



Contextual Equivalence for a Probabilistic Language with Continuous Random Variables and Recursion

MITCHELL WAND, Northeastern University, USA

RYAN CULPEPPER, Czech Technical University in Prague, Czech Republic

THEOPHILOS GIANNAKOPOULOS, BAE Systems, USA

ANDREW COBB, Northeastern University, USA

We present a complete reasoning principle for contextual equivalence in an untyped probabilistic language. The language includes continuous (real-valued) random variables, conditionals, and scoring. It also includes recursion, since the standard call-by-value fixpoint combinator is expressible.

We demonstrate the usability of our characterization by proving several equivalence schemas, including familiar facts from lambda calculus as well as results specific to probabilistic programming. In particular, we use it to prove that reordering the random draws in a probabilistic program preserves contextual equivalence. This allows us to show, for example, that

$$(\text{let } x = e_1 \text{ in let } y = e_2 \text{ in } e_0) =_{\text{ctx}} (\text{let } y = e_2 \text{ in let } x = e_1 \text{ in } e_0)$$

(provided x does not occur free in e_2 and y does not occur free in e_1) despite the fact that e_1 and e_2 may have sampling and scoring effects.

CCS Concepts: • **Mathematics of computing** → **Probability and statistics**; • **Theory of computation** → *Operational semantics*;

Additional Key Words and Phrases: probabilistic programming, logical relations, contextual equivalence

ACM Reference Format:

Mitchell Wand, Ryan Culpepper, Theophilos Giannakopoulos, and Andrew Cobb. 2018. Contextual Equivalence for a Probabilistic Language with Continuous Random Variables and Recursion. *Proc. ACM Program. Lang.* 2, ICFP, Article 87 (September 2018), 30 pages. <https://doi.org/10.1145/3236782>

1 INTRODUCTION

A *probabilistic programming language* is a programming language enriched with two features—*sampling* and *scoring*—that enable it to represent probabilistic models. We introduce these two features with an example program that models linear regression.

The first feature, *sampling*, introduces probabilistic nondeterminism. It is used to represent random variables. For example, let `normal`(m, s) be defined to nondeterministically produce a real number distributed according to a normal (Gaussian) distribution with mean m and scale s .

Here is a little model of linear regression that uses `normal` to randomly pick a slope and intercept for a line and then defines `f` as the resulting linear function:

Authors' addresses: Mitchell Wand, College of Computer and Information Science, Northeastern University, 360 Huntington Ave, Room 202WVH, Boston, MA, USA, 02115, wand@ccs.neu.edu; Ryan Culpepper, Faculty of Information Technology, Czech Technical University in Prague, Thákurova 9, Prague, 16000, Czech Republic, ryanc@ccs.neu.edu; Theophilos Giannakopoulos, FAST Labs, BAE Systems, 600 District Ave, Burlington, MA, USA, 01803, tgiannak@alum.wpi.edu; Andrew Cobb, College of Computer and Information Science, Northeastern University, 360 Huntington Ave, Room 202WVH, Boston, MA, USA, 02115, acobb@ccs.neu.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2018 Copyright held by the owner/author(s).

2475-1421/2018/9-ART87

<https://doi.org/10.1145/3236782>

```
A = normal(0, 10)
B = normal(0, 10)
f(x) = A*x + B
```

This program defines a distribution on lines, centered on $y = 0x + 0$, with high variance. This distribution is called the *prior*, since it is specified prior to considering any evidence.

The second feature, *scoring*, adjusts the likelihood of the current execution's random choices. It is used to represent conditioning on observed data.

Suppose we have the following data points: $\{(2.0, 2.4), (3.0, 2.7), (4.0, 3.0)\}$. The smaller the error between the result of f and the observed data, the better the choice of A and B . We express these observations with the following addition to our program:

```
factor normalpdf(f(2.0)-2.4; 0, 1)
factor normalpdf(f(3.0)-2.7; 0, 1)
factor normalpdf(f(4.0)-3.0; 0, 1)
```

Here scoring is performed by the `factor` form, which takes a positive real number to multiply into the current execution's likelihood. We use `normalpdf(_, 0, 1)`—the density function of the standard normal distribution—to convert the difference between predicted and observed values into the score. This scoring function assigns high likelihood when the error is near 0, dropping off smoothly to low likelihood for larger errors.

After incorporating the observations, the program defines a distribution centered near $y = 0.3x + 1.8$, with low variance. This distribution is often called the *posterior distribution*, since it represents the distribution after the incorporation of evidence.

Computing the posterior distribution—or a workable approximation thereof—is the task of *probabilistic inference*. We say that this program has *inferred* (or sometimes *learned*) the parameters A and B from the data. Probabilistic inference encompasses an arsenal of techniques of varied applicability and efficiency. Some inference techniques may benefit if the program above is transformed to the following shape:

```
A = normal(0, 10)
factor Z(A)
B = normal(M(A), S(A))
```

The transformation relies on the conjugacy relationship between the normal prior for B and the normal scoring function of the observations. A useful equational theory for probabilistic programming must incorporate facts from mathematics in addition to standard concerns such as function inlining.

In this paper we build a foundation for such an equational theory for a probabilistic programming language. In particular, our language supports

- sampling continuous random variables,
- scoring (soft constraints), and
- conditionals, higher-order functions, and recursion.

Other such languages include Church [Goodman et al. 2008], its descendants such as Venture [Mansinghka et al. 2014] and Anglican [Wood et al. 2014], and other languages [Kiselyov and Shan 2009; Narayanan et al. 2016; Paige and Wood 2014] and language models [Borgström et al. 2016; Huang and Morrisett 2016; Park et al. 2008]. Our framework is able to justify the transformation above.

In Section 2 we present our model of a probabilistic language, including its syntax and semantics. We then define our logical relation (Section 3), our CIU relation (Section 4), and contextual ordering (Section 5); and we prove that all three relations coincide. As usual, contextual ordering is powerful but difficult to prove directly. The virtue of the logical relation is that it eliminates the need to

reason about arbitrary syntactic contexts; they are boiled down to their essential components: substitutions and continuations (evaluation contexts). The CIU relation [Mason and Talcott 1991] is a further simplification of the logical relation; it offers the easiest way to prove relationships between specific terms. In Section 6 we define contextual equivalence and use the CIU relation to demonstrate a catalog of useful equivalence schemas, including β_v and let-associativity, as well as a method for importing first-order equivalences from mathematics. One unusual equivalence is let-commutativity:

$$(\text{let } x = e_1 \text{ in let } y = e_2 \text{ in } e_0) =_{\text{ctx}} (\text{let } y = e_2 \text{ in let } x = e_1 \text{ in } e_0)$$

(provided x does not occur free in e_2 and y does not occur free in e_1). This equivalence, while valid for a pure language, is certainly not valid for all effects (consider, for example, if there were an assignment statement in e_1 or e_2). In other words, sampling and scoring are commutative effects. We conclude with two related work sections: Section 7 demonstrates the correspondence between our language model and others, notably that of Borgström et al. [2016], and Section 8 informally discusses other related work.

Throughout the paper, we limit proofs mostly to high-level sketches and representative cases. Additional details and cases for some proofs can be found in the appendices of the long version of this paper [Wand et al. 2018].

2 PROBABILISTIC LANGUAGE MODEL

In this section we define our probabilistic language and its semantics. The semantics consists of three parts:

- A notion of *entropy* for modeling random behavior.
- An *evaluation* function that maps a program and entropy to a real-valued result and an importance weight. We define the evaluation function via an abstract machine. We then define a big-step semantics and prove it equivalent; the big-step formulation simplifies some proofs in Section 6.3 by making the structure of evaluation explicit.
- A mapping to *measures* over the real numbers, calculated by integrating the evaluation function with respect to the entropy space. A program with a finite, non-zero measure can be interpreted as an unnormalized probability distribution.

The structure of the semantics loosely corresponds to one inference technique for probabilistic programming languages: importance sampling. In an importance sampler, the entropy is approximated by a pseudo-random number generator (PRNG); the evaluation function is run many times with different initial PRNG states to produce a collection of weighted samples; and the weighted samples approximate the program's measure—either directly by conversion to a discrete distribution of results, or indirectly via computed statistical properties such as sample mean, variance, etc.

Our language is similar to that of Borgström et al. [2016], but with the following differences:

- Our language requires let-binding of nontrivial intermediate expressions; this simplifies the semantics. This restriction is similar to but looser than A-normal form [Sabry and Felleisen 1993].
- Our model of entropy is a finite measure space made of *splittable* entropy points, rather than an infinite measure space containing *sequences* of real numbers.
- Our sample operation models a standard uniform random variable, rather than being parameterized over a distribution.

We revisit these differences in Section 7.

$v ::= x \mid \lambda x.e \mid c_r$	Syntactic Values
$e ::= v \mid (v v) \mid \text{let } x = e \text{ in } e$	Expressions
$\mid op^n(v_1, \dots, v_n) \mid \text{if } v \text{ then } e \text{ else } e$	
$\mid \text{sample} \mid \text{factor } v$	
$op^1 ::= \log \mid \exp \mid \text{real?} \mid \dots$	Unary operations
$op^2 ::= + \mid - \mid \times \mid \div \mid < \mid \leq \mid \dots$	Binary operations
$op^3 ::= \text{normalinvcdf} \mid \text{normalpdf} \mid \dots$	Ternary operations
$K ::= \text{halt} \mid (x \rightarrow e)K$	Continuations

Fig. 1. Syntax of values, expressions, and continuations

2.1 Syntax

The syntax of our language is given in Figure 1. For simplicity, we require sequencing to be made explicit using `let`. There is a constant c_r for each real number r , and there are various useful primitive operations.

The `sample` form draws from a uniform distribution on $[0, 1]$. Any other standard real-valued distribution can be obtained by applying the appropriate *inverse cumulative distribution function*. For example, sampling from a normal distribution can be expressed as follows:

$$\text{normal}(m, s) \triangleq (\text{let } u = \text{sample} \text{ in normalinvcdf}(u; m, s))$$

where $\text{normalinvcdf}(u; m, s)$ is the least x such that if $X \sim \mathcal{N}(m, s^2)$ then $\Pr[X \leq x] = u$.

Finally, `factor` v weights (or “scores”) the current execution by the value v .

The language is untyped, but we express the scoping relations by rules like typing rules. We write $\Gamma \vdash e$ `exp` for the assertion that e is a well-formed expression whose free variables are contained in the set Γ , and similarly for values and continuations. The scoping rules are given in Figure 2.

2.2 Modeling Entropy

The semantics uses an *entropy* component as the source of randomness. We assume an entropy space \mathbb{S} along with its stock measure $\mu_{\mathbb{S}}$. We use σ and τ to range over values in \mathbb{S} . When we integrate over σ or τ , we implicitly use the stock measure; that is, we write $\int f(\sigma) d\sigma$ to mean $\int f(\sigma) \mu_{\mathbb{S}}(d\sigma)$. Following Culpepper and Cobb [2017], we assume that \mathbb{S} has the following properties:

PROPERTY 2.1 (PROPERTIES OF ENTROPY).

- (1) $\mu_{\mathbb{S}}(\mathbb{S}) = 1$
- (2) There is a function $\pi_U : \mathbb{S} \rightarrow [0, 1]$ such that for all measurable $f : [0, 1] \rightarrow \mathbb{R}^+$,

$$\int f(\pi_U(\sigma)) d\sigma = \int_0^1 f(x) \lambda(dx)$$

where λ is the Lebesgue measure. That is, π_U represents a standard uniform random variable.

- (3) There is a surjective pairing function $\text{‘::’} : \mathbb{S} \times \mathbb{S} \rightarrow \mathbb{S}$, with projections π_L and π_R , all measurable.
- (4) The projections are measure-preserving: for all measurable $g : \mathbb{S} \times \mathbb{S} \rightarrow \mathbb{R}^+$,

$$\int g(\pi_L(\sigma), \pi_R(\sigma)) d\sigma = \iint g(\sigma_1, \sigma_2) d\sigma_1 d\sigma_2$$

Since $\mathbb{S} \cong \mathbb{S} \times \mathbb{S}$ and thus $\mathbb{S} \cong \mathbb{S}^n$ ($n \geq 1$), we can also use entropy to encode non-empty sequences of entropy values.

One model that satisfies these properties is the space of infinite sequences of real numbers in $[0, 1]$; π_L and π_R take the odd- and even-indexed subsequences, respectively, and π_U takes the first element in the sequence. Another model is the space of infinite sequences of bits, where π_L and π_R

$$\begin{array}{c}
\frac{x \in \Gamma}{\Gamma \vdash x \text{ val}} \qquad \frac{\Gamma, x \vdash e \text{ exp}}{\Gamma \vdash \lambda x. e \text{ val}} \qquad \Gamma \vdash c_r \text{ val} \\
\\
\frac{\Gamma \vdash v \text{ val}}{\Gamma \vdash v \text{ exp}} \qquad \frac{\Gamma \vdash v_1 \text{ val} \quad \Gamma \vdash v_2 \text{ val}}{\Gamma \vdash (v_1 \ v_2) \text{ exp}} \qquad \frac{\Gamma \vdash e_1 \text{ exp} \quad \Gamma, x \vdash e_2 \text{ exp}}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 \text{ exp}} \\
\\
\frac{\Gamma \vdash v_i \text{ val} \quad (\forall i \in \{1, \dots, n\})}{\Gamma \vdash \text{op}^n (v_1, \dots, v_n) \text{ exp}} \qquad \frac{\Gamma \vdash v \text{ val} \quad \Gamma \vdash e_1 \text{ exp} \quad \Gamma \vdash e_2 \text{ exp}}{\Gamma \vdash \text{if } v \text{ then } e_1 \text{ else } e_2 \text{ exp}} \\
\\
\Gamma \vdash \text{sample exp} \qquad \frac{\Gamma \vdash v \text{ val}}{\Gamma \vdash \text{factor } v \text{ exp}} \\
\\
\vdash \text{halt cont} \qquad \frac{\{x\} \vdash e \text{ exp} \quad \vdash K \text{ cont}}{\vdash (x \rightarrow e)K \text{ cont}}
\end{array}$$

Fig. 2. Scoping rules for values, expressions, and continuations

take odd and even subsequences and π_U interprets the entire sequence as the binary expansion of a number in $[0, 1]$. It is tempting to envision entropy as infinite binary trees labeled with real numbers in $[0, 1]$, but the pairing function is not surjective.

