* (@ r_1 1) (@ r_1 2))))$		\emptyset	$(\mathbf{free} r_2) :: \mathbf{mt}$
$(\mathbf{letregion} r_3 (@ r_2 (* (@ r_3 3) (@ r_3 4))))$			
$\rightarrow_{\lambda^r} (* (\mathbf{letregion} r_1 (@ r_2 (* (@ r_1 1) (@ r_1 2))))$		\emptyset	$(@ r_0 \cdot) :: (\mathbf{free} r_2) :: \mathbf{mt}$
$(\mathbf{letregion} r_3 (@ r_2 (* (@ r_3 3) (@ r_3 4))))$			
$\rightarrow_{\lambda^r} \dots$			
$\rightarrow_{\lambda^r} (* (r_2.o_1) (r_2.o_2))$	$\left[\begin{array}{l} [(r_2.o_2) \mapsto 2] [(r_2.o_1) \mapsto 12] \\ [(r_2.o_2) \mapsto 2] [(r_2.o_1) \mapsto 12] \\ [(r_2.o_2) \mapsto 2] [(r_2.o_1) \mapsto 12] [(r_0.o_3) \mapsto 24] \\ [(r_0.o_3) \mapsto 24] \end{array} \right.$		$(@ r_0 \cdot) :: (\mathbf{free} r_2) :: \mathbf{mt}$
$\rightarrow_{\lambda^r} 24$			$(@ r_0 \cdot) :: (\mathbf{free} r_2) :: \mathbf{mt}$
$\rightarrow_{\lambda^r} (r_0.o_3)$			$(\mathbf{free} r_2) :: \mathbf{mt}$
$\rightarrow_{\lambda^r} (r_0.o_3)$			mt

Using this abstract machine, we can measure some simple allocation behaviour, such as:

- (1) What is the maximum number of regions live at once ($\mathbf{MAX-REGIONS}\llbracket _ \rrbracket$)?
- (2) What is the maximum number of live addresses, in all regions, at once ($\mathbf{MAX-MEMORY}\llbracket _ \rrbracket$)?

We omit their formal definitions, but these metafunctions are essentially similar to definition of $\mathbf{MAX-STACK}\llbracket _ \rrbracket$: they take a trace and measuring the maximum value in the trace. For the above example, the $\mathbf{MAX-REGIONS}\llbracket _ \rrbracket$ of the trace is 3 and the $\mathbf{MAX-MEMORY}\llbracket _ \rrbracket$ is 4.

We extend A-normalization to λ^r with the following two A-reductions.

$E[(\mathbf{letregion} r e)]$	\xrightarrow{A}	$(\mathbf{letregion} r E[e])$	A_4
		where $E \neq \cdot$	
$E[(@ r N)]$	\xrightarrow{A}	$(\mathbf{let} (x' (@ r N)) E[x'])$	A_5
		where $E \neq \cdot$, $E \neq E'[(\mathbf{let} (x \cdot) e')]$, $E \neq E'[(@ r \cdot)]$, and fresh x' ,	

The side condition $E \neq E' [(@ r \cdot)]$ must also be added to the A_3 reduction. The reduction system is fully implemented in the anonymous supplementary materials.

A_4 extends the lifetime of region r , but this is necessary for A_3 from Figure 1 to lift an intermediate computation. Concretely, A -normalizing the running example, we get the term in Listing 15, which runs with $\text{MAX-REGIONS}[\llbracket _ \rrbracket]$ 4 and a $\text{MAX-MEMORY}[\llbracket _ \rrbracket]$ of 7, compared to the original 3 and 4, respectively.

```

547 (letregion r2
548   (letregion r1
549     (let (x (@ r1 1))
550       (let (x1 (@ r1 2))
551         (let (x2 (@ r2 (* x x1)))
552           (letregion r3
553             (let (x3 (@ r3 3))
554               (let (x4 (@ r3 4))
555                 (let (x5 (@ r2 (* x3 x4))
556                   (@ r0 (* x2 x5))))))))))

```

Listing 15. A -normalization of λ^r Example

```

547 (letregion r2
548   (let (x4 (letregion r1
549     (let (x2 (@ r1 1))
550       (let (x3 (@ r1 2))
551         (@ r2 (* x2 x3))))))
552     (let (x5 (letregion r3
553       (let (x (@ r3 3))
554         (let (x1 (@ r3 4))
555           (@ r2 (* x x1))))))
556       (@ r0 (* x4 x5))))))

```

Listing 16. B -normalization of λ^r Example

THEOREM 3.2 (A-NORMALIZATION IS UNSAFE FOR SCOPE).

If $C \xrightarrow{*} M$, $\mathcal{D}_B : (C; \emptyset; \mathbf{mt} \xrightarrow{*}_{\lambda^r} a; S'; \mathbf{mt})$ and $\mathcal{D}_A : (M; \emptyset; \mathbf{mt} \xrightarrow{*}_{\lambda^r} a; S''; \mathbf{mt})$, then $\text{MAX-REGIONS}[\llbracket \mathcal{D}_A \rrbracket] \geq \text{MAX-REGIONS}[\llbracket \mathcal{D}_B \rrbracket]$, and $\text{MAX-MEMORY}[\llbracket \mathcal{D}_A \rrbracket] \geq \text{MAX-MEMORY}[\llbracket \mathcal{D}_B \rrbracket]$. Furthermore, there exists C and \mathcal{D}_A such that $\text{MAX-MEMORY}[\llbracket \mathcal{D}_A \rrbracket] > \text{MAX-MEMORY}[\llbracket \mathcal{D}_B \rrbracket]$ and $\text{MAX-REGIONS}[\llbracket \mathcal{D}_A \rrbracket] > \text{MAX-REGIONS}[\llbracket \mathcal{D}_B \rrbracket]$.

By contrast, we can extend B -normalization to λ^r with the following reductions.

$$\begin{array}{lcl}
 E^b [(\mathbf{letregion} \ r \ e)] & \xrightarrow{B} & \begin{array}{l} \dots \\ (\mathbf{let} \ (x \ (\mathbf{letregion} \ r \ e)) \ E^b [x]) \\ E^b \neq \cdot \end{array} & B_4 \\
 E^b [(@ \ r \ N)] & \xrightarrow{B} & \begin{array}{l} (\mathbf{let} \ (x' \ (@ \ r \ N)) \ E^b [x']) \\ E^b \neq \cdot, E^b \neq E_1^b [(\mathbf{let} \ (x \cdot) \ e)], E^b \neq E_2^b [(@ \ r \cdot)], \text{fresh } x' \end{array} & B_5
 \end{array}$$

Since $\mathbf{letregion}$ is effectful, B_4 binds it instead of lifting it, preserving the original lifetime. The B -normalization of our running example is given in Listing 16.

THEOREM 3.3 (B-NORMALIZATION IS SAFE-FOR-SCOPE).

If $e \xrightarrow{*} C$, $\mathcal{D}_\lambda : (e; \emptyset; \mathbf{mt} \xrightarrow{*}_{\lambda^r} a; S'; \mathbf{mt})$ and $\mathcal{D}_B : (C; \emptyset; \mathbf{mt} \xrightarrow{*}_{\lambda^r} a; S''; \mathbf{mt})$, then $\text{MAX-REGIONS}[\llbracket \mathcal{D}_\lambda \rrbracket] = \text{MAX-REGIONS}[\llbracket \mathcal{D}_B \rrbracket]$, and $\text{MAX-MEMORY}[\llbracket \mathcal{D}_\lambda \rrbracket] = \text{MAX-MEMORY}[\llbracket \mathcal{D}_B \rrbracket]$.