We also use Tonelli's Theorem:

LEMMA 2.2 (TONELLI). *Let $f : \mathbb{S} \times \mathbb{S} \rightarrow \mathbb{R}^+$ be measurable. Then*

$$\int \left(\int f(\sigma_1, \sigma_2) d\sigma_1 \right) d\sigma_2 = \int \left(\int f(\sigma_1, \sigma_2) d\sigma_2 \right) d\sigma_1$$

2.3 Operational Semantics

2.3.1 Small-Step Semantics. We define evaluation via an abstract machine with a small-step operational semantics. The semantics rewrites configurations $\langle \sigma \mid e \mid K \mid \tau \mid w \rangle$ consisting of:

- an entropy σ (representing the “current” value of the entropy),
- a closed expression e ,
- a closed continuation K ,
- an entropy τ (encoding a stack of entropies, one for each frame of K), and
- a positive real number w (representing the weight of the current run)

The rules for the semantics are given in Figure 3.

The semantics uses continuations for sequencing and substitutions for procedure calls. Since $\text{let } x = e_1 \text{ in } e_2$ is the only sequencing construct, there is only one continuation-builder. The first rule recurs into the right-hand side of a let , using the left half of the entropy as its entropy, and saving the right half for use with e_2 . The second rule (“return”) substitutes the value of the expression into the body of the let and restores the top saved entropy value for use in the body. More precisely, we view the third component as an encoded pair of an entropy value and an encoded

$\langle \sigma \mid \text{let } x = e_1 \text{ in } e_2 \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \pi_L(\sigma) \mid e_1 \mid (x \rightarrow e_2)K \mid \pi_R(\sigma)::\tau \mid w \rangle$
$\langle \sigma \mid v \mid (x \rightarrow e_2)K \mid \sigma'::\tau \mid w \rangle$	$\rightarrow \langle \sigma' \mid e_2[v/x] \mid K \mid \tau \mid w \rangle$
$\langle \sigma \mid ((\lambda x.e) v) \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \sigma \mid e[v/x] \mid K \mid \tau \mid w \rangle$
$\langle \sigma \mid \text{sample} \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \pi_R(\sigma) \mid c_{\pi_U(\pi_L(\sigma))} \mid K \mid \tau \mid w \rangle$
$\langle \sigma \mid \text{op}^n(v_1, \dots, v_n) \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \sigma \mid \delta(\text{op}^n, v_1, \dots, v_n) \mid K \mid \tau \mid w \rangle$ (if defined)
$\langle \sigma \mid \text{if } c_r \text{ then } e_1 \text{ else } e_2 \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \sigma \mid e_1 \mid K \mid \tau \mid w \rangle$ (if $r > 0$)
$\langle \sigma \mid \text{if } c_r \text{ then } e_1 \text{ else } e_2 \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \sigma \mid e_2 \mid K \mid \tau \mid w \rangle$ (if $r \leq 0$)
$\langle \sigma \mid \text{factor } c_r \mid K \mid \tau \mid w \rangle$	$\rightarrow \langle \sigma \mid c_r \mid K \mid \tau \mid r \times w \rangle$ (provided $r > 0$)

Fig. 3. Small-step operational semantics

entropy stack, as mentioned in Section 2.2.¹ The entropy stack τ and continuation K are always updated simultaneously. The return rule can be written using explicit projections as follows:

$$\langle \sigma \mid v \mid (x \rightarrow e_2)K \mid \tau \mid w \rangle \rightarrow \langle \pi_L(\tau) \mid e_2[v/x] \mid K \mid \pi_R(\tau) \mid w \rangle$$

Note that in the return rule the current entropy σ is dead. Except for the entropy and weight, these rules are standard for a continuation-passing interpreter for the λ -calculus with `let`.

The δ partial function interprets primitive operations. We assume that all the primitive operations are measurable partial functions returning real values, and with the exception of `real?`, they are undefined if any of their arguments is a closure. A conditional expression evaluates to its first branch if the condition is a positive real constant, its second branch if nonpositive; if the condition is a closure, evaluation is stuck. Comparison operations and the `real?` predicate return 1 for truth and 0 for falsity.

The rule for `sample` uses π_U to extract from the entropy a real value in the interval $[0, 1]$. The entropy is split first, to make it clear that entropy is never reused, but the leftover entropy is dead per the return rule. The rule for `factor` v weights the current execution by v , provided v is a positive number; otherwise, evaluation is stuck.

When reduction of an initial configuration halts properly, there are two relevant pieces of information in the final configuration: the result value and the weight. Furthermore, we are only interested in real-valued final results. We define *evaluation* as taking an extra parameter A , a measurable set of reals. Evaluation produces a positive weight only if the result value is in the expected set.

$$\text{eval}(\sigma, e, K, \tau, w, A) = \begin{cases} w' & \text{if } \langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma' \mid r \mid \text{halt} \mid \tau' \mid w' \rangle, \\ & \text{where } r \in A \\ 0 & \text{otherwise} \end{cases}$$

We will also need approximants to `eval`:

$$\text{eval}^{(n)}(\sigma, e, K, \tau, w, A) = \begin{cases} w' & \text{if } \langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma' \mid r \mid \text{halt} \mid \tau' \mid w' \rangle \\ & \text{in } n \text{ or fewer steps, where } r \in A \\ 0 & \text{otherwise} \end{cases}$$

The following lemmas are clear from inspection of the small-step semantics.

LEMMA 2.3. *If $\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow \langle \sigma' \mid e' \mid K' \mid \tau' \mid w' \rangle$ then*

¹We defer the explanation of the initial entropy stack to Section 2.4.

$$\begin{array}{c}
\sigma \vdash \lambda x.e \Downarrow \lambda x.e, 1 \qquad \qquad \qquad \sigma \vdash c_r \Downarrow c_r, 1 \\
\\
\frac{\sigma \vdash e[v/x] \Downarrow v', w}{\sigma \vdash ((\lambda x.e) v) \Downarrow v', w} \qquad \frac{\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1 \quad \pi_R(\sigma) \vdash e_2[v_1/x] \Downarrow v_2, w_2}{\sigma \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow v_2, w_2 \times w_1} \\
\\
\frac{\delta(\text{op}^n, v_1, \dots, v_n) = v}{\sigma \vdash \text{op}^n(v_1, \dots, v_n) \Downarrow v, 1} \\
\\
\frac{\sigma \vdash e_1 \Downarrow v, w \quad r > 0}{\sigma \vdash \text{if } c_r \text{ then } e_1 \text{ else } e_2 \Downarrow v, w} \qquad \frac{\sigma \vdash e_2 \Downarrow v, w \quad r \leq 0}{\sigma \vdash \text{if } c_r \text{ then } e_1 \text{ else } e_2 \Downarrow v, w} \\
\\
\sigma \vdash \text{sample} \Downarrow c_{\pi_U(\pi_L(\sigma))}, 1 \qquad \frac{r > 0}{\sigma \vdash \text{factor } c_r \Downarrow c_r, r}
\end{array}$$

Fig. 4. Big-step operational semantics

- (1) $\text{eval}^{(p+1)}(\sigma, e, K, \tau, w, A) = \text{eval}^{(p)}(\sigma', e', K', \tau', w', A)$
- (2) $\text{eval}(\sigma, e, K, \tau, w, A) = \text{eval}(\sigma', e', K', \tau', w', A)$

LEMMA 2.4 (WEIGHTS ARE LINEAR).

- (1) *Weights can be factored out of reduction sequences. That is,*

$$\langle \sigma \mid e \mid K \mid \tau \mid 1 \rangle \rightarrow^* \langle \sigma' \mid e' \mid K' \mid \tau' \mid w' \rangle,$$

if and only if for any $w > 0$

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma' \mid e' \mid K' \mid \tau' \mid w' \times w \rangle$$

- (2) *Weights can be factored out of evaluation. That is, for all $w > 0$,*

$$\text{eval}(\sigma, e, K, \tau, w, A) = w \times \text{eval}(\sigma, e, K, \tau, 1, A),$$

and similarly for $\text{eval}^{(n)}$.

2.3.2 Big-Step Semantics. We regard the small-step semantics as normative, and we use it for our primary soundness and completeness results. However, for program transformations it is useful to have a big-step semantics as well. In this section, we define a big-step semantics and characterize its relation to the small-step semantics.

The big-step semantics is given in Figure 4. It has judgments of the form $\sigma \vdash e \Downarrow v, w$, where σ is a value of the entropy, e is a closed expression, v is a closed value, and w is a weight (a positive real number). Its intention is that when e is supplied with entropy σ , it returns v with weight w , consuming some portion (possibly empty) of the given entropy σ . The rules are those of a straightforward call-by-value λ -calculus, modified to keep track of the entropy and weight.

The translation from big-step to small-step semantics is straightforward:

THEOREM 2.5 (BIG-STEP TO SMALL-STEP). *If $\sigma \vdash e \Downarrow v, w$, then for any K and τ , there exists a σ' such that*

$$\langle \sigma \mid e \mid K \mid \tau \mid 1 \rangle \rightarrow^* \langle \sigma' \mid v \mid K \mid \tau \mid w \rangle$$

PROOF. By induction on the definition of \Downarrow . We will show selected cases.

Case $\sigma \vdash \lambda x.e \Downarrow \lambda x.e, 1$: The required small-step reduction is empty. Similarly for c_r .

Case $\sigma \vdash \text{sample} \Downarrow c_{\pi_U(\pi_L(\sigma))}, 1$: The required reduction is the single step reduction

$$\langle \sigma \mid \text{sample} \mid K \mid \tau \mid 1 \rangle \rightarrow \langle \pi_R(\sigma) \mid c_{\pi_U(\pi_L(\sigma))} \mid K \mid \tau \mid 1 \rangle$$

Similarly for factor c_r and the op^n rules.

Case $((\lambda x.e) v)$: The rule is

$$\frac{\sigma \vdash e[v/x] \Downarrow v', w}{\sigma \vdash ((\lambda x.e) v) \Downarrow v', w}$$

By inversion, we have $\sigma \vdash e[v/x] \Downarrow v', w$. So the reduction sequence is:

$$\begin{aligned} & \langle \sigma \mid ((\lambda x.e) v) \mid K \mid \tau \mid 1 \rangle \\ & \rightarrow \langle \sigma \mid e[v/x] \mid K \mid \tau \mid 1 \rangle \\ & \rightarrow^* \langle \sigma' \mid v' \mid K \mid \tau \mid w \rangle \quad \text{by the induction hypothesis} \end{aligned}$$

Similarly for the if rules.

Case $\text{let } x = e_1 \text{ in } e_2$: The rule is

$$\frac{\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1 \quad \pi_R(\sigma) \vdash e_2[v_1/x] \Downarrow v_2, w_2}{\sigma \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow v_2, w_2 \times w_1}$$

By inversion, we have $\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1$ and $\pi_R(\sigma) \vdash e_2[v_1/x] \Downarrow v_2, w_2$. So the required reduction sequence is:

$$\begin{aligned} & \langle \sigma \mid \text{let } x = e_1 \text{ in } e_2 \mid K \mid \tau \mid 1 \rangle \\ & \rightarrow \langle \pi_L(\sigma) \mid e_1 \mid (x \rightarrow e_2)K \mid \pi_R(\sigma)::\tau \mid w \rangle \\ & \rightarrow^* \langle \sigma' \mid v_1 \mid (x \rightarrow e_2)K \mid \pi_R(\sigma)::\tau \mid w_1 \rangle \\ & \rightarrow \langle \pi_R(\sigma) \mid e_2[v_1/x] \mid K \mid \tau \mid w_1 \rangle \\ & \rightarrow^* \langle \sigma'' \mid v_2 \mid K \mid \tau \mid w_2 \times w_1 \rangle \end{aligned}$$

where the third line follows from the induction hypothesis, and the last line follows from the other induction hypothesis and the linearity of weights (Lemma 2.4). \square

Note that the weak quantifier (“there exists a σ'' ”) corresponds to the fact that the entropy is dead in the return rule.

In order to prove a converse, we need some additional results about the small-step semantics.