4 IMPERATIVE A-NORMALIZATION

The problems with normalizing commuting conversion are merely with the syntax of ANF, and not with (machine) semantics of commuting conversion. By designing a syntax based on the ANF abstract machine, we get an imperative language very similar to ANF. Working backwards, un-normalizing commuting conversions, we design an imperative monadic language in which A -normalization is safe, neither extending scope, nor duplicating continuations.

Recall from Figure 6 that each transition in the monadic CK machine uniquely maps to an A -normal form expression. In Figure 9, we present an imperative A -normal form designed from the transitions in the abstract machine, with \mathbf{begin} for sequencing imperative statements.

589 $t ::= (\mathbf{begin} \vec{s} t) \mid v \mid (\mathbf{if0} v t t) \mid (\mathbf{call} v \vec{v}) \mid (op \vec{v})$
 590 $s ::= (\mathbf{begin} \vec{s}) \mid (\mathbf{set!} x v) \mid (\mathbf{set!} x (op \vec{v})) \mid (\mathbf{set!} x (\mathbf{call} v \vec{v}))$
 591 $v ::= (\lambda (x) t) \mid \iota$
 592
 593
 594
 595

Fig. 9. Imperative ANF Syntax (AB-normal form)

596 Intuitively, a program is a sequence of statements s followed by a tail-position operation (tail)
 597 t producing the final value. We interpret the ANF **let** as an effectful **set!** operation; rather than
 598 performing substitution, we set a variable in the machine state. Tails t include operations for
 599 machine transitions that correspond to all expression in M position in ANF. For example, the
 600 TAIL-CALL transition is realized by the **call** statement in tail position. All non-tail expressions in
 601 ANF must appear on the right-hand side of a **let** in ANF, and therefore appear on the right-hand side
 602 of **set!** in imperative ANF. These transitions—CALL, MOVE, and PRIMOP—are realized by statements
 603 s . The result is similar to a register-transfer language, although one with an infinite set of registers
 604 x , higher-order functions, and primitive call and return.

605 The translation from ANF into this syntax is straightforward; we give a definition later in
 606 Figure 18a. But for now, we want to work backwards from this syntax to imperative monadic form,
 607 and demonstrate that an imperative interpretation of A-normalization within imperative monadic
 608 form is safe and easy. So what must imperative monadic form be?

609 If **set!** is **let**, then we also need to support **(set! x t)**, that is, a tail computation on the right-hand
 610 side an assignment. This corresponds to **(let (x C) C)** in monadic form. We see this also by doing
 611 code generation of Equation Associativity and Equation Commute, which would look like the
 612 following.
 613

$$\begin{aligned}
 &(\mathbf{begin} (\mathbf{set!} y (\mathbf{begin} (\mathbf{set!} x t) t)) t) \equiv (\mathbf{begin} (\mathbf{set!} x t) (\mathbf{set!} y t) t) && \text{(Imp. Associativity)} \\
 &(\mathbf{begin} (\mathbf{set!} y (\mathbf{if} v t_1 t_2)) t) \equiv (\mathbf{if0} v (\mathbf{begin} (\mathbf{set!} y t_1) t) (\mathbf{begin} (\mathbf{set!} y t_2) t)) && \text{(Imp. Commute (1))}
 \end{aligned}$$

614
 615
 616
 617
 618 To support both sides of the equations, we must allow **begin** and **if0** on the right-hand side of a
 619 **set!**. This requires one new statement.
 620

621 $t ::= (\mathbf{begin} \vec{s} t) \mid v \mid (\mathbf{if0} v t t) \mid (\mathbf{call} v \vec{v}) \mid (op \vec{v})$
 622 $s ::= (\mathbf{begin} \vec{s}) \mid (\mathbf{set!} x v) \mid (\mathbf{set!} x t) \mid (\mathbf{set!} x (op \vec{v})) \mid (\mathbf{set!} x (\mathbf{call} v \vec{v}))$
 623 $v ::= (\lambda (x) t) \mid \iota$
 624

625 Equation Imp. Commute (1) still duplicates the continuation t , but only because the syntax
 626 supports conditional *expressions*, but not conditional *statements*. If we add a conditional statement,
 627 we could rephrase the equation as the following.
 628

$$(\mathbf{begin} (\mathbf{set!} y (\mathbf{if} v t_1 t_2)) t) \equiv (\mathbf{begin} (\mathbf{if0} v (\mathbf{set!} y t_1) (\mathbf{set!} y t_2)) t) \quad \text{(Imp. Commute (2))}$$

630 We perform the commuting conversion by duplicating the assignment to x rather than duplicating
 631 the context in which x is bound. This isn't possible in (high-level) ANF, since lexical binding cannot
 632 export x , but it is possible when all lexical binding has been transformed into imperative statements.
 633

634 Adding the conditional statement and **set!** with a complex right-hand side, we get the final
 635 version of the imperative monadic form is in Figure 10.
 636

637 In this imperative monadic syntax, we can easily perform the imperative equivalent of A-
 638 normalization into imperative ANF by taking the congruent closure of the following set of reductions
 639

638 $t ::= (\mathbf{begin} \vec{s} \ t) \mid v \mid (\mathbf{if0} \ v \ t \ t) \mid (\mathbf{call} \ v \ v) \mid (op \ \vec{v})$
 639 $s ::= (\mathbf{begin} \vec{s}) \mid (\mathbf{set!} \ x \ v) \mid (\mathbf{if0} \ v \ s \ s) \mid (\mathbf{set!} \ x \ t) \mid (\mathbf{set!} \ x \ (op \ \vec{v})) \mid (\mathbf{set!} \ x \ (\mathbf{call} \ v \ v))$
 640 $v ::= (\lambda \ (x) \ t) \mid \iota$
 641
 642

Fig. 10. Imperative Monadic Syntax

644 $hv ::= \iota \mid (\mathbf{closure} \ \Sigma \ (\lambda \ (x) \ t))$ $\Sigma ::= \emptyset \mid \Sigma[x \mapsto hv]$
 645 $wv ::= x \mid \iota$ $F ::= (\mathbf{begin} \ (\mathbf{set!} \ x \cdot) \ \vec{s} \ t)$
 646 $K ::= \mathbf{mt} \mid F :: K$
 647
 648

Fig. 11. CEK Machine Syntax

650 $AB = \{AB_1, AB_2\}$, yielding AB -normal form⁴.

652 $(\mathbf{set!} \ x \ (\mathbf{if0} \ v \ t_1 \ t_2)) \longrightarrow_{AB} (\mathbf{if0} \ v \ (\mathbf{set!} \ x \ t_1) \ (\mathbf{set!} \ x \ t_2)) \quad AB_1$
 653 $(\mathbf{set!} \ x \ (\mathbf{begin} \ \vec{s} \ t)) \longrightarrow_{AB} (\mathbf{begin} \ \vec{s} \ (\mathbf{set!} \ x \ t)) \quad AB_2$
 654

655 AB -normalization is straightforward. Unlike A -normalization and B -normalization, it does not
 656 shuffle the evaluation contexts at all. It is a completely local transformation. We never need to deal
 657 with join points, or CPS'd compilers, or code duplication. We gain all the stack optimization effect
 658 on A -normalization (as we formalize next), the advantages of a register-transfer syntax for code
 659 generation and machine implementation, etc.

660 *The key idea* in AB -normalization is that lexical expressions must be elaborated into imperative
 661 statements before A -normalization.