Definition 2.6. Define \geq to be the smallest relation defined by the following rules:

$$\begin{array}{l} \text{RULE 1:} \\ (K, \tau) \geq (K, \tau) \end{array} \quad \begin{array}{l} \text{RULE 2:} \\ \frac{(K', \tau') \geq (K, \tau)}{((x \rightarrow e)K', \sigma::\tau') \geq (K, \tau)} \end{array}$$

LEMMA 2.7. *Let*

$$\langle \sigma_1 \mid e_1 \mid K_1 \mid \tau_1 \mid w_1 \rangle \rightarrow \langle \sigma_2 \mid e_2 \mid K_2 \mid \tau_2 \mid w_2 \rangle \rightarrow \dots$$

be a reduction sequence in the operational semantics. Then for each i in the sequence either

- a. *there exists a smallest $j \leq i$ such that e_j is a value and $K_j = K_1$ and $\tau_j = \tau_1$, or*
- b. $(K_i, \tau_i) \geq (K_1, \tau_1)$

PROOF. See the long version [Wand et al. 2018, Appendix A]. \square

The next result is an interpolation theorem, which imposes structure on reduction sequences: any terminating computation starting with an expression e begins by evaluating e to a value v and then sending that value to the continuation K .

THEOREM 2.8 (INTERPOLATION THEOREM). *If*

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma'' \mid v'' \mid \text{halt} \mid \tau'' \mid w'' \rangle$$

then there exists a smallest n such that for some quantities σ' , v , and w' ,

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^n \langle \sigma' \mid v \mid K \mid \tau \mid w' \times w \rangle \rightarrow^* \langle \sigma'' \mid v'' \mid \text{halt} \mid \tau'' \mid w'' \rangle$$

PROOF. If $K = \text{halt}$, then the result is trivial. Otherwise, apply the invariant of the preceding lemma, observing that $(\text{halt}, \tau') \not\prec (K, \tau)$ and that weights are multiplicative. \square

Note that both Lemma 2.7 and Theorem 2.8 would be false if our language contained jumping control structures like `call/cc`.

Finally, we show that in the interpolation theorem, σ' , v , and w' are independent of K .

THEOREM 2.9 (GENERICITY THEOREM). *Let $w_1 > 0$ and let n be the smallest integer such that for some quantities σ' , v , and w' ,*

$$\langle \sigma \mid e \mid K_1 \mid \tau_1 \mid w_1 \rangle \rightarrow^n \langle \sigma' \mid v \mid K_1 \mid \tau_1 \mid w' \times w_1 \rangle$$

then for any K_2 , τ_2 , and w_2 ,

$$\langle \sigma \mid e \mid K_2 \mid \tau_2 \mid w_2 \rangle \rightarrow^n \langle \sigma' \mid v \mid K_2 \mid \tau_2 \mid w' \times w_2 \rangle$$

PROOF. Let R be the smallest relation defined by the rules

$$((K_1, \tau_1), (K_2, \tau_2)) \in R \quad \frac{((K, \tau), (K', \tau')) \in R}{(((x \rightarrow e)K, \sigma::\tau), ((x \rightarrow e)K', \sigma::\tau')) \in R}$$

Extend R to be a relation on configurations by requiring the weights to be related by a factor of w_2/w_1 and the remaining components of the configurations to be equal. It is easy to see, by inspection of the small-step rules, that R is a bisimulation over the first n steps of the given reduction sequence. \square

We are now ready to state the converse of Theorem 2.5.

Definition 2.10. We say that a configuration $\langle \sigma \mid e \mid K \mid \tau \mid w \rangle$ *halts iff*

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma' \mid v \mid \text{halt} \mid \tau' \mid w' \rangle$$

for some σ' , v , τ' and w' .

THEOREM 2.11 (SMALL-STEP TO BIG-STEP). *If*

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma'' \mid v'' \mid \text{halt} \mid \tau'' \mid w'' \rangle,$$

then there exist σ' , v' and w' such that

$$\sigma \vdash e \Downarrow v', w'$$

and

$$\langle \sigma'' \mid v' \mid K \mid \tau \mid w' \times w \rangle \rightarrow^* \langle \sigma' \mid v' \mid \text{halt} \mid \tau'' \mid w'' \rangle$$

PROOF. Given

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma'' \mid v'' \mid \text{halt} \mid \tau'' \mid w'' \rangle$$

apply the Interpolation Theorem (Theorem 2.8) to get n , σ' , v , and w' such that

$$\langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^n \langle \sigma' \mid v \mid K \mid \tau \mid w' \times w \rangle \rightarrow^* \langle \sigma'' \mid v'' \mid \text{halt} \mid \tau'' \mid w'' \rangle$$

This gives us the second part of the conclusion. To get the first part, we proceed by (course-of-values) induction on n , and then by cases on e .

Case $\lambda x.e$: For configurations of the form $\langle \sigma \mid \lambda x.e \mid K \mid \tau \mid w \rangle$, the expression is already a value, so n is 0. So set $v = \lambda x.e$ and $w' = 1$, and observe that $\sigma \vdash \lambda x.e \Downarrow \lambda x.e, 1$, as desired. The case of constants c_r is similar.

Case `sample`: We know

$$\langle \sigma \mid \text{sample} \mid K \mid \tau \mid w \rangle \rightarrow \langle \pi_R(\sigma) \mid c_{\pi_U(\pi_L(\sigma))} \mid K \mid \tau \mid w \rangle$$

so the value length is 1, and we also have $\sigma \vdash \text{sample} \Downarrow c_{\pi_U(\pi_L(\sigma))}, 1$, as desired. The cases of `factor` and of op^n are similar.

Case $((\lambda x.e) v)$: Assume that the value length of $\langle \sigma \mid ((\lambda x.e) v) \mid K \mid \tau \mid w \rangle$ is $n + 1$. So we have

$$\langle \sigma \mid ((\lambda x.e) v) \mid K \mid \tau \mid w \rangle \rightarrow \langle \sigma \mid e[v/x] \mid K \mid \tau \mid w \rangle \rightarrow^n \langle \sigma' \mid v' \mid K \mid \tau \mid w' \times w \rangle$$

By induction, we have $\sigma \vdash e[v/x] \Downarrow v', w'$. Hence, by the big-step rule for λ -expressions, we have $\sigma \vdash ((\lambda x.e) v) \Downarrow v', w'$, as desired. The cases for conditionals are similar.

Case `let x = e1 in e2`: Assume the value length of $\langle \sigma \mid \text{let } x = e_1 \text{ in } e_2 \mid K \mid \tau \mid w \rangle$ is n . Then the first n steps of its reduction sequence must be

$$\begin{aligned} & \langle \sigma \mid \text{let } x = e_1 \text{ in } e_2 \mid K \mid \tau \mid w \rangle \\ & \rightarrow \langle \pi_L(\sigma) \mid e_1 \mid (x \rightarrow e_2)K \mid \pi_R(\sigma)::\tau \mid w \rangle \\ & \rightarrow^m \langle \sigma' \mid v_1 \mid (x \rightarrow e_2)K \mid \pi_R(\sigma)::\tau \mid w_1 \times w \rangle \\ & \rightarrow \langle \pi_R(\sigma)::\tau \mid e_2[v_1/x] \mid K \mid \tau \mid w_1 \times w \rangle \\ & \rightarrow^p \langle \sigma'' \mid v \mid K \mid \tau \mid w_2 \times w_1 \times w \rangle \end{aligned}$$

where m and p are the value lengths of the configurations on the second and fourth lines, respectively. So $n = m + p + 2$, and we can apply the induction hypothesis to the two relevant configurations. Applying the induction hypothesis twice, we get

$$\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1 \quad \text{and} \quad \pi_R(\sigma) \vdash e_2[v_1/x] \Downarrow v_2, w_2 .$$

Hence, by the big-step rule for `let`, we conclude that

$$\sigma \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow v, w_2 \times w_1$$

as desired. □

2.4 From Evaluations to Measures

Up to now, we have considered only single runs of the machine, using particular entropy values. To obtain the overall meaning of the program we need to integrate over all possible values of the entropies σ and τ :

Definition 2.12. The *measure* of e and K is the measure on the reals defined by

$$\mu(e, K, A) = \iint \text{eval}(\sigma, e, K, \tau, 1, A) \, d\sigma \, d\tau$$

for each measurable set A of the reals.

This measure is similar to both Culpepper and Cobb's $\mu_e(A)$ and Borgström et al.'s $\llbracket e \rrbracket_S(A)$, but whereas they define measures on arbitrary syntactic values, our $\mu(e, K, -)$ is a measure on the reals. Furthermore, whereas their measures represent the meanings of intermediate expressions, our measure—due to the inclusion of the continuation argument K —represents the meanings of whole programs.

The simplicity of the definition above relies on the mathematical trick of encoding entropy stacks as entropy values; if we represented stacks directly the number of integrals would depend on the stack depth. Note that even for the base continuation ($K = \text{halt}$) we still integrate with respect

to both σ and τ . Since $\mathbb{S} \not\cong \mathbb{S}^0$, there is no encoding for an empty stack as an entropy value; we cannot just choose a single arbitrary τ_{init} because $\mu_{\mathbb{S}}(\{\tau_{\text{init}}\}) = 0$. But since evaluation respects the stack discipline, it produces the correct result for any initial τ_{init} . So we integrate over *all* choices of τ_{init} , and since $\mu_{\mathbb{S}}(\mathbb{S}) = 1$ the empty stack “drops out” of the integral.

As before, we will also need the approximants:

$$\mu^{(n)}(e, K, A) = \iint \text{eval}^{(n)}(\sigma, e, K, \tau, 1, A) d\sigma d\tau$$

For these integrals to be well-defined, of course, we need to know that eval and its approximants are measurable.

LEMMA 2.13 (EVAL IS MEASURABLE). *For any $e, K, w \geq 0, A \in \Sigma_{\mathbb{R}}$, and n , $\text{eval}(\sigma, e, K, \tau, w, A)$ and $\text{eval}^{(n)}(\sigma, e, K, \tau, w, A)$ are measurable in σ and τ .*

PROOF. The proof is based on the proof from [Borgström et al. \[2017\]](#). See the long version [[Wand et al. 2018, Appendix A](#)] for more details. \square

The next lemma establishes some properties of μ and the approximants $\mu^{(n)}$. In particular, it shows that μ is the limit of the approximants.

LEMMA 2.14 (MEASURES ARE MONOTONIC). *In the following, e and K range over closed expressions and continuations, and let A range over measurable sets of reals.*

- (1) $\mu(e, K, A) \geq 0$
- (2) for any m , $\mu^{(m)}(e, K, A) \geq 0$
- (3) if $m \leq n$, then $\mu^{(m)}(e, K, A) \leq \mu^{(n)}(e, K, A) \leq \mu(e, K, A)$
- (4) $\mu(e, K, A) = \sup_n \{\mu^{(n)}(e, K, A)\}$

Finally, the next lemma’s equations characterize how the approximant and limit measures, $\mu^{(n)}$ and μ , behave under the reductions of the small-step machine. Almost all the calculations in Section 3 depend only on these equations.

LEMMA 2.15. *The following equations hold for approximant measures:*

$$\begin{aligned} \mu^{(p+1)}(\text{let } x = e_1 \text{ in } e_2, K, A) &= \mu^{(p)}(e_1, (x \rightarrow e_2)K, A) \\ \mu^{(p+1)}(v, (x \rightarrow e)K, A) &= \mu^{(p)}(e[v/x], K, A) \\ \mu^{(p+1)}((\lambda x. e \ v), K, A) &= \mu^{(p)}(e[v/x], K, A) \\ \mu^{(p+1)}(\text{op}^n(v_1, \dots, v_n), K, A) &= \mu^{(p)}(\delta(\text{op}^n, v_1, \dots, v_n), K, A) \quad \text{if defined} \\ \mu^{(p+1)}(\text{if } c_r \text{ then } e_1 \text{ else } e_2, K, A) &= \mu^{(p)}(e_1, K, A) \quad \text{if } r > 0 \\ \mu^{(p+1)}(\text{if } c_r \text{ then } e_1 \text{ else } e_2, K, A) &= \mu^{(p)}(e_2, K, A) \quad \text{if } r \leq 0 \\ \mu^{(p+1)}(\text{sample } c_r, K, A) &= \int_0^1 \mu^{(p)}(c_r, K, A) dr \\ \mu^{(p+1)}(\text{factor } c_r, K, A) &= r \times \mu^{(p)}(c_r, K, A) \quad \text{if } r > 0 \end{aligned}$$

In addition, the analogous index-free equations hold for the unapproximated (limit) measure $\mu(-, -, -)$.

In general, the proofs of the equations of Lemma 2.15 involve unfolding the definition of the measure and applying Lemma 2.3 under the integral. The proof for `let` is representative:

PROOF FOR let.

$$\begin{aligned}
& \mu^{(p+1)}(\text{let } x = e_1 \text{ in } e_2, K, A) \\
&= \iint \text{eval}^{(p+1)}(\sigma, \text{let } x = e_1 \text{ in } e_2, K, \tau, 1, A) \, d\sigma \, d\tau \\
&= \iint \text{eval}^{(p)}(\pi_L(\sigma), e_1, (x \rightarrow e_2)K, \pi_R(\sigma)::\tau, 1, A) \, d\sigma \, d\tau && \text{(Lemma 2.3)} \\
&= \iiint \text{eval}^{(p)}(\sigma', e_1, (x \rightarrow e_2)K, \sigma''::\tau, 1, A) \, d\sigma' \, d\sigma'' \, d\tau && \text{(Property 2.1.4 on } \sigma) \\
&= \iint \text{eval}^{(p)}(\sigma', e_1, (x \rightarrow e_2)K, \pi_L(\tau')::\pi_R(\tau'), 1, A) \, d\sigma' \, d\tau' && \text{(Property 2.1.4 on } \tau') \\
&= \iint \text{eval}^{(p)}(\sigma', e_1, (x \rightarrow e_2)K, \tau', 1, A) \, d\sigma' \, d\tau' && (\pi_L(\tau')::\pi_R(\tau') = \tau') \\
&= \mu^{(p)}(e_1, (x \rightarrow e_2)K, A)
\end{aligned}$$

□

The proof for factor additionally uses linearity (Lemma 2.4), and the proof for sample additionally uses Property 2.1.2.