662 4.1 Machine Semantics of AB -normal form

663 We use a CEK machine [5–7] to define a machine semantics for imperative monadic form in
 664 Figure 11 and Figure 12. We then show that AB -normalization optimizes the stack in the same way
 665 as A -normalization.

666 The environment Σ (to distinguish it from evaluation contexts) represents the register file,
 667 mapping variables to their values. We separate word values wv , which represent run-time values
 668 that can be eliminated, and heap values hv , which must be bound to a name and intuitively would
 669 be allocated in memory. We use the syntax $\Sigma(wv)$ get the value of wv , either dereferencing wv in
 670 Σ if wv is a variable, or returning wv if not. The code continues to represent the program counter,
 671 and the continuation continues to represent the call stack. W.l.o.g., we assume that all λ s appear
 672 on the right-hand side of a $\mathbf{set!}$, to simplify implementing closures. This can be implemented by
 673 another B -reduction that binds λ , treating it as a computation rather than a value (which it is, in
 674 this machine; it performs allocation).
 675

676 This machine is essentially similar to the monadic CK machine in Figure 6. The main difference
 677 is the use of an environment rather than substitution, becomes necessary with imperative inter-
 678 pretation of \mathbf{let} as $\mathbf{set!}$. We also require two administrative reductions, ADMIN1 and ADMIN2,
 679 to enable composing internally with \mathbf{begin} , although these could be normalized to simplify the final
 680 machine semantics.

681 There is one key difference: the AB transition, which cannot occur in ABNF but can in imperative
 682 monadic form. This rule, and \mathbf{CALL} for a non-tail call, both push a frame $(\mathbf{begin} \ (\mathbf{set!} \ x \cdot) \ t)$. After
 683

684 ⁴We choose the name AB for two reasons: (1) this normalization yields the best of both A - and B -normalization (2) AB -
 685 normalization follows B -normalization, which was defined second after A -normalization, so AB is 3... in big-endian binary,
 686 anyway.

687	$t; \Sigma; K \rightarrow_{CEK} t; \Sigma; K$	
688		
689	$(\mathbf{begin}(\mathbf{set!} x wv) \vec{s} t); \Sigma; K \rightarrow_{CEK} t; \Sigma[x := wv]; K$	MOVE
690	$(\mathbf{begin}(\mathbf{set!} x (\lambda(x) t)) \vec{s} t); \Sigma; K \rightarrow_{CEK} t; \Sigma[x := (\mathbf{closure} \Sigma(\lambda(x) t))]; K$	CLOSURE
691	$(\mathbf{begin}(\mathbf{set!} x (\mathbf{call} wv_1 wv_2)) \vec{s} t); \Sigma; K \rightarrow_{CEK} t; \Sigma'[y := wv_2]; (\mathbf{begin}(\mathbf{set!} x \cdot) \vec{s} t) :: K$	CALL
692		
693	$v; \Sigma; F :: K \rightarrow_{CEK} F[v]; \Sigma; K$	RETURN
694	$(\mathbf{begin}(\mathbf{if0} wv s_1 s_2) t); \Sigma; K \rightarrow_{CEK} (\mathbf{begin} s_1 t); \Sigma; K$	IFZ-S
695		where $\Sigma(wv) = 0$
696	$(\mathbf{begin}(\mathbf{if0} wv s_1 s_2) t); \Sigma; K \rightarrow_{CEK} (\mathbf{begin} s_2 t); \Sigma; K$	IFNZ-S
697		where $\Sigma(wv) \neq 0$
698	$(\mathbf{begin}(\mathbf{set!} x (op \vec{wv})) t); \Sigma; K \rightarrow_{CEK} t; \Sigma[x := \llbracket op \Sigma(sv) \rrbracket]; K$	PRIMOP
699	$(op \vec{wv}); \Sigma; K \rightarrow_{CEK} \llbracket op \Sigma(wv) \rrbracket; \Sigma; K$	TAIL-PRIMOP
700	$(\mathbf{if0} wv t_1 t_2); \Sigma; K \rightarrow_{CEK} t_1; \Sigma; K$	IFZ-T
701		where $\Sigma(wv) = 0$
702	$(\mathbf{if0} wv t_1 t_2); \Sigma; K \rightarrow_{CEK} t_2; \Sigma; K$	IFNZ-T
703		where $\Sigma(wv) \neq 0$
704	$(\mathbf{begin}(\mathbf{set!} x t_1) \vec{s} t_2); \Sigma; K \rightarrow_{CEK} t_1; \Sigma; (\mathbf{begin}(\mathbf{set!} x \cdot) \vec{s} t_2) :: K$	AB
705	$(\mathbf{begin}(\mathbf{begin} \vec{s}_1) \vec{s}_2 t); \Sigma; K \rightarrow_{CEK} (\mathbf{begin} \vec{s}_1 \vec{s}_2 t); \Sigma; K$	ADMIN1
706	$(\mathbf{begin} t); \Sigma; K \rightarrow_{CEK} t; \Sigma; K$	ADMIN2
707		

Fig. 12. CEK Machine for Imperative Monadic Form

AB-normalization, the AB rule cannot occur, so we regain the useful property that only non-tail calls push a stack frame.

It's straightforward to show that AB-normalization is an optimization of the stack, just like A-normalization.

THEOREM 4.1 (AB-NORMALIZATION OPTIMIZES THE STACK). *If $t \xrightarrow{*}_{AB} t'$ where $\mathcal{D}_B : (t; \emptyset; \mathbf{mt} \rightarrow_{CEK}^* v; S; K)$ and $\mathcal{D}_A : (t'; \emptyset; \mathbf{mt} \rightarrow_{CEK}^* v'; S'; K')$, then trace \mathcal{D}_A uses no more frames than \mathcal{D}_B . That is, $\text{MAX-STACK}[\mathcal{D}_B] \geq \text{MAX-STACK}[\mathcal{D}_A]$. Furthermore, there exists a program t for which a trace $\text{MAX-STACK}[\mathcal{D}_B] > \text{MAX-STACK}[\mathcal{D}_A]$.*

PROOF. The proof proceeds by induction on the trace \mathcal{D}_B . The two interesting cases correspond to the AB reductions; all other cases preserve max stack size. We sketch the case for $t = (\mathbf{begin}(\mathbf{set!} x (\mathbf{begin} \vec{s} t_1)) t_2)$ as a diagram; the other case is similar.

$$\begin{array}{ccc}
 (\mathbf{begin}(\mathbf{set!} x (\mathbf{begin} \vec{s} t_1)) t_2); \Sigma; K & \xrightarrow{CEK^*} & (\mathbf{begin} \vec{s} t_1); \Sigma; (\mathbf{begin}(\mathbf{set!} x \cdot) t_2) :: K \\
 \downarrow AB_2 & & \downarrow CEK-AB \\
 & & t_1; \Sigma'; K' + (\mathbf{begin}(\mathbf{set!} x \cdot) t_2) :: K \\
 & & \vdots \\
 & & > \text{max-stack} \\
 (\mathbf{begin} \vec{s} (\mathbf{set!} x t_1) t_2); \Sigma; K & \xrightarrow{CEK^*} & (\mathbf{begin}(\mathbf{set!} x t_1) t_2); \Sigma; K' + K
 \end{array}$$

Before AB-normalization, a program of this shape must push the stack frame $F = (\mathbf{begin}(\mathbf{set!} x \cdot) t_2)$, since the right-hand side of the instruction is a complex tail. But AB-normalization, via the AB_2 rule, eliminates that one stack frame by reassociating the $\mathbf{set!}$. In the machine state after evaluating the AB-normalized program, the instructions \vec{s} evaluate in a subtrace, producing a stack $K' + K$ (where $+$ is the append operation on stacks). The resulting stack is one frame smaller than $K' + F :: K$. The

$t ::= (\mathbf{begin} \vec{s} t) \mid rv \mid (\mathbf{alloc} r v) \mid (\mathbf{if0} rv t t) \mid (\mathbf{call} rv rv) \mid (\mathbf{alloc} r (op \vec{rv}))$
 $s ::= (\mathbf{begin} \vec{s}) \mid (\mathbf{if0} rv s s) \mid (\mathbf{set!} x (\mathbf{alloc} r v)) \mid (\mathbf{set!} x rv) \mid (\mathbf{set!} x t)$
 $\quad \mid (\mathbf{set!} x (\mathbf{alloc} r (op \vec{rv}))) \mid (\mathbf{set!} x (\mathbf{call} rv rv)) \mid (\mathbf{ralloc} r) \mid (\mathbf{rfree} r)$
 $v ::= \iota \mid x \mid a \mid (\lambda (x) t) \quad hv ::= (\mathbf{closure} \Sigma (\lambda (x) t) \mid wv)$
 $a ::= (r.o) \quad F ::= (\mathbf{begin} (\mathbf{set!} x \cdot) \vec{s} t)$
 $rv ::= a \mid x \quad S ::= \emptyset \mid S[a \mapsto hv]$
 $wv ::= \iota \mid a \quad \Sigma ::= \emptyset \mid \Sigma[x := a]$

Fig. 13. Imperative Monadic w/ Regions Syntax

result follows by the induction hypotheses, which guarantees that the subtraces for \vec{s} , t_1 , and t_2 do not increase the stack size. \square

4.2 AB-normalization and Regions

Recall our key idea claims that we must elaborate lexical expressions into imperative statements to solve A -normalization. Let us test this idea against lexically scoped regions.

Unlike **let**⁵, **letregion** has an effect at the beginning and end of its scope. We can compile lexical regions using something like the following code generator.

$$\begin{aligned}
 \mathbf{CG}[_] &: C \rightarrow t \\
 \mathbf{CG}[(\mathbf{letregion} r C)] &= (\mathbf{begin} (\mathbf{ralloc} r) (\mathbf{set!} x \mathbf{CG}[C]) (\mathbf{rfree} r) x) \\
 \mathbf{CG}[(\mathbf{@} r N)] &= (\mathbf{begin} (\mathbf{set!} x (\mathbf{alloc} r \mathbf{CG}[N])) x)
 \end{aligned}$$

This syntax begins to look suspiciously like monadic regions [9, 10].

We extend the CEK machine into a CESK machine [5–7], using the store to model regions. The machine essentially smashes together the previous two machines. We give only the key rules, for brevity, but a complete implementation is available in the anonymous supplementary materials.

We extend our imperative monadic syntax with imperative regions in Figure 13. There are two main differences. First, the additional instructions for **ralloc** and **rfree**, which correspond to the **RALLOC** and **RFREE** transitions of the CSK machine for the monadic region calculus. Second, as all values must be passed by reference, all the prior value positions and primitive operators are wrapped in **alloc**, which corresponds to the **ALLOC** transition of the CSK machine in Figure 8. To ensure all values are passed by reference, all value operands are now register values rv —either a variable x or an address. Heap values are now explicitly allocated in the store.

We define the key rules CESK machine in Figure 14. The **MOVE** transition dereferences a register value rv , either reading the value from a register x or if rv is an a using it directly. Then the register file is updated to map x to that value. The **ALLOC** transition allocates a word value wv in memory to a fresh address in region r , updating the register file with x mapped to that address. There are analogous instructions for allocating closures and allocating the result of a primop, as well as analogous instructions for allocating in tail position. The **RALLOC** transition allocates the region r , which in this machine is a no-op. The **RFREE** transition frees a region. These two corresponds to the transitions of the same name in the CSK machine, but as they are now instructions, they do not use the stack. The only transitions that push a stack frame are the non-tail **CALL** transition, and the **AB** transition for not AB -normal terms.

⁵Strictly, there is an effect at the end of **let**'s scope: the variable is dead. We could introduce a marker for the end-of-lifetime of a variable into the imperative language. But this sort of information is simple to infer via liveness analysis.

$t; \Sigma; S; K \rightarrow_{CESK} t; \Sigma; S; K$		
785		
786		
787	...	
788	$(\mathbf{begin} (\mathbf{set!} \ x \ rv) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} (\mathbf{begin} \ \vec{s} \ t); \Sigma[x := \Sigma(rv)]; S; K$	MOVE
789	$(\mathbf{begin} (\mathbf{set!} \ x (\mathbf{alloc} \ r \ wv)) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} (\mathbf{begin} \ \vec{s} \ t); \Sigma[x := (r.o)]; S[(r.o) := v]; K$	ALLOC
790		fresh o
791	$(\mathbf{begin} (\mathbf{ralloc} \ r) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} (\mathbf{begin} \ \vec{s} \ t); \Sigma; S; K$	RALLOC
792	$(\mathbf{begin} (\mathbf{rfree} \ r) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} (\mathbf{begin} \ \vec{s} \ t); \Sigma; \text{free}(S, r); K$	RFREE
793	$(\mathbf{begin} (\mathbf{set!} \ x (\mathbf{call} \ rv_1 \ rv_2)) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} t_1; \Sigma_1[y := rv_2]; S; ((\mathbf{begin} (\mathbf{set!} \ x \cdot) \ \vec{s} \ t):: K)$	CALL
794		where $(\mathbf{closure} \ \Sigma_1 (\lambda (y) \ t_1)) = S(\Sigma(rv_1))$
795	$(\mathbf{begin} (\mathbf{set!} \ x \ t_1) \ \vec{s} \ t); \Sigma; S; K \rightarrow_{CESK} t_1; \Sigma; S; K$	AB
796		

Fig. 14. CESK Machine for Imperative Monadic w/ Regions

With this machine, we can now *AB*-normalize and run our earlier region example *without* extending lifetimes, but *with* optimized stack usage. Recall from Listing 15 that *A*-normalization increased $\text{MAX-REGIONS}[\llbracket _ \rrbracket]$ and $\text{MAX-MEMORY}[\llbracket _ \rrbracket]$. In Listing 17, we perform code generation followed by *AB*-normalization on the monadic example from Listing 16. In Listing 17, we perform code generation of the ANF example from Listing 15.

<pre> 804 1 (begin 805 2 (ralloc r2) 806 3 (ralloc r1) 807 4 (set! x2 (alloc r1 1)) 808 5 (set! x3 (alloc r1 2)) 809 6 (set! xt1 (alloc r2 (* x2 x3))) 810 7 (rfree r1) 811 8 (set! x4 xt1) 812 9 (ralloc r3) 813 10 (set! x (alloc r3 3)) 814 11 (set! x1 (alloc r3 4)) 815 12 (set! xt2 (alloc r2 (* x x1))) 816 13 (rfree r3) 817 14 (set! x5 xt2) 818 15 (set! xt3 (alloc r0 (* x4 x5))) 819 16 (rfree r2) 820 17 xt3) </pre>	<pre> 1 (begin 2 (ralloc r2) 3 (ralloc r1) 4 (set! x (alloc r1 1)) 5 (set! x1 (alloc r1 2)) 6 (set! x2 (alloc r2 (* x x1))) 7 (ralloc r3) 8 (set! x3 (alloc r3 3)) 9 (set! x4 (alloc r3 4)) 10 (set! x5 (alloc r2 (* x3 x4))) 11 (set! xt1 (alloc r0 (* x2 x5))) 12 (rfree r3) 13 (set! xt2 xt1) 14 (rfree r1) 15 (set! xt3 xt2) 16 (rfree r2) 17 xt3) </pre>
---	---

Listing 17. AB-Normalized · Code Gen · Monadic

Listing 17. Code Generation · ANF

The lifetime of $r1$ is extended by *A*-normalization, but not *B*-normalization or *AB*-normalization, and *AB*-normalization has eliminated all the stack usage of monadic form. The formal proof is similar to that of Theorem 3.2; note that *AB*-normalization does not change the order of instructions.