So far, our semantics speaks directly only about the meanings of whole programs. In the following sections, we develop a collection of relations for expressions and ultimately show that they respect the contextual ordering relation on expression induced by the semantics of whole programs.

3 THE LOGICAL RELATION

In this section, we define a step-indexed logical relation on values, expressions, and continuations, and we prove the Fundamental Property (a form of reflexivity) for our relation.

We begin by defining step-indexed logical relations on *closed* values, *closed* expressions, and continuations (which are always closed) as follows:

$$\begin{aligned}
(v_1, v_2) \in \mathbb{V}_n &\iff v_1 = v_2 = c_r \text{ for some } r \\
&\quad \vee (v_1 = \lambda x. e \wedge v_2 = \lambda x. e' \\
&\quad \quad \wedge (\forall m < n)(\forall v, v')[(v, v') \in \mathbb{V}_m \implies (e[v/x], e'[v'/x]) \in \mathbb{E}_m]) \\
(e, e') \in \mathbb{E}_n &\iff (\forall m \leq n)(\forall K, K')(\forall A \in \Sigma_{\mathbb{R}}) \\
&\quad [(K, K') \in \mathbb{K}_m \implies \mu^{(m)}(e, K, A) \leq \mu^{(m)}(e', K', A)] \\
(K, K') \in \mathbb{K}_n &\iff (\forall m \leq n)(\forall v, v')(\forall A \in \Sigma_{\mathbb{R}}) \\
&\quad [(v, v') \in \mathbb{V}_m \implies \mu^{(m)}(v, K, A) \leq \mu^{(m)}(v', K', A)]
\end{aligned}$$

The definitions are well-founded because \mathbb{V}_- refers to \mathbb{E}_- at strictly smaller indexes. Note that for all n , $\mathbb{V}_n \supseteq \mathbb{V}_{n+1} \supseteq \dots$, and similarly for \mathbb{E} and \mathbb{K} . That is, at higher indexes the relations make more distinctions and thus relate fewer things.

We use γ to range over substitutions of closed values for variables, and we define \mathbb{G}_n by lifting \mathbb{V}_n to substitutions as follows:

$$\begin{aligned}
(\gamma, \gamma') \in \mathbb{G}_n^\Gamma &\iff \text{dom}(\gamma) = \text{dom}(\gamma') = \Gamma \\
&\quad \wedge \forall x \in \Gamma, (\gamma(x), \gamma'(x)) \in \mathbb{V}_n
\end{aligned}$$

Last, we define the logical relations on open terms. In each case, the relation is on terms of the specified sort that are well-formed with free variables in Γ :

$$\begin{aligned}
(v, v') \in \mathbb{V}^\Gamma &\iff (\forall n)(\forall \gamma, \gamma')[(\gamma, \gamma') \in \mathbb{G}_n^\Gamma \implies (v\gamma, v'\gamma') \in \mathbb{V}_n] \\
(e, e') \in \mathbb{E}^\Gamma &\iff (\forall n)(\forall \gamma, \gamma')[(\gamma, \gamma') \in \mathbb{G}_n^\Gamma \implies (e\gamma, e'\gamma') \in \mathbb{E}_n] \\
(K, K') \in \mathbb{K} &\iff (\forall n)(K, K') \in \mathbb{K}_n
\end{aligned}$$

$$\begin{array}{c}
\frac{x \in \Gamma}{(x, x) \in \mathbb{V}^\Gamma} \quad \frac{(e, e') \in \mathbb{E}^{\Gamma, x}}{(\lambda x.e, \lambda x.e') \in \mathbb{V}^\Gamma} \quad (c_r, c_r) \in \mathbb{V}^\Gamma \quad \frac{(v, v') \in \mathbb{V}^\Gamma}{(v, v') \in \mathbb{E}^\Gamma} \\
\\
\frac{(v_1, v'_1) \in \mathbb{V}^\Gamma \quad (v_2, v'_2) \in \mathbb{V}^\Gamma}{((v_1 v_2), (v'_1 v'_2)) \in \mathbb{E}^\Gamma} \quad \frac{(e_1, e'_1) \in \mathbb{E}^\Gamma \quad (e_2, e_2) \in \mathbb{E}^{\Gamma, x}}{(\text{let } x = e_1 \text{ in } e_2, \text{let } x = e'_1 \text{ in } e'_2) \in \mathbb{E}^\Gamma} \\
\\
\frac{(v_i, v'_i) \in \mathbb{V}^\Gamma \quad (i \in \{1, \dots, k\})}{(\text{op}^k (v_1, \dots, v_k), \text{op}^k (v'_1, \dots, v'_k)) \in \mathbb{E}^\Gamma} \\
\\
\frac{(v, v') \in \mathbb{V}^\Gamma \quad (e_1, e'_1) \in \mathbb{E}^\Gamma \quad (e_2, e_2) \in \mathbb{E}^\Gamma}{(\text{if } v \text{ then } e_1 \text{ else } e_2, \text{if } v' \text{ then } e'_1 \text{ else } e'_2) \in \mathbb{E}^\Gamma} \\
\\
(\text{sample, sample}) \in \mathbb{E} \quad \frac{(v, v') \in \mathbb{V}^\Gamma}{(\text{factor } v, \text{factor } v') \in \mathbb{E}^\Gamma} \\
\\
(\text{halt, halt}) \in \mathbb{K} \quad \frac{(e_1, e_2) \in \mathbb{E}^{\{x\}} \quad (K, K') \in \mathbb{K}}{((x \rightarrow e)K, (x \rightarrow e')K') \in \mathbb{K}}
\end{array}$$

Fig. 5. Compatibility rules for the logical relation

The limit relation \mathbb{K} is not indexed by Γ because we work only with closed continuations.

Our first goal is to show the so-called fundamental property of logical relations:

$$\Gamma \vdash e \text{ exp} \implies (e, e) \in \mathbb{E}^\Gamma$$

We begin with a series of compatibility lemmas. These show that the logical relations form a congruence under (“are compatible with”) the scoping rules of values, expressions, and continuations. Note the correspondence between the scoping rules of Figure 2 and the compatibility rules of Figure 5.

LEMMA 3.1 (COMPATIBILITY). *The implications summarized as inference rules in Figure 5 hold.*

Most parts of the lemma follow by general reasoning about the λ -calculus, the definitions of the logical relations, and calculations involving $\mu^{(n)}$ and μ using Lemma 2.15. The proof for application is representative:

PROOF FOR APP. We must show that if $(v_1, v'_1) \in \mathbb{V}^\Gamma$ and $(v_2, v'_2) \in \mathbb{V}^\Gamma$, then $((v_1 v_2), (v'_1 v'_2)) \in \mathbb{E}^\Gamma$.

Choose n , and assume $(\gamma, \gamma') \in \mathbb{G}_n^\Gamma$. Then $(v_1 \gamma, v'_1 \gamma') \in \mathbb{V}_n$ and $(v_2 \gamma, v'_2 \gamma') \in \mathbb{V}_n$. We must show $((v_1 \gamma v_2 \gamma), (v'_1 \gamma' v'_2 \gamma')) \in \mathbb{E}_n$.

If $v_1 \gamma$ is of the form c_r , then $\mu^{(m)}(v_1 \gamma, K, A) = 0$ for any m, K , and A , so the conclusion holds by Lemma 2.14.

Otherwise, assume $v_1 \gamma$ is of the form $\lambda x.e$, and so $v'_1 \gamma'$ is of the form $\lambda x.e'$. So choose $m \leq n$ and A , and let $(K, K') \in \mathbb{K}_m$. We must show that

$$\mu^{(m)}((\lambda x.e \gamma v_2 \gamma), K, A) \leq \mu((\lambda x.e' \gamma' v'_2 \gamma'), K', A).$$

If $m = 0$ the left-hand side is 0 and the inequality holds trivially. So consider $m \geq 1$. Since all the relevant terms are closed and the relations on closed terms are antimonotonic in the index, we have $(\lambda x.e\gamma, \lambda x.e'\gamma') \in \mathbb{V}_m$ and $(v_1\gamma, v_1'\gamma') \in \mathbb{V}_{m-1}$. Therefore $(e\gamma[v_2\gamma/x], e'\gamma'[v_2'\gamma'/x]) \in \mathbb{E}_{m-1}$.

Now, $\langle \sigma \mid (\lambda x.e\gamma \ v_2\gamma) \mid K \mid \tau \mid w \rangle \rightarrow \langle \sigma \mid e\gamma[v_2\gamma/x] \mid K \mid \tau \mid w \rangle$, and similarly for the primed side. So we have

$$\begin{aligned} \mu^{(m)}((\lambda x.e\gamma \ v_2\gamma), K, A) &= \mu^{(m-1)}(e\gamma[v_2\gamma/x], K, A) && \text{(Lemma 2.15)} \\ &\leq \mu(e'\gamma'[v_2'\gamma'/x], K', A) && \text{(by } (e\gamma[v_2\gamma/x], e'\gamma'[v_2'\gamma'/x]) \in \mathbb{E}_{m-1}\text{)} \\ &= \mu((\lambda x.e'\gamma' \ v_2'\gamma'), K', A) \end{aligned}$$

□

More detailed proofs can be found in the long version [Wand et al. 2018, Appendix A].

Now we can prove the Fundamental Property:

THEOREM 3.2 (FUNDAMENTAL PROPERTY).

- (1) $\Gamma \vdash e \text{ exp} \implies (e, e) \in \mathbb{E}^\Gamma$
- (2) $\Gamma \vdash v \text{ val} \implies (v, v) \in \mathbb{V}^\Gamma$
- (3) $\vdash K \text{ cont} \implies \forall n, (K, K) \in \mathbb{K}_n$

PROOF. By induction on the derivation of $\Gamma \vdash e \text{ exp}$, etc, applying the corresponding compatibility rule from Lemma 3.1 at each point. □

The essential properties of the logical relation we wish to hold are soundness and completeness with respect to the contextual ordering. We address these properties in Section 5 after taking a detour to define another useful intermediate relation, CIU^Γ , and establish its equivalence to \mathbb{E}^Γ .

4 CIU ORDERING

The CIU (“closed instantiation of uses”) ordering of two terms asserts that they yield related observable behavior under a single substitution and a single continuation. We take “observable behavior” to be a program’s measure over the reals, as we did for the logical relations.

Definition 4.1.

- (1) If e and e' are closed expressions, then $(e, e') \in \text{CIU}$ iff for all closed K and measurable A , $\mu(e, K, A) \leq \mu(e', K, A)$.
- (2) If $\Gamma \vdash e \text{ exp}$ and $\Gamma \vdash e' \text{ exp}$, then $(e, e') \in \text{CIU}^\Gamma$ iff for all closing substitutions γ , $(e\gamma, e'\gamma) \in \text{CIU}$.

Since it requires considering only a single substitution and a single continuation rather than related pairs, it is often easier to prove particular expressions related by CIU^Γ . But in fact, this relation coincides with the logical relation, as we demonstrate now. One direction is an easy consequence of the Fundamental Property.

LEMMA 4.2 ($\mathbb{E} \subseteq \text{CIU}$). *If $(e, e') \in \mathbb{E}^\Gamma$ then $(e, e') \in \text{CIU}^\Gamma$.*

PROOF. Choose a closing substitution γ , a closed continuation K , and $A \in \Sigma_{\mathbb{R}}$. By the Fundamental Property, we have for all n , $(\gamma, \gamma) \in \mathbb{G}_n^\Gamma$ and $(K, K) \in \mathbb{K}_n$. Therefore, for all n , $\mu^{(n)}(e\gamma, K, A) \leq \mu(e'\gamma, K, A)$. So

$$\mu(e\gamma, K, A) = \sup_n \{\mu^{(n)}(e\gamma, K, A)\} \leq \mu(e'\gamma, K, A).$$

□

In the other direction:

LEMMA 4.3 ($\mathbb{E}^\Gamma \circ \text{CIU}^\Gamma \subseteq \mathbb{E}^\Gamma$). *If $(e_1, e_2) \in \mathbb{E}^\Gamma$ and $(e_2, e_3) \in \text{CIU}^\Gamma$, then $(e_1, e_3) \in \mathbb{E}^\Gamma$.*

PROOF. Choose n and $(\gamma, \gamma') \in \mathbb{G}_\Gamma^n$. We must show that $(e_1\gamma, e_3\gamma') \in \mathbb{E}_n^\Gamma$. So choose $m \leq n$, $(K, K') \in \mathbb{K}_m$, and $A \in \Sigma_{\mathbb{R}}$. Now we must show $\mu^{(m)}(e_1\gamma, K, A) \leq \mu(e_3\gamma', K', A)$.

We have $(e_1, e_2) \in \mathbb{E}^\Gamma$ and $(\gamma, \gamma') \in \mathbb{G}_\Gamma^n$, so $(e_1\gamma, e_2\gamma') \in \mathbb{E}_n$, and by $m \leq n$ we have $(e_1\gamma, e_2\gamma') \in \mathbb{E}_m$. So

$$\begin{aligned} \mu^{(n)}(e_1\gamma, K, A) &\leq \mu(e_2\gamma', K', A) && \text{(by } (e_1\gamma, e_2\gamma') \in \mathbb{E}_m) \\ &\leq \mu(e_3\gamma', K', A) && \text{(by } (e_2, e_3) \in \text{CIU}) \end{aligned}$$

Therefore $(e_1, e_3) \in \mathbb{E}^\Gamma$. □

LEMMA 4.4 ($\text{CIU} \subseteq \mathbb{E}$). *If $(e, e') \in \text{CIU}^\Gamma$ then $(e, e') \in \mathbb{E}^\Gamma$.*

PROOF. Assume $(e, e') \in \text{CIU}^\Gamma$. By the Fundamental Property, we know $(e, e) \in \mathbb{E}^\Gamma$. So we have $(e, e) \in \mathbb{E}^\Gamma$ and $(e, e') \in \text{CIU}^\Gamma$. Hence, by Lemma 4.3, $(e, e') \in \mathbb{E}^\Gamma$. □

THEOREM 4.5. *$(e, e') \in \text{CIU}^\Gamma$ iff $(e, e') \in \mathbb{E}^\Gamma$.*

PROOF. Immediate from Lemmas 4.2 and 4.4. □

5 CONTEXTUAL ORDERING

Finally, we arrive at the contextual order relation. We define the contextual ordering as the largest preorder that is both *adequate*—that is, it distinguishes terms that have different observable behavior by themselves—and *compatible*—that is, closed under context formation, and we show that the contextual ordering, the CIU ordering, and the logical relation all coincide. Thus in order to show two terms contextually ordered, it suffices to use the friendlier machinery of the CIU relation.

Definition 5.1 (CTX^Γ). CTX^Γ is the largest family of relations R^Γ such that:

- (1) R is adequate, that is, if $\Gamma = \emptyset$, then $(e, e') \in R^\Gamma$ implies that for all measurable subsets A of the reals, $\mu(e, \text{halt}, A) \leq \mu(e', \text{halt}, A)$.
- (2) For each Γ , R^Γ is a preorder.
- (3) The family of relations R is compatible, that is, it is closed under the type rules for expressions:
 - (a) If $(e, e') \in R^{\Gamma, x}$, then $(\lambda x.e, \lambda x.e') \in R^\Gamma$.
 - (b) If $(v_1, v'_1) \in R^\Gamma$ and $(v_2, v'_2) \in R^\Gamma$, then $((v_1 v_2), (v'_1 v'_2)) \in R^\Gamma$.
 - (c) If $(v, v') \in R^\Gamma$, then $(\text{factor } v, \text{factor } v') \in R^\Gamma$.
 - (d) If $(e_1, e'_1) \in R^\Gamma$ and $(e_2, e'_2) \in R^{\Gamma, x}$,
then $(\text{let } x = e_1 \text{ in } e_2, \text{let } x = e'_1 \text{ in } e'_2) \in R^\Gamma$.
 - (e) If $(v_1, v'_1) \in R^\Gamma, \dots, (v_n, v'_n) \in R^\Gamma$,
then $(\text{op}^n(v_1, \dots, v_n), \text{op}^n(v'_1, \dots, v'_n)) \in R^\Gamma$.
 - (f) If $(v, v') \in R^\Gamma, (e_1, e'_1) \in R^\Gamma$, and $(e_2, e'_2) \in R^\Gamma$,
then $(\text{if } v \text{ then } e_1 \text{ else } e_2, \text{if } v' \text{ then } e'_1 \text{ else } e'_2) \in R^\Gamma$.

Note, as usual, that the union of any family of relations satisfying these conditions also satisfies these conditions, so the union of all of them is the largest such family of relations.

We prove that \mathbb{E}^Γ , CIU^Γ , and CTX^Γ by first showing that $\mathbb{E}^\Gamma \subseteq \text{CTX}^\Gamma$ and then that $\text{CTX}^\Gamma \subseteq \text{CIU}^\Gamma$. Then, having caught CTX^Γ between \mathbb{E}^Γ and CIU^Γ —two relations that we have already proven equivalent—we conclude that all of the relations coincide.

First, we must show that $\mathbb{E}^\Gamma \subseteq \text{CTX}^\Gamma$. The heart of that proof is showing that \mathbb{E}^Γ is compatible in the sense of Definition 5.1. That is *nearly* handled by the existing compatibility rules for \mathbb{E}^Γ (Lemma 3.1), except for an occasional mismatch between expressions and values—that is, between

\mathbb{E}^Γ and \mathbb{V}^Γ in the rules. So we need a lemma to address the mismatch (Lemma 5.3), which itself needs the following lemma due to Pitts [2010].

LEMMA 5.2. *If $(K, K') \in \mathbb{K}_n$ and $(v, v') \in \mathbb{V}_n$, then*

$$((z \rightarrow (z v))K, (z \rightarrow (z v'))K') \in \mathbb{K}_{n+2}$$

PROOF. See the long version [Wand et al. 2018, Appendix A]. \square

LEMMA 5.3. *For all closed values v , if $(v, v') \in \mathbb{E}$, then $(v, v') \in \mathbb{V}$.*

PROOF. We will show that for all closed values v, v' , if $(v, v') \in \mathbb{E}_{n+3}$, then $(v, v') \in \mathbb{V}_n$, from which the lemma follows.

If $v = c_r$ and $v' = c_{r'}$, then $r = r'$ and thus $(c_r, c_{r'}) \in \mathbb{V}$ because otherwise we would have $\mu(c_r, \text{halt}, \{r\}) = I_{\{r\}}(r) = 1$ and $\mu(c_{r'}, \text{halt}, \{r\}) = I_{\{r\}}(r') = 0$, violating the assumption $(c_r, c_{r'}) \in \mathbb{E}$.

If only one of v and v' is a constant, then $(v, v') \in \mathbb{E}_{n+3}$ is impossible, since constants and lambda-expressions are distinguishable by real? (which requires 3 steps to do so).

So assume $v = \lambda x.e$ and $v' = \lambda x.e'$. To establish $(v, v') \in \mathbb{V}_n$, choose $m < n$ and $(u, u') \in \mathbb{V}_m$. We must show that $(e[u/x], e'[u'/x]) \in \mathbb{E}_m$. To do that, choose $p \leq m$, $(K, K') \in \mathbb{K}_p$, and $A \in \Sigma_{\mathbb{R}}$. We must show that

$$\mu^{(p)}(e[u/x], K, A) \leq \mu(e'[u'/x], K', A)$$

Let $K_1 = (f \rightarrow (f u))K$ and $K'_1 = (f \rightarrow (f u'))K'$. By monotonicity, $(u, u') \in \mathbb{V}_p$. By Lemma 5.2, $(K'_1, K_1) \in \mathbb{K}_{p+2}$. Furthermore, $p \leq m < n$, so $p + 2 \leq n + 1$ and therefore $(\lambda x.e, \lambda x.e') \in \mathbb{E}_{p+2}$. And furthermore, we have

$$\langle \sigma \mid \lambda x.e \mid K_1 \mid \tau \mid w \rangle \rightarrow \langle \sigma \mid (\lambda x.e u) \mid K \mid \tau \mid w \rangle \rightarrow \langle \sigma \mid e[u/x] \mid K \mid \tau \mid w \rangle$$

and similarly on the primed side.

We can put the results together to get

$$\begin{aligned} \mu^{(p)}(e[u/x], K, A) &= \mu^{(p+2)}(\lambda x.e, K_1, A) \\ &\leq \mu(\lambda x.e', K'_1, A) \\ &= \mu(e'[u'/x], K', A) \end{aligned}$$

\square

THEOREM 5.4. $\mathbb{E}^\Gamma \subseteq \text{CTX}^\Gamma$.

PROOF. We will show that \mathbb{E} forms a family of reflexive preorders that is adequate and compatible. Each \mathbb{E}^Γ is reflexive by the Fundamental Property, and is a preorder because it is equal to CIU^Γ , which is a preorder. To show that it is adequate, observe that $(\text{halt}, \text{halt}) \in \mathbb{K}$ by Lemma 3.1, hence for any measurable subset A of reals, $(e, e') \in \mathbb{E}^\Gamma$ implies $\mu(e, \text{halt}, A) = \mu(e', \text{halt}, A)$.

The \mathbb{E} -compatibility rules (Lemma 3.1) are almost exactly what is needed for CTX -compatibility. The exceptions are in the application, operation, if, and factor rules where their hypotheses refer to \mathbb{V}^Γ rather than \mathbb{E}^Γ . We fill the gap with Lemma 5.3. We show how this is done for factor v ; the other cases are similar.

$$\begin{aligned}
(v, v') \in \mathbb{E}^\Gamma &\implies (v, v') \in \mathbb{CIU}^\Gamma \\
&\implies (\forall \gamma)((v\gamma, v'\gamma) \in \mathbb{CIU}^0) \\
&\implies (\forall \gamma)((v\gamma, v'\gamma) \in \mathbb{E}^0) \\
&\implies (\forall \gamma)((v\gamma, v'\gamma) \in \mathbb{V}^0) && \text{(Lemma 5.3)} \\
&\implies (\forall \gamma)((\text{factor } v\gamma, \text{factor } v'\gamma) \in \mathbb{E}^0) && \text{(Lemma 3.1)} \\
&\implies (\forall \gamma)((\text{factor } v\gamma, \text{factor } v'\gamma) \in \mathbb{CIU}^0) \\
&\implies (\text{factor } v, \text{factor } v') \in \mathbb{CIU}^\Gamma \\
&\implies (\text{factor } v, \text{factor } v') \in \mathbb{E}^\Gamma
\end{aligned}$$

□

Next, we must show that $\mathbb{CTX}^\Gamma \subseteq \mathbb{CIU}^\Gamma$ by induction on the closing substitution and then induction on the continuation. We use the following two lemmas to handle the closing substitution.

LEMMA 5.5. *If $\Gamma, x \vdash e \text{ exp}$ and $\Gamma \vdash v \text{ exp}$, then*

$$(e[v/x], (\lambda x. e v)) \in \mathbb{CIU}^\Gamma \text{ and } ((\lambda x. e v), e[v/x]) \in \mathbb{CIU}^\Gamma.$$

PROOF. Let γ be a closing substitution for Γ . Then for any σ , closed K , and w , by Lemmas 2.15 and 2.14.4 we have

$$\langle \sigma \mid (\lambda x. e\gamma v\gamma) \mid K \mid \tau \mid w \rangle \rightarrow \langle \sigma \mid e\gamma[v\gamma/x] \mid K \mid \tau \mid w \rangle$$

Therefore for any $A \in \Sigma_{\mathbb{R}}$, $\mu((\lambda x. e\gamma v\gamma), K, A) = \mu(e\gamma[v\gamma/x], K, A)$. □

LEMMA 5.6. *If $(e, e') \in \mathbb{CTX}^{\Gamma, x}$, and $(v, v') \in \mathbb{CTX}^\Gamma$, then $(e[v/x], e'[v'/x]) \in \mathbb{CTX}^\Gamma$.*

PROOF. From the assumptions and the compatibility of \mathbb{CTX} , we have

$$((\lambda x. e v), (\lambda x. e' v')) \in \mathbb{CTX}^\Gamma \tag{1}$$

So now we have:

$$\begin{aligned}
(e[v/x], (\lambda x. e v)) &\in \mathbb{CIU}^\Gamma && \text{(Lemma 5.5)} \\
\implies (e[v/x], (\lambda x. e v)) &\in \mathbb{CTX}^\Gamma && (\mathbb{CIU}^\Gamma \subseteq \mathbb{CTX}^\Gamma) \\
\implies (e[v/x], (\lambda x. e' v')) &\in \mathbb{CTX}^\Gamma && \text{(Equation (1) and transitivity of } \mathbb{CTX}^\Gamma) \\
\implies (e[v/x], e'[v'/x]) &\in \mathbb{CTX}^\Gamma && \text{(Lemma 5.5 and transitivity of } \mathbb{CTX}^\Gamma)
\end{aligned}$$

□

Now we are ready to complete the theorem. Here we need to use \mathbb{CIU} rather than \mathbb{E} , so that we can deal with only one continuation rather than two.

THEOREM 5.7 ($\mathbb{CTX}^\Gamma \subseteq \mathbb{CIU}^\Gamma$). *If $(e, e') \in \mathbb{CTX}^\Gamma$, then $(e, e') \in \mathbb{CIU}^\Gamma$*

PROOF. By the preceding lemma, we have $(e\gamma, e'\gamma) \in \mathbb{CTX}$. So it suffices to show that for all $A \in \Sigma_{\mathbb{R}}$, if $(e, e') \in \mathbb{CTX}^0$ and $\vdash K \text{ cont}$, then $\mu(e, K, A) = \mu(e', K, A)$.

The proof proceeds by induction on K such that $\vdash K \text{ cont}$. The induction hypothesis on K is: for all closed e, e' , if $(e, e') \in \mathbb{CTX}^0$, then $\mu(e, K, A) = \mu(e', K, A)$.