THEOREM 4.2 (AB-NORMALIZATION IS SAFE-FOR-SCOPE).

If $t_M \xrightarrow{*}_{AB} t_{AB}$, $\mathcal{D}_M : (t_M; \emptyset; \emptyset; \mathbf{mt} \rightarrow^*_{CESK} a; \Sigma; S; \mathbf{mt})$ and $\mathcal{D}_{AB} : (t_{AB}; \emptyset; \emptyset; \mathbf{mt} \rightarrow^*_{CESK} a; \Sigma'; S'; \mathbf{mt})$, then $\text{MAX-REGIONS}[\llbracket \mathcal{D}_M \rrbracket] = \text{MAX-REGIONS}[\llbracket \mathcal{D}_{AB} \rrbracket]$, and $\text{MAX-MEMORY}[\llbracket \mathcal{D}_M \rrbracket] = \text{MAX-MEMORY}[\llbracket \mathcal{D}_{AB} \rrbracket]$, and

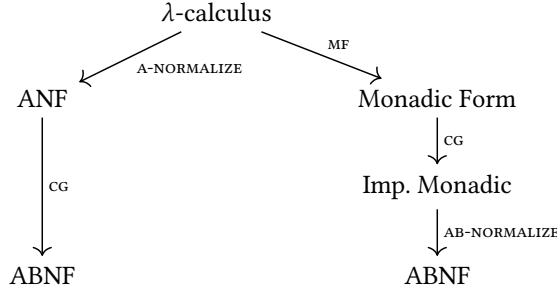


Fig. 15. Compiler Architecture

834
835
836
837
838
839
840
841
842
843
844
845
846
847 $MAX-STACK[\mathcal{D}_M] \geq MAX-STACK[\mathcal{D}_{AB}]$. Furthermore, there exists a program t with trace \mathcal{D}_{AB} such
848 that $MAX-STACK[\mathcal{D}_M] > MAX-STACK[\mathcal{D}_{AB}]$.
849

850 5 AN AB-NORMAL COMPILER

851 Formalizing these normal forms as reduction systems is useful for metatheory, allowing us to prove
852 generic properties of any compiler that targets these normal forms, but the formalism ignores
853 important details about how to implement a compiler using these normal forms. In this section,
854 we define two compilers: the ABnormal compiler, and the ANF compiler. Their pipelines are
855 summarized in Figure 15. The ABnormal compiler targets AB-normal form via monadic form and
856 AB-normalization, while the ANF compiler target AB-normal form via A-normalization. We show
857 that the compiler design is simplified using AB-normalization compared to A-normalization.
858

859 5.1 ANF Compiler

860
861 There are pragmatic problems to implementing an effective compiler into ANF. This is one of
862 the major disadvantages of target ANF directly: normalizing commuting conversions while also
863 sequencing computations introduces (unnecessary) complexity.

864 We define ANF translation in Figure 16. This implementation requires managing some complexity.

865 The main issue in the definition is that the result of ANF translation is an M , which cannot be
866 composed internally with another M . Some other external technique is needed to compose two M s.
867 This is a result of the requirement that we normalize Equation Associativity and Equation Commute.
868 We use the original technique of Flanagan et al. [8], in which the compiler reifies the evaluation
869 context of the reduction system as a meta-language continuation that builds target language terms.
870 The compiler is indexed by this continuation, and builds up the translation in the continuation,
871 rather than directly returning the translated term. The continuation has type $V \rightarrow M$. The original
872 Flanagan et al. [8] version did not actually produce ANF terms, but monadic terms, to avoid code
873 duplication. The typical solution in the literature to generating ANF without code duplication is
874 to introduce a *join point* [1, 3, 11, 14]. This implementation technique has been formalized and
875 verified to be correct with respect to typing, whole program compilation, and separate compilation
876 by Koronkevich et al. [13] in the context of an ANF compiler for dependent types.

877 The first bit of complexity is that we write the compiler in CPS. This adds some complexity to
878 the implementation, and may negatively impact compile-time performance. The compiler $ANF[e]\kappa$
879 takes a source term e and an ANF evaluation context (continuation) κ . When e is a value v , it's
880 translated by calling the continuation with the value $\kappa[v]$, forming a complete ANF program.
881 Otherwise, the subterms are translated, calling the compiler in CPS: a subterm is translated with a
882

$$\kappa ::= \cdot \mid (\mathbf{let} (x \cdot) M)$$

$$\boxed{\text{ANF}[[e]]\kappa = M}$$

$$\begin{aligned} \text{ANF}[[e]] &= \text{ANF}[[e]]\cdot \\ \text{ANF}[[t]]\kappa &= \kappa[t] \\ \text{ANF}[[x]]\kappa &= \kappa[x] \\ \text{ANF}[[\lambda (x) e]]\kappa &= \kappa[(\lambda (x) \text{ANF}[[e]])] \\ \text{ANF}[[e_1 e_2]]\kappa &= \text{ANF}[[e_1]](\mathbf{let} (x_1 \cdot) \text{ANF}[[e_2]](\mathbf{let} (x_2 \cdot) \kappa[(x_1 x_2)])) \\ \text{ANF}[[\mathbf{if0} e e_1 e_2]]\kappa &= \text{ANF}[[e]](\mathbf{let} (f(\lambda (x) \kappa[x])) \\ &\quad (\mathbf{let} (x \cdot) (\mathbf{if0} x \text{ANF}[[e_1]](\mathbf{let} (x_1 \cdot) (f x_1)) \\ &\quad \quad \text{ANF}[[e_2]](\mathbf{let} (x_2 \cdot) (f x_2)))))) \\ \text{ANF}[[op \vec{e}]]\kappa &= \text{ANF}[[e_i]](\mathbf{let} (x_i \cdot) \text{ANF}[[e_{i1}]] (\mathbf{let} (x_{i1} \cdot) \dots \kappa[(op \vec{x})])) \\ \text{ANF}[[\mathbf{let} (x e_1) e_2]]\kappa &= \text{ANF}[[e_1]](\mathbf{let} (x \cdot) \text{ANF}[[e_2]]\kappa) \end{aligned}$$

Fig. 16. λ -calculus to ANF Compiler with Join Points

new continuation, which when called with a ANF value, produces an ANF term. Compiling n-ary operators requires a fold over the list of operands, which is slightly informally specified.

It's possible to avoid this CPSed compiler, but empirically this approach appears to be easier to reason about. An alternative involves returning two values, the tail of the computation M and the list of introduced bindings, and eventually merging them⁶. Bowman [2] attempted a proof of compiler correctness for ANF using a technique that relied on reasoning about lists of introduced bindings in this way, but the formalism did not scale to join points or branching constructs in general. By contrast, Koronkevich et al. [13] provided an extension of the ANF translation, with join points and branching, proving correctness entirely by a dependent typing of the compiler's continuation. The history with compiler verification suggests that if we care about reasoning, the CPSed compiler is the right approach, as it leads to a scalable, compositional, verifiable compiler.