If $K = \text{halt}$ and $(e, e') \in \mathbb{CTX}^0$, then $\mu(e, \text{halt}, A) = \mu(e', \text{halt}, A)$ by the adequacy of \mathbb{CTX}^0 .

$$\begin{aligned}
((\lambda x.e) v) &=_{\text{ctx}} e[v/x] && (\beta_v) \\
\text{let } x = v \text{ in } e &=_{\text{ctx}} e[v/x] && (\text{let}_v) \\
\text{let } x = e \text{ in } x &=_{\text{ctx}} e && (\text{let}_{id}) \\
op(v_1, \dots, v_n) &=_{\text{ctx}} v \quad \text{where } \delta(op, v_1, \dots, v_n) = v && (\delta) \\
\text{let } x_2 = (\text{let } x_1 = e_1 \text{ in } e_2) \text{ in } e_3 &=_{\text{ctx}} \text{let } x_1 = e_1 \text{ in } (\text{let } x_2 = e_2 \text{ in } e_3) && (\text{assoc}) \\
\text{let } x_1 = e_1 \text{ in } \text{let } x_2 = e_2 \text{ in } e_3 &=_{\text{ctx}} \text{let } x_2 = e_2 \text{ in } \text{let } x_1 = e_1 \text{ in } e_3 && (\text{commut})
\end{aligned}$$

In (assoc), $x_1 \notin FV(e_3)$. In (commut), $x_1 \notin FV(e_2)$ and $x_2 \notin FV(e_1)$.

Fig. 6. An incomplete catalog of equivalences

For the induction step, consider $(x \rightarrow e_1)K$, where $x \vdash e_1$ exp. Choose $(e, e') \in \text{CTX}^0$. We must show $\mu(e, (x \rightarrow e_1)K, A) \leq \mu(e', (x \rightarrow e_1)K, A)$.

By the compatibility of CTX , we have

$$(\text{let } x = e \text{ in } e_1, \text{let } x = e' \text{ in } e_1) \in \text{CTX}^0 \quad (2)$$

Then we have

$$\begin{aligned}
\mu(e, (x \rightarrow e_1)K, A) &= \mu(\text{let } x = e \text{ in } e_1, K, A) && (\text{Lemma 2.15}) \\
&\leq \mu(\text{let } x = e' \text{ in } e_1, K, A) && (\text{by IH at } K, \text{ applied to (2)}) \\
&= \mu(e', (x \rightarrow e_1)K, A) && (\text{Lemma 2.15})
\end{aligned}$$

Thus completing the induction step. \square

Summarizing the results:

THEOREM 5.8. *For all Γ , $\text{CIU}^\Gamma = \mathbb{E}^\Gamma = \text{CTX}^\Gamma$.*

PROOF. $\text{CIU}^\Gamma = \mathbb{E}^\Gamma \subseteq \text{CTX}^\Gamma \subseteq \text{CIU}^\Gamma$ by Theorems 4.5, 5.4, and 5.7, respectively. \square

6 CONTEXTUAL EQUIVALENCE

Definition 6.1. If $\Gamma \vdash e$ exp and $\Gamma \vdash e'$ exp, we say e and e' are *contextually equivalent* ($e =_{\text{ctx}} e'$) if both $(e, e') \in \text{CTX}^\Gamma$ and $(e', e) \in \text{CTX}^\Gamma$.

In this section we use the machinery from the last few sections to prove several equivalence schemes. The equivalences fall into three categories:

- (1) provable directly using CIU and Theorem 5.8
- (2) dependent on “entropy-shuffling”
- (3) mathematical properties of \mathbb{R} , probability distributions, etc

Some equivalences of the first and second kinds are listed in Figure 6; Section 6.4 gives some examples of the third kind.

6.1 β_v , let_v , and δ

The proof for β_v demonstrates the general pattern of equivalence proofs using CIU: first we prove the equation holds for closed expressions, then we generalize to open terms by considering all closing substitutions.

LEMMA 6.2. *If $\vdash ((\lambda x.e) v)$ exp, then $((\lambda x.e) v) =_{\text{ctx}} e[v/x]$.*

PROOF. By Lemma 2.15, the definition of CIU, and Theorem 5.8. \square

COROLLARY 6.3 (β_v). If $\Gamma \vdash ((\lambda x.e) v) \text{ exp}$, then $((\lambda x.e) v) =_{\text{ctx}} e[v/x]$.

PROOF. By Lemma 6.2, any closed instances of these expressions are contextually equivalent and thus CIU-equivalent. Hence the open expressions are CIU-equivalent and thus contextually equivalent. \square

The proofs of let_v and δ are similar.

6.2 Rearranging Entropy

The remaining equivalences from Figure 6 involve non-trivial changes to the entropy access patterns of their subexpressions. In this section we characterize a class of transformations on the entropy space that are measure-preserving. In the next section we use these functions to justify reordering and rearranging subexpression evaluation.

Definition 6.4 (measure-preserving). A function $\phi : \mathbb{S} \rightarrow \mathbb{S}$ is measure-preserving when for all measurable $g : \mathbb{S} \rightarrow \mathbb{R}^+$,

$$\int g(\phi(\sigma)) d\sigma = \int g(\sigma) d\sigma$$

Note that this definition is implicitly specific to the stock entropy measure $\mu_{\mathbb{S}}$, which is sufficient for our needs.

More specifically, the kinds of functions we are interested in are ones that break apart the entropy into independent pieces using π_L and π_R and then reassemble the pieces of entropy using $::$. Pieces may be discarded, but no piece may be used more than once.

For example, the following function is measure-preserving:

$$\phi_c(\sigma_1 :: (\sigma_2 :: \sigma_3)) = \sigma_2 :: (\sigma_1 :: \sigma_3)$$

Or equivalently, written using explicit projections:

$$\phi_c(\sigma) = \pi_L(\pi_R(\sigma)) :: (\pi_L(\sigma) :: \pi_R(\pi_R(\sigma)))$$

We will use this function in Theorem 6.9 to justify let -reordering.

To characterize such functions, we need some auxiliary definitions:

- A *path* $p = [d_1, \dots, d_n]$ is a (possibly empty) list of directions (L or R). It represents a sequence of projections, and it can be viewed as a function from \mathbb{S} to \mathbb{S} .

$$[d_1, \dots, d_n](\sigma) = (\pi_{d_1} \circ \dots \circ \pi_{d_n})(\sigma)$$

- A *finite shuffling function* (FSF) ϕ is either a path or $\phi_1 :: \phi_2$ where ϕ_1 and ϕ_2 are FSFs. It represents the disassembly and reassembly of entropy, and it can be viewed as a recursively defined function from \mathbb{S} to \mathbb{S} .

$$\phi(\sigma) = \begin{cases} p(\sigma) & \text{if } \phi = p \\ \phi_1(\sigma) :: \phi_2(\sigma) & \text{if } \phi = \phi_1 :: \phi_2 \end{cases}$$

- A sequence of paths is said to be *non-duplicating* if no path is the suffix of another path in the sequence.
- An FSF is said to be *non-duplicating* if the sequence of paths appearing in its definition is non-duplicating.

LEMMA 6.5. Let p_1, \dots, p_n be a non-duplicating sequence of paths and $g : \mathbb{S}^n \rightarrow \mathbb{R}^+$. Then

$$\int g(p_1(\sigma), \dots, p_n(\sigma)) d\sigma = \int \dots \int g(\sigma_1, \dots, \sigma_n) d\sigma_1 \dots d\sigma_n$$

PROOF. By strong induction on the length of the longest path in the sequence, and by the definition of non-duplicating and Lemma 2.2 (Tonelli). \square

THEOREM 6.6. *If ϕ is a non-duplicating FSF then ϕ is measure preserving.*

PROOF. We need to show that for any $g : \mathbb{S} \rightarrow \mathbb{R}^+$,

$$\int g(\phi(\sigma)) d\sigma = \int g(\sigma'') d\sigma''$$

If ϕ has paths p_1, \dots, p_n , then we can decompose ϕ using $s : \mathbb{S}^n \rightarrow \mathbb{S}$ such that

$$\phi(\sigma) = s(p_1(\sigma), \dots, p_n(\sigma))$$

where the p_i are non-duplicating. Then by Lemma 6.5 it is enough to show that

$$\int \dots \int g(s(\sigma_1, \dots, \sigma_n)) d\sigma_1 \dots d\sigma_n = \int g(\sigma'') d\sigma''$$

We proceed by induction on ϕ .

- case $\phi = p$. This means that $n = 1$ and s is the identity function, so the equality holds trivially.
- case $\phi = \phi_1 :: \phi_2$. If m is the number of paths in ϕ_1 , then there must be $s_1 : \mathbb{S}^m \rightarrow \mathbb{S}$ and $s_2 : \mathbb{S}^{n-m} \rightarrow \mathbb{S}$ such that

$$s(\sigma_1, \dots, \sigma_m, \sigma_{m+1}, \dots, \sigma_n) = s_1(\sigma_1, \dots, \sigma_m) :: s_2(\sigma_{m+1}, \dots, \sigma_n)$$

We can conclude that

$$\begin{aligned} & \int \dots \int g(s(\sigma_1, \dots, \sigma_n)) d\sigma_1 \dots d\sigma_n \\ &= \int \dots \int g(s_1(\sigma_1, \dots, \sigma_m) :: s_2(\sigma_{m+1}, \dots, \sigma_n)) d\sigma_1 \dots d\sigma_n \\ &= \iint g(\sigma :: \sigma') d\sigma d\sigma' && \text{(IH twice)} \\ &= \int g(\sigma'') d\sigma'' && \text{(Property 2.1(4))} \end{aligned}$$

□

6.3 Equivalences That Depend on Rearranging Entropy

We first prove a general theorem relating value-preserving transformations on the entropy space:

THEOREM 6.7. *Let e and e' be closed expressions, and let $\phi : \mathbb{S} \rightarrow \mathbb{S}$ be a measure-preserving transformation such that for all σ, K, τ , and A*

$$\text{eval}(\sigma, e, K, \tau, 1, A) \leq \text{eval}(\phi(\sigma), e', K, \tau, 1, A)$$

Then $(e, e') \in \text{CTX}$.

PROOF. Without loss of generality, assume e and e' are closed (otherwise apply a closing substitution). By Theorem 5.8, it is sufficient to show that for any K and A , $\mu(e, K, A) \leq \mu(e', K, A)$. We calculate:

$$\begin{aligned} \mu(e, K, A) &= \iint \text{eval}(\sigma, e, K, \tau, 1, A) d\sigma d\tau \\ &\leq \iint \text{eval}(\phi(\sigma), e', K, \tau, 1, A) d\sigma d\tau \\ &= \iint \text{eval}(\sigma, e', K, \tau, 1, A) d\sigma d\tau && (\phi \text{ is measure-preserving}) \\ &= \mu(e', K, A) \end{aligned}$$

□

THEOREM 6.8. *Let e and e' be closed expressions, and let $\phi : \mathbb{S} \rightarrow \mathbb{S}$ be a measure-preserving transformation such that for all v and w ,*

$$\sigma \vdash e \Downarrow v, w \implies \phi(\sigma) \vdash e' \Downarrow v, w.$$

Then $(e, e') \in \text{CTX}$.

PROOF. We will use Theorem 6.7. Assume $\text{eval}(\sigma, e, K, \tau, 1, A) = r > 0$. Hence by Theorem 2.11, there exist quantities v', w', σ' , and τ' such that

$$\langle \sigma \mid e \mid K \mid \tau \mid 1 \rangle \rightarrow^* \langle \sigma' \mid v' \mid \text{halt} \mid \tau' \mid r \rangle$$

with $v' \in A$. By Theorem 2.11 there exist v'', σ'' , and w'' such that

$$\sigma \vdash e \Downarrow v'', w'' \text{ and } \langle \sigma'' \mid v'' \mid K \mid \tau \mid w'' \rangle \rightarrow^* \langle \sigma' \mid v' \mid \text{halt} \mid \tau' \mid r \rangle$$

By the assumption of the theorem, we have $\phi(\sigma) \vdash e' \Downarrow v'', w''$.

Therefore, by Theorem 2.5, there is a σ''' such that

$$\langle \sigma' \mid e \mid K \mid t \mid 1 \rangle \rightarrow^* \langle \sigma''' \mid v'' \mid K \mid \tau \mid w'' \rangle$$

We claim that $\text{eval}(\phi(s), e', K, \tau, 1, A) = r$. Proceed by cases on K . We have $\langle \sigma'' \mid v'' \mid K \mid \tau \mid w'' \rangle \rightarrow^* \langle \phi(s) \mid v' \mid \text{halt} \mid \tau' \mid r \rangle$. If $K = \text{halt}$, this reduction must have length 0. Therefore $v'' = v' \in A$ and $w'' = r$, so $\text{eval}(\phi(s), e', K, \tau, 1, A) = r$.

Otherwise assume $K = (x \rightarrow e_3)K'$. Then both $\langle \sigma'' \mid v'' \mid K \mid t \mid w'' \rangle$ and $\langle \sigma''' \mid v'' \mid K \mid t \mid w'' \rangle$ take a step to $\langle \pi_L(\tau) \mid e_3[v''/x] \mid K' \mid \pi_R(\tau) \mid w'' \rangle$, so $\text{eval}(\phi(s), e', K, \tau, 1, A) = \text{eval}(\sigma, e, K, \tau, 1, A) = r$, as desired, thus establishing the requirement of Theorem 6.7. \square

Now we can finally prove the commutativity theorem promised at the beginning.