The second bit of complexity is the representation of the join point in the target language. If we just use λ , as we do in the above translation, we introduce a procedure call for every branch. Worse, it's a closure, since there are free variables in those branches, introducing allocation and memory indirects. Instead, we need an intermediate language with some notion of continuation, ideally one that introduces no allocation, allowing registers and frames to be shared between the caller and callee. This is possible—for example, Kennedy [11] and Tolmach and Oliva [21] gives intermediate languages with constructs for introducing a local continuation, which could be compiled separately, and contrast this with ANF. Maurer et al. [14] add a primitive join point form and equations for reasoning about them, citing problems with ANF. Cong et al. [3] add control operators, and types to distinguish different kinds of continuations, to enable optimizing using either or both ANF or CPS. All these choices increases the complexity of the IL, with new abstractions for procedures and continuations.

While managing all this complexity and using ANF successfully is possible, it is unnecessary. Our ANF compiler is more complex and unsafe for scope.

5.2 Monadic Compiler

We define the monadic form translation in Figure 17. The monadic form compiler merely sequences intermediate computation, but does not normalize commuting conversions.

⁶<http://siek.blogspot.com/2012/07/my-new-favorite-abstract-machine-ecc-on.html>

MF[[e]] = C	
MF[[ι]]	= ι
MF[[x]]	= x
MF[[λ (x) e]]	= (λ (x) MF[[e]])
MF[[e ₁ e ₂]]	= (let (x ₁ MF[[e ₁]]) (let (x ₂ MF[[e ₂]]) (x ₁ x ₂)))
MF[[if0 e e ₁ e ₂]]	= (let (x MF[[e]]) (if0 x MF[[e ₁]] MF[[e ₂]]))
MF[[op \vec{e}_i]]	= (let (x _i MF[[e _i]]) (op \vec{x}_i))
MF[[let (x e ₁) e ₂]]	= (let (x MF[[e ₁]]) MF[[e ₂]])

Fig. 17. λ-calculus to Monadic Form Compiler

The monadic form translation is straightforward; every expression can be locally transformed to sequence computation. The simplicity is because the compiler return a C , and a C can be composed with any other C using **let**. That is, the simplicity relies on *not* normalizing [Equation Associativity](#) and [Equation Commute](#).

This monadic form compiler does not produce the *same* B -normal form that the B -reductions produce. Instead, they're equivalent up to [Equation Associativity](#). However, they AB -normalize to equal terms.

5.3 Code Generation

We define two versions of code generation: one for ANF and one for monadic form. Performing code generation from monadic form relies the ability to generate (**set!** x t), while the ANF code generator should never generate this instruction. Otherwise, the output of the ANF compiler would need to be AB -normalized, defeating the purpose of A -normalizing in the first place.

The code generator for ANF is given in [Figure 18a](#). Code generation from ANF is straightforward, since all the complexity is in getting into ANF in the first place. We give the types of each translation function, although they are imprecise. The function $\text{ACG}[_]_N$ does not return an arbitrary t , but only a value, call, or primitive operation, which are also valid in non-tail position. The definitions could be collapsed, but since we must be careful to ensure well-formedness of **set!** to avoid extra stack frames, we separate the traversal of different syntactic categories.

The code generator for monadic form is defined in [Figure 18b](#). It is simpler than the ANF code generator, since it can freely compose all ts in (**set!** x t).

5.4 AB-normalizer

Finally, we define the AB -normalizer [Figure 19](#). The entire compiler is essentially searching for a (**set!** x t) form, then performing the AB -reductions, recursively. There is nothing complex in the definition; all rules but the last two are just traversals. The compiler assumes all λ s are bound to variable with **set!**, which is necessary in the region calculi; this also simplifies the compilation of operands.

The compiler $\text{AB-NORMAL-COMPILE} = \text{ABNF} \cdot \text{MCG} \cdot \text{MF}$ is simpler than the $\text{ANF-COMPILE} = \text{ACG} \cdot \text{ANF}$ compiler. AB-NORMAL-COMPILE does not require CPS to implement, does not require implementing (or managing the implementation of) join points, and it achieves the same or better performance characteristics as ANF-COMPILE , by [Theorem 4.2](#).

$$\begin{array}{ll}
981 & \text{ACG}[_] : M \rightarrow t \\
982 & \text{ACG}[\text{(let } (x N) M)] = (\text{begin } (\text{set! } x \text{ ACG}[N]_N) \\
983 & \quad \text{ACG}[M]) \\
984 & \text{ACG}[\text{(if0 } V M_1 M_2)] = (\text{if0 } \text{ACG}[V]_V \text{ ACG}[M_1] \\
985 & \quad \text{ACG}[M_2]) \\
986 & \text{ACG}[V] = \text{ACG}[V]_V \\
987 & \text{ACG}[N] = \text{ACG}[N]_N \\
988 & \text{ACG}[_]_N : N \xrightarrow{\quad} t \\
989 & \text{ACG}[(O \vec{V})]_N = (O \text{ACG}[V]_V) \\
990 & \text{ACG}[V]_N = \text{ACG}[V]_V \\
991 & \text{ACG}[_]_V : V \rightarrow t \\
992 & \text{ACG}[x]_V = x \\
993 & \text{ACG}[t]_V = t \\
994 & \text{ACG}[(\lambda (x) M)]_V = (\lambda (x) \text{ACG}[M]) \\
995 & \\
996 & \\
997 & \text{(a) From ANF} \\
998 & \\
999 & \\
1000 & \\
1001 & \\
1002 & \text{MCG}[_] : C \rightarrow t \\
1003 & \text{MCG}[\text{(let } (x C_1) C_2)] = (\text{begin } (\text{set! } x \text{ MCG}[C_1]) \\
1004 & \quad \text{MCG}[C_2]) \\
1005 & \text{MCG}[\text{(if0 } U C_1 C_2)] = (\text{if0 } \text{MCG}[U] \text{ MCG}[C_1] \\
1006 & \quad \text{MCG}[C_2]) \\
1007 & \text{MCG}[(O \vec{U})] = (O \text{MCG}[U]) \\
1008 & \text{MCG}[x] = x \\
1009 & \text{MCG}[t] = t \\
1010 & \text{MCG}[(\lambda (x) C)] = (\lambda (x) \text{MCG}[C]) \\
1011 & \\
1012 & \\
1013 & \text{(b) From Monadic Form} \\
1014 & \\
1015 & \\
1016 & \\
1017 & \\
1018 & \\
1019 & \\
1020 & \\
1021 & \\
1022 & \\
1023 & \\
1024 & \\
1025 & \\
1026 & \\
1027 & \\
1028 & \\
1029 &
\end{array}$$