THEOREM 6.9. *Let e_1 and e_2 be closed expressions, and $\{x_1, x_2\} \vdash e_0 \text{ exp}$. Then the expressions*

$$\text{let } x_1 = e_1 \text{ in let } x_2 = e_2 \text{ in } e_0$$

and

$$\text{let } x_2 = e_2 \text{ in let } x_1 = e_1 \text{ in } e_0$$

are contextually equivalent.

PROOF USING BIG-STEP SEMANTICS. Let e and e' denote the two expressions of the theorem. We will use Theorem 6.8 with the function $\phi_c(\sigma_1::(\sigma_2::\sigma_3)) = \sigma_2::(\sigma_1::\sigma_3)$, which preserves entropy as shown in the preceding section. We will show that if $\sigma \vdash e \Downarrow v, w$, then $\phi(\sigma) \vdash e \Downarrow v, w$.

Inverting $\sigma \vdash e \Downarrow v, w$, we know there must be a derivation

$$\frac{\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1 \quad \frac{\pi_L(\pi_R(\sigma)) \vdash e_2 \Downarrow v_2, w_{21} \quad \pi_R(\pi_R(\sigma)) \vdash e_0[v_1/x_1][v_2/x_2] \Downarrow v, w_{22}}{\pi_R(\sigma) \vdash \text{let } x_2 = e_2 \text{ in } e_0 \Downarrow v, w_2}}{\sigma \vdash \text{let } x_1 = e_1 \text{ in let } x_2 = e_2 \text{ in } e_0 \Downarrow v, w}$$

where $w = w_1 \times w_2 = w_1 \times (w_{21} \times w_{22})$.

Since e_1 and e_2 are closed, they evaluate to closed v_1 and v_2 , and so the substitutions $[v_1/x_1]$ and $[v_2/x_2]$ commute. Using that and the associativity and commutativity of multiplication, we can rearrange the pieces to get

$$\frac{\pi_L(\pi_R(\sigma)) \vdash e_2 \Downarrow v_2, w_{21} \quad \frac{\pi_L(\sigma) \vdash e_1 \Downarrow v_1, w_1 \quad \pi_R(\pi_R(\sigma)) \vdash e_0[v_2/x_2][v_1/x_1] \Downarrow v, w_{22}}{\pi_L(\sigma)::\pi_R(\pi_R(\sigma)) \vdash \text{let } x_1 = e_1 \text{ in } e_0 \Downarrow v, w_1 \times w_{22}}}{\pi_L(\pi_R(\sigma))::\pi_L(\sigma)::\pi_R(\pi_R(\sigma)) \vdash \text{let } x_2 = e_2 \text{ in let } x_1 = e_1 \text{ in } e_0 \Downarrow v, w}$$

The entropy in the last line is precisely $\phi(\sigma)$, so the requirement of Theorem 6.8 is established. \square

Theorem 6.9 can also be proven directly from the small-step semantics using the interpolation and genericity theorems (2.8 and 2.9) to recover the structure that the big-step semantics makes explicit. The proof may be found in the long version [Wand et al. 2018, Appendix A].

COROLLARY 6.10 (COMMUTATIVITY). *Let e_1 and e_2 be expressions such that x_1 is not free in e_2 and x_2 is not free in e_1 . Then*

$$(\text{let } x_1 = e_1 \text{ in let } x_2 = e_2 \text{ in } e_0) =_{ctx} (\text{let } x_2 = e_2 \text{ in let } x_1 = e_1 \text{ in } e_0)$$

PROOF. Same as Corollary 6.3: since all of the closed instances are equivalent by Theorem 6.9, the open expressions are equivalent. \square

The proofs of let-associativity and let_{id} follow the same structure, except that associativity uses $\phi_a((\sigma_1::\sigma_2)::\sigma_3) = \sigma_1::(\sigma_2::\sigma_3)$ and let_{id} uses $\phi_i(\sigma_1::\sigma_2) = \sigma_1$.

6.4 Quasi-Denotational Reasoning

In this section we give a powerful “quasi-denotational” reasoning tool that shows that if two expressions denote the same measure, they are contextually equivalent. This allows us to import mathematical facts about real arithmetic and probability distributions.

To support this kind of reasoning, we need a notion of measure for a (closed) expression independent of a program continuation. We define $\hat{\mu}(e, -)$ as a measure over arbitrary *syntactic values*—not just real numbers as with $\mu(e, K, -)$. This measure corresponds directly to the μ_e of Culpepper and Cobb [2017] and $\llbracket e \rrbracket_S$ of Borgström et al. [2017]. The definition of $\hat{\mu}$ uses a generalization of eval from measurable sets of reals (A) to measurable sets of syntactic values (V). This requires a measurable space for syntactic values; we take the construction of Borgström et al. [2017, Figure 5] mutatis mutandis.

Definition 6.11.

$$\hat{\mu}(e, V) = \iint \text{eval}(\sigma, e, \text{halt}, \tau, 1, V) d\sigma d\tau$$

$$\text{eval}(\sigma, e, K, \tau, w, V) = \begin{cases} w' & \text{if } \langle \sigma \mid e \mid K \mid \tau \mid w \rangle \rightarrow^* \langle \sigma' \mid v \mid \text{halt} \mid \tau' \mid w' \rangle, \\ & \text{where } v \in V \\ 0 & \text{otherwise} \end{cases}$$

Our goal is to relate an expression’s measure $\hat{\mu}(e, -)$ with the measure of that expression with a program continuation ($\mu(e, K, -)$). Then if two expressions have the same measures, we can use CIU to show them contextually equivalent.

First we need a lemma about decomposing evaluations. It is easiest to state if we define the value and weight projections of evaluation:

$$\text{ev}(\sigma, e, K, \tau) = \begin{cases} v & \text{when } \langle \sigma \mid e \mid K \mid \tau \mid 1 \rangle \rightarrow^* \langle \sigma' \mid v \mid \text{halt} \mid \tau' \mid w \rangle \\ \perp & \text{otherwise} \end{cases}$$

$$\text{ew}(\sigma, e, K, \tau) = \begin{cases} w & \text{when } \langle \sigma \mid e \mid K \mid \tau \mid 1 \rangle \rightarrow^* \langle \sigma' \mid v \mid \text{halt} \mid \tau' \mid w \rangle \\ 0 & \text{otherwise} \end{cases}$$

Note that

$$\text{eval}(\sigma, e, K, \tau, 1, V) = I_V(\text{ev}(\sigma, e, K, \tau)) \times \text{ew}(\sigma, e, K, \tau)$$

$$\hat{\mu}(e, A) = \mu(e, \text{halt}, A) \quad \text{for } A \in \Sigma_{\mathbb{R}}$$

where I_V is the characteristic function of V .

LEMMA 6.12.

$$\text{ev}(\sigma, e, K, \tau) = \text{ev}(\sigma', \text{ev}(\sigma, e, \text{halt}, \tau'), K, \tau)$$

$$\text{ew}(\sigma, e, K, \tau) = \text{ew}(\sigma', \text{ev}(\sigma, e, \text{halt}, \tau'), K, \tau) \times \text{ew}(\sigma, e, \text{halt}, \tau')$$

PROOF. By reduction-sequence surgery using Theorem 2.9. Note that the primed variables are dead: σ' because `ev` returns a value and τ' because `halt` does not use its entropy stack. \square

Next we need a lemma from measure theory:

LEMMA 6.13. *If μ and ν are measures and $\nu(A) = \int I_A(f(x)) \times w(x) \mu(dx)$, then*

$$\int g(y) \nu(dy) = \int g(f(x)) \times w(x) \mu(dx)$$

PROOF. By the pushforward and Radon-Nikodym lemmas from measure theory. \square

Now we are ready for the main theorem, which says that $\mu(e, K, -)$ can be expressed as an integral over $\hat{\mu}(e, -)$ where K appears only in the integrand and e appears only in the measure of integration.

THEOREM 6.14.

$$\mu(e, K, A) = \iiint \text{eval}(\sigma, v, K, \tau, 1, A) \hat{\mu}(e, dv) d\sigma d\tau$$

PROOF. By integral calculations and Lemma 6.12:

$$\begin{aligned} \mu(e, K, A) &= \iint \text{eval}(\sigma, e, K, \tau, 1, A) d\sigma d\tau \\ &= \iint I_A(\text{ev}(\sigma, e, K, \tau)) \times \text{ew}(\sigma, e, K, \tau) d\sigma d\tau \\ &= \iiint I_A(\text{ev}(\sigma, e, K, \tau)) \times \text{ew}(\sigma, e, K, \tau) d\sigma' d\tau' d\sigma d\tau && (\mu_{\mathbb{S}}(\mathbb{S}) = 1) \\ &= \iiint I_A(\text{ev}(\sigma', \text{ev}(\sigma, e, \text{halt}, \tau'), K, \tau)) \times \text{ew}(\sigma', \text{ev}(\sigma, e, \text{halt}, \tau'), K, \tau) \\ &\quad \times \text{ew}(\sigma, e, \text{halt}, \tau') d\sigma' d\tau' d\sigma d\tau && (\text{Lemma 6.12}) \\ &= \iint I_A(\text{ev}(\sigma', v, K, \tau)) \times \text{ew}(\sigma', v, K, \tau) \hat{\mu}(e, dv) d\sigma' d\tau && (\text{Lemma 6.13}) \\ &= \iiint \text{eval}(\sigma', v, K, \tau, 1, A) \hat{\mu}(e, dv) d\sigma' d\tau \end{aligned}$$

\square

As a consequence, two real-valued expressions are contextually equivalent if their expression measures agree:

THEOREM 6.15 ($\hat{\mu}$ IS QUASI-DENOTATIONAL). *If e and e' are closed expressions such that*

- *e and e' are almost always real-valued—that is, $\hat{\mu}(e, \text{Values} - \mathbb{R}) = 0$ and likewise for e' —and*
- *for all $A \in \Sigma_{\mathbb{R}}$, $\hat{\mu}(e, A) = \hat{\mu}(e', A)$*

then $e =_{\text{ctx}} e'$.

PROOF. The two conditions together imply that $\hat{\mu}(e, -) = \hat{\mu}(e', -)$.

We use Theorem 5.8; we must show $(e, e') \in \mathbb{C}\mathbb{I}\mathbb{U}$ and $(e', e) \in \mathbb{C}\mathbb{I}\mathbb{U}$. Choose a continuation K and a measurable set $A \in \Sigma_{\mathbb{R}}$. Then

$$\begin{aligned} \mu(e, K, A) &= \iiint \text{eval}(\sigma, v, K, \tau, 1, A) \hat{\mu}(e, dv) d\sigma d\tau && (\text{by Lemma 6.14}) \\ &= \iiint \text{eval}(\sigma, v, K, \tau, 1, A) \hat{\mu}(e', dv) d\sigma d\tau && (\hat{\mu}(e, -) = \hat{\mu}(e', -)) \\ &= \mu(e', K, A) && (\text{by Lemma 6.14 again}) \end{aligned}$$

The proof of $(e', e) \in \mathbb{C}\mathbb{I}\mathbb{U}$ is symmetric. \square

Theorem 6.15 allows us to import many useful facts from mathematics about real numbers, real operations, and real-valued probability distributions. For example, here are a few equations useful in the transformation of the linear regression example from Section 1:

- $x + y = y + x$

- $(y + x) - z = x - (z - y)$
- (factor x); (factor y) = (factor $x*y$); y
- $\text{normalpdf}(x - y; 0, s) = \text{normalpdf}(y; x, s)$
- The closed-form posterior and normalizer for a normal observation with normal conjugate prior [Murphy 2007]:

```

let m = normal(m0, s0) in
let _ = factor normalpdf(d; m, s) in
m

let m = normal( (1/s02 + 1/s2)-1 (m0/s02 + d/s2), (1/s02 + 1/s2)-1/2 ) in
= let _ = factor normalpdf(d; m0, (s02 + s2)1/2) in
m

```

Note that we must keep the normalizer (the marginal likelihood of d); it is needed to score the hyper-parameters m_0 and s_0 .

Section 7.2 contains an additional application of Theorem 6.15.

6.5 An Application

Recall the example program from the introduction and the proposed transformation:

```

A = normal(0, 10)
B = normal(0, 10)
f(x) = A*x + B
factor normalpdf(f(2) - 2.4; 0, 1)
factor normalpdf(f(3) - 2.7; 0, 1)
factor normalpdf(f(4) - 3.0; 0, 1)

```

→

```

A = normal(0, 10)
factor Z(A)
B = normal(M(A), S(A))

```

The core of the transformation is the last equivalence from Section 6.4, which transforms an observation with a conjugate prior into the posterior and normalizer (which scores the prior's hyperparameters). But applying that transformation requires auxiliary steps to focus the program into the requisite shape:

- Inline f to expose the dependence of the observations on B .
- Rewrite the observations to the form $\text{normalpdf}(_; B, _)$ using facts about arithmetic and normalpdf .
- Reassociate the (implicit) lets to isolate the definition of B and the first observation from the rest of the program's main let chain.

That sets the stage for the application of the conjugacy transformation for one observation. Additional shuffling is required to process subsequent observations. Moreover, each of the mathematical rewrite rules needs help from the rules of Figure 6 to manage the intermediate let bindings required by our language's syntax.