Fig. 18. Code Generation

$$\begin{array}{ll}
1001 & \text{ABNF}[_]_t : t \rightarrow t \\
1002 & \text{ABNF}[\text{(begin } \vec{s} t)]_t = (\text{begin } \text{ABNF}[\vec{s}]_s \text{ABNF}[t]_t) \\
1003 & \text{ABNF}[\text{(if0 } v t_1 t_2)]_t = (\text{if0 } v \text{ABNF}[t_1]_t \text{ABNF}[t_2]_t) \\
1004 & \text{ABNF}[(op \vec{v})]_t = (op \vec{v}) \\
1005 & \text{ABNF}[(call v_1 v_2)]_t = (\text{call } v_1 v_2) \\
1006 & \text{ABNF}[_]_v : v \rightarrow v \\
1007 & \text{ABNF}[x]_v = x \\
1008 & \text{ABNF}[t]_v = t \\
1009 & \text{ABNF}[(\lambda (x) t)]_v = (\lambda (x) \text{ABNF}[t]_t) \\
1010 & \text{ABNF}[_]_s : s \rightarrow s \\
1011 & \text{ABNF}[\text{(begin } \vec{s})]_s = (\text{begin } \text{ABNF}[\vec{s}]_s) \\
1012 & \text{ABNF}[\text{(set! } x v)]_s = (\text{set! } x \text{ABNF}[v]_v) \\
1013 & \text{ABNF}[\text{(set! } x (op \vec{v}))]_s = (\text{set! } x (op \vec{v})) \\
1014 & \text{ABNF}[\text{(set! } x (\text{call } v_1 v_2))]_s = (\text{set! } x (\text{call } v_1 v_2)) \\
1015 & \text{ABNF}[\text{(set! } x (\text{begin } \vec{s} t))]_s = (\text{begin } \text{ABNF}[\vec{s}]_s \text{ABNF}[(\text{set! } x t)]_s) \\
1016 & \text{ABNF}[\text{(set! } x (\text{if0 } v t_1 t_2))]_s = (\text{if0 } v \text{ABNF}[(\text{set! } x t_1)]_s \text{ABNF}[(\text{set! } x t_2)]_s) \\
1017 & \\
1018 & \\
1019 & \\
1020 & \\
1021 & \\
1022 & \\
1023 & \\
1024 & \\
1025 & \\
1026 & \\
1027 & \\
1028 & \\
1029 &
\end{array}$$

Fig. 19. AB-normalization

6 CONCLUSIONS

6.1 What about CPS

We ignore CPS throughout this paper because, in our view, CPS solves a different problem than ANF.

A-normalization, and our AB-normalization, solve the problem of *local* control: how to explicate the data and control flow of subexpressions without control effects. But neither address the problem

of *non-local* control, such as returning from a function call (a control effect). The problems with ANF begin when conflating these two, using a solution for non-local control (continuations) to solve a local control problem (commuting conversions).

Non-local control is still important. In this paper, we never notice the non-local control problem since we never compile `CALL` and `RETURN`. To compile these, we need something like a join point or a control operator—*i.e.*, we need object-language continuations by any other name. Whatever the technique, it should reify the continuation introduced by the non-tail `CALL` instruction, and used by the `RETURN` instruction, into data, which is stored in the heap, and optimized by some set of equations. This would let us transform our CESK machine into a CES machine, which more closely models a real machine, with a code pointer, a register file, and memory. For example, if we add a **let/cc** instruction of the form (**let/cc** k t) to our ABNF, we can compile `CALL` and `RETURN` as something like the following.

$$\begin{aligned} \text{MCG}[\langle\langle \text{call } v_1 \ v_2 \rangle\rangle_C] &= (\text{let/cc } k \ (\text{call } \text{MCG}[\langle v_1 \rangle] \ \text{MCG}[\langle v_2 \rangle] \ k)) \\ \text{MCG}[\langle v \rangle_C] &= (\text{call } k \ \text{MCG}[\langle v \rangle_V]) \end{aligned}$$

That is, when we find a non-tail call (a call in C position), we capture the current continuation and generate a tail-call (since the body of **let/cc** accepts a tail) that explicitly passes the captured continuation as its return label. A value in C position is being returned, so we transform it into a tail-call to the continuation. The **let/cc** operation can be compiled to a labelled instruction that pushes and pops caller saved variables around the call, avoiding allocating continuations altogether.

This calculus starts to look, morally, like Tolmach and Oliva [21]’s SIL or Cong et al. [3]’s IL. Chez Scheme implements this transformation, suggesting it works well in practice (Subsection 6.4). We could probably perform a rational reconstruction of this pattern from the machine semantics, as we did with ANF and monadic form.

6.2 Case-of-Case

The AB-normal compiler fails to optimize case-of-case commuting conversions. Consider the example (**if0** (**if0** e 1 0) 5 6). From this, the ANF compiler generates (**if0** e (**if0** 1 5 6) (**if0** 0 5 6)). We can see the branches have been duplicated, which looks like a problem. However, now a simple partial evaluator can optimize this to (**if0** e 6 5). By contrast, the AB-normal compiler produces first (**let** (x (**if0** e 1 0)) (**if0** x 5 6)) then (**begin** (**if0** e (**set!** x 1) (**set!** x 0)) (**if0** x 5 6)). No code is duplicated, but we cannot easily perform the same optimization. The join point calculus of Maurer et al. [14] handles this example very well, although at the cost of introducing join points.

We believe this can be supported without join points by contextually separating values and boolean position. For example, Danvy [4] present a standard monadic translation (such as we define in Figure 17) extended with the following rules (more or less).

$$\begin{aligned} \text{MF}[\langle\langle \text{if0 } e \ e_1 \ e_2 \rangle\rangle] &= (\text{if } \llbracket e \rrbracket_P \ \text{MF}[\langle e_1 \rangle] \ \text{MF}[\langle e_2 \rangle]) \\ \llbracket _ \rrbracket_P &: C \rightarrow \mathbb{B} \\ \llbracket \langle\langle \text{if0 } e \ e_1 \ e_2 \rangle\rangle \rrbracket_{\mathbb{B}} &= (\text{if } \llbracket e \rrbracket_{\mathbb{B}} \ \llbracket e_1 \rrbracket_{\mathbb{B}} \ \llbracket e_2 \rrbracket_{\mathbb{B}}) \\ \llbracket \langle v \rangle \rrbracket_{\mathbb{B}} &= (\text{equal? } \text{MF}[\langle v \rangle] \ 0) \end{aligned}$$

This compiles pattern matching to a sublanguage of boolean expressions. With this translation, the example is transformed to (**if** (**if** (**equal?** e 0) (**equal?** 1 0) (**equal?** 0 0)) 5 6), which a simple boolean optimizer could simplify as (**if** (**if** (**equal?** e 0) **false true**) 5 6), then (**if** (**equal?** e 0) 6 5).

It’s not clear that this completely solves the case-of-case problem that Maurer et al. [14] investigate, but it solves their simple example, and does so without join points. Chez Scheme implements this transformation, suggesting it works well in practice (Subsection 6.4).

6.3 AB-normalizing with Monadic Effects

We work in an imperative language, but that’s a choice based on our thinking in machines. AB-normalization could be first performed in a high-level monadic language by transforming lexical binding into a state monad to normalize commuting conversions. For example, perhaps an A-normalizing monadic compiler would normalize with respect to the following rule.

$$E[(\mathbf{if} \ v \ e_1 \ e_2)] \longrightarrow_B \ (\mathbf{do} \ (\mathbf{if} \ v \ (\mathbf{do} \ (v_1 \leftarrow e_1) \ (\mathbf{set} \ x \ v_1)) \ (\mathbf{do} \ (v_2 \leftarrow e_2) \ (\mathbf{set} \ x \ v_2)))) \quad B_A$$

$$E[(\mathbf{get} \ x)]$$

We use the **do** notation to interact with the state monad in the target language. This lets us to express an **if** statement in our high-level language. If **do** is compiled to **begin**, **set** to **set!**, and **get** to variable reference, then this ends up in the same AB-normal form as our compiler.

Similarly, regions could first be elaborated into monadic regions [9, 10] rather than directly to low-level machine instructions.

Of course, if we consider **let** to be the bind of partiality monad, for explicating local control, and **do** to be the bind for the state monad, we may run into problems with composing these two monads. The problem is made worse if we throw in a region monad. Perhaps we need a single “compilation effects” monad; or perhaps not—what’s wrong with an imperative language that behaves monadically? We leave this question to future work.

6.4 AB-normalization in Practice

The AB-reductions are the essence of a transformation done in Chez Scheme. The can be seen, indirectly, in the difference between two of Chez’s language definitions: L10⁷, and L7⁸. L7 includes the expression (**set!** *lvalue* *e*), which is our (**set!** *x e*), and L6 generates this from a **let** expression. L10 rewrite this into (**set!** *lvalue rhs*), where *rhs* is a heavily restricted set of operations very similar to our final ABNF. In fact, Chez uses both of the transformations mentioned above: transforming boolean position into a predicate sublanguage, to optimize case-of-case transformations⁹, and compiling non-tail calls to tail calls with **let/cc** (called **return-point** in Chez)¹⁰.

The Chez compiler is renowned for its performance both at compile-time and in the run-time code it generates. Our paper, in part, is an attempt to understand why and how Chez avoids CPS, ANF, and monadic form—the standard choices in the published literature—in favour of an apriori strange imperative intermediate representation. The answer is that they choose ABNF before it had a name.

ACKNOWLEDGMENTS

We gratefully acknowledge the feedback of the many reviewers of this work. Also some other people.

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC), funding reference number RGPIN-9999-9999. Cette recherche a été financée par le Conseil de recherches en sciences naturelles et en génie du Canada (CRSNG), numéro de référence RGPIN-9999-9999.

⁷<https://github.com/cisco/ChezScheme/blob/9576b83dd757cf1494933c9fbc80cb6aff022295/s/np-languages.ss#L836>

⁸<https://github.com/cisco/ChezScheme/blob/9576b83dd757cf1494933c9fbc80cb6aff022295/s/np-languages.ss#L561>

⁹<https://github.com/cisco/ChezScheme/blob/9576b83dd757cf1494933c9fbc80cb6aff022295/s/np-languages.ss#L901>

¹⁰<https://github.com/cisco/ChezScheme/blob/9576b83dd757cf1494933c9fbc80cb6aff022295/s/np-languages.ss#L1055>

DATA AVAILABILITY STATEMENT

The software artifact with mechanized models, containing the full definitions and details, is available in anonymous supplementary materials and will be made publically available.

REFERENCES

- [1] Nick Benton, Andrew Kennedy, and George Russell. 1998. Compiling Standard ML to Java Bytecodes. In *International Conference on Functional Programming (ICFP)*. <https://doi.org/10.1145/289423.289435>
- [2] William J. Bowman. 2018. *Compiling with Dependent Types*. Ph.D. Dissertation. Northeastern University. <https://doi.org/10.17760/D20316239>
- [3] Youyou Cong, Leo Osvald, Grégory M. Essertel, and Tiark Rumpf. 2019. Compiling with continuations, or without? whatever. *PACMPL* 3, ICFP (2019), 79:1–79:28. <https://doi.org/10.1145/3341643>
- [4] Olivier Danvy. 2003. A New One-Pass Transformation into Monadic Normal Form. In *Compiler Construction, 12th International Conference, CC 2003, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2003, Warsaw, Poland, April 7-11, 2003, Proceedings (Lecture Notes in Computer Science, Vol. 2622)*, Görel Hedin (Ed.), Springer, 77–89. https://doi.org/10.1007/3-540-36579-6_6
- [5] Matthias Felleisen. 1987. *The Calculi of Lambda-v-cs Conversion: A Syntactic Theory of Control and State in Imperative Higher-order Programming Languages*. phdthesis. <https://www2.ccs.neu.edu/racket/pubs/dissertation-felleisen.pdf>
- [6] Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. 2009. *Semantics Engineering with PLT Redex*. MIT Press. <https://mitpress.mit.edu/books/semantics-engineering-plt-redex>
- [7] Matthias Felleisen and Daniel P. Friedman. 1987. Control operators, the SECD-machine, and the λ -calculus. In *Formal Description of Programming Concepts III*, Martin Wirsing (Ed.), North-Holland, 193–222. <https://legacy.cs.indiana.edu/ftp/techreports/TR197.pdf>
- [8] Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In *International Conference on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/155090.155113>
- [9] Matthew Fluet and Greg Morrisett. 2006. Monadic regions. *J. Funct. Program.* 16, 4-5 (2006), 485–545. <https://doi.org/10.1017/S095679680600596X>
- [10] Matthew Fluet, Greg Morrisett, and Amal J. Ahmed. 2006. Linear Regions Are All You Need. In *European Symposium on Programming (ESOP)*. Springer. https://doi.org/10.1007/11693024_2
- [11] Andrew Kennedy. 2007. Compiling with Continuations, Continued. In *International Conference on Functional Programming (ICFP)*. <https://doi.org/10.1145/1291220.1291179>
- [12] Casey Klein, John Clements, Christos Dimoulas, Carl Eastlund, Matthias Felleisen, Matthew Flatt, Jay A. McCarthy, Jon Raskind, Sam Tobin-Hochstadt, and Robert Bruce Findler. 2012. Run your research: on the effectiveness of lightweight mechanization. In *Symposium on Principles of Programming Languages (POPL)*. ACM. <https://doi.org/10.1145/2103656.2103691>
- [13] Paulette Koronkevich, Ramon Rakow, Amal Ahmed, and William J. Bowman. 2022. ANF Preserves Dependent Types up to Extensional Equality. *Journal of Functional Programming (JFP)* 32 (2022). <https://doi.org/10.1017/s0956796822000090>
- [14] Luke Maurer, Paul Downen, Zena M. Ariola, and Simon Peyton Jones. 2017. Compiling without Continuations. In *International Conference on Programming Language Design and Implementation (PLDI)*. <https://doi.org/10.1145/3062341.3062380>
- [15] Eugenio Moggi. 1991. Notions of Computation and Monads. *Inf. Comput.* 93, 1 (1991), 55–92. [https://doi.org/10.1016/0890-5401\(91\)90052-4](https://doi.org/10.1016/0890-5401(91)90052-4)
- [16] Zoe Paraskevopoulou. 2020. *Verified Optimizations for Functional Languages*. Ph.D. Dissertation. Princeton University. https://zoep.github.io/thesis_final.pdf
- [17] Zoe Paraskevopoulou and Anvay Grover. 2021. Compiling with continuations, correctly. *Proc. ACM Program. Lang.* 5, OOPSLA (2021), 1–29. <https://doi.org/10.1145/3485491>
- [18] Zoe Paraskevopoulou, John M. Li, and Andrew W. Appel. 2021. Compositional optimizations for CertiCoq. *Proc. ACM Program. Lang.* 5, ICFP (2021), 1–30. <https://doi.org/10.1145/3473591>
- [19] Zhong Shao and Andrew W. Appel. 1994. Space-Efficient Closure Representations. In *Proceedings of the 1994 ACM Conference on LISP and Functional Programming, Orlando, Florida, USA, 27-29 June 1994*, Robert R. Kessler (Ed.), ACM, 150–161. <https://doi.org/10.1145/182409.156783>
- [20] Mads Tofte and Jean-Pierre Talpin. 1994. Implementation of the typed call-by-value λ -calculus using a stack of regions. In *Symposium on Principles of Programming Languages (POPL)*. ACM Press. <https://doi.org/10.1145/174675.177855>
- [21] Andrew P. Tolmach and Dino Oliva. 1998. From ML to Ada: Strongly-Typed Language Interoperability via Source Translation. *J. Funct. Program.* 8, 4 (1998), 367–412. <https://doi.org/10.1017/S0956796898003086>