An alternative transformation strategy is to combine the observations beforehand using equations about products of normal densities. The same preliminary transformations are necessary, but the observation-processing loop is eliminated.

6.6 Other Equivalences

The list of equivalences presented in this section is not exhaustive. On the λ -calculus side, we focused on a few broadly applicable rules that involve only syntactic restrictions on the sub-expressions—specifically, constraints on free variables. There are other equivalences that require additional

semantic constraints. For example, a `let`-binding is useless and can be dropped if the right-hand side has an expression measure of weight 1; that is, it nearly always terminates without an error and it does not (effectively) use `factor`. Similarly, hoisting an argument-invariant expression out of a function body requires the same conditions, and the expression must also be deterministic.

On the domain-specific side, Theorem 6.15 works well for programs that contain first-order islands of sampling, scoring, and mathematical operations. But programs with mathematics tangled up with higher-order code, it would be necessary to find either a method of detangling them or a generalization to higher-order expression measures.

7 FORMALLY RELATED WORK

Our language model differs from other models of probabilistic languages, such as that of Borgström et al. [2016], in the following ways. Our language

- uses *splitting* rather than *sequenced* entropy,
- requires `let`-binding of nontrivial intermediate expressions, and
- directly models only the standard uniform distribution.

These differences, while they make our proofs easier, do *not* amount to fundamental differences in the meaning of probabilistic programs. In this section, we show how our semantics corresponds to other formulations.

7.1 Splitting versus Sequenced Entropy

Let the sequenced entropy space \mathbb{T} be the space of finite sequences (“traces”) of real numbers [Borgström et al. 2017, Section 3.3]:

$$\mathbb{T} = \bigcup_{n \geq 0} \mathbb{R}^n$$

Its stock measure $\mu_{\mathbb{T}}$ is the sum of the standard Lebesgue measures on \mathbb{R}^n (but restricted to the Borel algebras on \mathbb{R}^n rather than their completions with negligible sets). Note that $\mu_{\mathbb{T}}$ is infinite.

We write ϵ for the empty sequence and $r::t$ for the sequence consisting of r followed by the elements of t . Integration with respect to $\mu_{\mathbb{T}}$ has the following property:

$$\int f(t) \mu_{\mathbb{T}}(dt) = f(\epsilon) + \iint f(r::t) \mu_{\mathbb{T}}(dt) \lambda(dr)$$

We define $\dot{\rightarrow}$, $\text{e}\dot{\text{v}}\text{al}(t, e, K, w, A)$, and $\dot{\mu}(e, K, A)^2$ as the sequenced-entropy analogues of \rightarrow , `eval`, and μ . Here are some representative rules of $\dot{\rightarrow}$:

$$\begin{aligned} \langle t \mid \text{let } x = e_1 \text{ in } e_2 \mid K \mid w \rangle &\dot{\rightarrow} \langle t \mid e_1 \mid (x \rightarrow e_2)K \mid w \rangle \\ \langle t \mid v \mid (x \rightarrow e_2)K \mid w \rangle &\dot{\rightarrow} \langle t \mid e_2[v/x] \mid K \mid w \rangle \\ \langle r::t \mid \text{sample} \mid K \mid w \rangle &\dot{\rightarrow} \langle t \mid c_r \mid K \mid w \rangle \quad (\text{when } 0 \leq r \leq 1) \\ \langle t \mid \text{factor } c_r \mid K \mid w \rangle &\dot{\rightarrow} \langle t \mid c_r \mid K \mid w \times r \rangle \quad (\text{when } r > 0) \end{aligned}$$

and here are the definitions of $\text{e}\dot{\text{v}}\text{al}$ and $\dot{\mu}$:

$$\text{e}\dot{\text{v}}\text{al}(t, e, K, w, A) = \begin{cases} w' & \text{if } \langle t \mid e \mid K \mid w \rangle \dot{\rightarrow}^* \langle \epsilon \mid r \mid \text{halt} \mid w' \rangle, \text{ where } r \in A \\ 0 & \text{otherwise} \end{cases}$$

$$\dot{\mu}(e, K, A) = \int \text{e}\dot{\text{v}}\text{al}(t, e, K, 1, A) \mu_{\mathbb{T}}(dt)$$

Note that an evaluation counts only if it completely exhausts its entropy sequence t . The approximants $\text{e}\dot{\text{v}}\text{al}^{(n)}$ and $\dot{\mu}^{(n)}$ are defined as before; in particular, they are indexed by number of steps, not by random numbers consumed.

²The dots are intended as a mnemonic for sequencing.

In general, the entropy access pattern is so different between split and sequenced entropy models that there is no correspondence between individual evaluations, and yet the resulting measures are equivalent.

LEMMA 7.1. *If $\langle t \mid e \mid K \mid w \rangle \rightsquigarrow \langle t' \mid e' \mid K' \mid w' \rangle$, then $\text{e\ddot{v}a}l^{(p+1)}(t, e, K, w, A) = \text{e\ddot{v}a}l^{(p)}(t', e', K', w', A)$.*

PROOF. By the definition of $\text{e\ddot{v}a}l^{(p+1)}$. \square

LEMMA 7.2. *The equations of Lemma 2.15 also hold for $\ddot{\mu}^{(n)}$ and $\ddot{\mu}$.*

PROOF. By definition of $\ddot{\mu}$ and Lemma 7.1. In fact, in contrast to Lemma 2.15, most of the cases are utterly straightforward, because no entropy shuffling is necessary. The sample case is different, because it relies on the structure of the entropy space:

$$\begin{aligned} \ddot{\mu}(e, K, A) &= \int \text{e\ddot{v}a}l(t, \text{sample}, K, 1, A) d(t) \\ &= \text{e\ddot{v}a}l(\epsilon, \text{sample}, K, 1, A) + \iint \text{e\ddot{v}a}l(r::t, \text{sample}, K, 1, A) \mu_{\mathbb{T}}(dt) \lambda(dr) \\ &= 0 + \iint I_{[0,1]}(r) \times \text{e\ddot{v}a}l(t, c_r, K, 1, A) \mu_{\mathbb{T}}(dt) \lambda(dr) \\ &= \int_0^1 \ddot{\mu}(c_r, K, A) \lambda(dr) \end{aligned}$$

\square

THEOREM 7.3 ($\ddot{\mu} = \mu$). *For all e, K , and $A \in \Sigma_{\mathbb{R}}$, $\ddot{\mu}(e, K, A) = \mu(e, K, A)$.*

PROOF. We first show $\ddot{\mu}^{(n)} = \mu^{(n)}$ by induction on n . The base case is $\ddot{\mu}^{(0)}(e, K, A) = \mu^{(0)}(e, K, A)$. There are two subcases: if $e = r$ and $K = \text{halt}$, then both results are $I_A(r)$. Otherwise, both measures are 0. Lemma 7.2 handles the inductive case. Finally, since the approximants are pointwise equivalent, their limits are equivalent. \square

7.2 Distributions

The language of Borgström et al. [2016] supports multiple real-valued distributions with real parameters; sampling from a distribution, in addition to consuming a random number, multiplies the current execution weight by the *density* of the distribution at that point. In this section we show that `sample` is equally expressive, given the inverse-CDF operations.

For each real-valued distribution of interest with n real-valued parameters, we add the following to the language: a sampling form $D(v_1, \dots, v_n)$ and operations $D\text{pdf}$, $D\text{cdf}$, and $D\text{invcdf}$ representing the distribution's density function, cumulative distribution function, and inverse cumulative distribution function, respectively. The operations take $n + 1$ arguments; by convention we write a semicolon before the parameters. For example, $\text{gammapdf}(x; k, s)$ represents the density at x of the gamma distribution with shape k and scale s .

We define the semantics of D using the sequenced-entropy framework by extending \rightsquigarrow with the following rule schema:

$$\langle r::t \mid D(r_1, \dots, r_n) \mid K \mid w \rangle \rightsquigarrow \langle t \mid r \mid K \mid w \times w' \rangle \quad \text{where } w' = D\text{pdf}(r; r_1, \dots, r_n) > 0$$

THEOREM 7.4. *$D(v_1, \dots, v_n)$ and $D\text{invcdf}(\text{sample}; v_1, \dots, v_n)$ are CIU-equivalent (and thus contextually equivalent).*

PROOF. By Theorem 6.15. Both expressions are real-valued. We must show that their real measures are equal. We abbreviate the parameters as \vec{v} . The result follows from the relationship between the density function and the cumulative density function.

$$\hat{\mu}(D\text{invcdf}(\text{sample}; \vec{v}), A) = \int_0^1 I_A(D\text{invcdf}(x; \vec{v})) dx$$

We change the variable of integration with $x = \text{Dcdf}(t; \vec{v})$ and $\frac{dx}{dt} = \text{Dpdf}(t; \vec{v})$:

$$\begin{aligned} &= \int_{-\infty}^{\infty} I_A(\text{Dinvcdf}(\text{Dcdf}(t; \vec{v}); \vec{v})) \times \text{Dpdf}(t; \vec{v}) dt \\ &= \int_{-\infty}^{\infty} I_A(t) \times \text{Dpdf}(t; \vec{v}) dt \\ &= \hat{\mu}(\text{D}(\vec{v}), A) \end{aligned}$$

□

7.3 From let-Style to Direct-Style

Let us call the language of Section 2.1 \mathcal{L} (for “let”) and the direct-style analogue \mathcal{D} (for “direct”). Once again following [Borgström et al. \[2016\]](#), we give the semantics of \mathcal{D} using a CS-style abstract machine, in contrast to the CSK-style machines we have used until now [\[Felleisen et al. 2009\]](#).

Here are the definitions of expressions and evaluation contexts for \mathcal{D} :

$$\begin{aligned} e ::= & v \mid \text{let } x = e \text{ in } e \mid (e e) \mid \text{op}(e, \dots, e) \mid \text{if } e \text{ then } e \text{ else } e \\ E ::= & [] \mid \text{let } x = E \text{ in } e \mid (E e) \mid (v E) \mid \text{op}(v, \dots, E, e, \dots) \mid \text{if } E \text{ then } e \text{ else } e \end{aligned}$$

Here are some representative rules for its abstract machine:

$$\begin{aligned} \langle t \mid E[(\lambda x. e) v] \mid w \rangle &\rightarrow_{\mathcal{D}} \langle t \mid E[e[v/x]] \mid w \rangle \\ \langle r :: t \mid E[\text{sample}] \mid w \rangle &\rightarrow_{\mathcal{D}} \langle t \mid E[c_r] \mid w \rangle \end{aligned}$$

And here are the corresponding definitions of evaluation and measure:

$$\begin{aligned} \text{eval}_{\mathcal{D}}(t, e, w, A) &= \begin{cases} w' & \text{if } \langle t \mid e \mid w \rangle \rightarrow_{\mathcal{D}}^* \langle \epsilon \mid r \mid w' \rangle, \text{ where } r \in A \\ 0 & \text{otherwise} \end{cases} \\ \mu_{\mathcal{D}}(e, A) &= \int \text{eval}_{\mathcal{D}}(t, e, 1, A) \mu_{\mathbb{T}}(dt) \end{aligned}$$

To show that our \mathcal{L} corresponds with \mathcal{D} , we define a translation $\text{tr}[-] : \mathcal{D} \rightarrow \mathcal{L}$. More precisely, $\text{tr}[-]$ translates \mathcal{D} -expressions to \mathcal{L} -expressions, such as

$$\begin{aligned} \text{tr}[[r]] &= r \\ \text{tr}[[\lambda x. e]] &= \lambda x. \text{tr}[[e]] \\ \text{tr}[[e_1 e_2]] &= \text{let } x_1 = \text{tr}[[e_1]] \text{ in let } x_2 = \text{tr}[[e_2]] \text{ in } (x_1 x_2) \\ \text{tr}[[\text{let } x = e_1 \text{ in } e_2]] &= \text{let } x = \text{tr}[[e_1]] \text{ in } \text{tr}[[e_2]] \end{aligned}$$

and it translates \mathcal{D} -evaluation contexts to \mathcal{L} -continuations, such as

$$\begin{aligned} \text{tr}[[[]]] &= \text{halt} \\ \text{tr}[[E[([] e_2)]]] &= (x_1 \rightarrow \text{let } x_2 = e_2 \text{ in } (x_1 x_2)) \text{tr}[[E]] \\ \text{tr}[[E[(v_1 [])]]] &= (x_2 \rightarrow (v_1 x_2)) \text{tr}[[E]] \\ \text{tr}[[E[\text{let } x = [] \text{ in } e]]] &= (x \rightarrow e) \text{tr}[[E]] \end{aligned}$$

Now we demonstrate the correspondence of evaluation and then lift it to measures.

LEMMA 7.5 (SIMULATION). $\text{eval}_{\mathcal{D}}(t, E[e], w, A) = \text{eval}(t, \text{tr}[[e]], \text{tr}[[E]], w, A)$

PROOF. From the CS-style machine above we can derive a corresponding CSK machine (call it $\rightarrow_{\mathcal{D}_{\text{CSK}}}$); the technique is standard [\[Felleisen et al. 2009\]](#). Then it is straightforward to show that

$$\langle t \mid e \mid K \mid w \rangle \rightarrow_{\mathcal{D}_{\text{CSK}}} \langle t' \mid e' \mid K' \mid w \rangle \implies \langle t \mid \text{tr}[[e]] \mid \text{tr}[[K]] \mid w \rangle \xrightarrow{*} \langle t' \mid \text{tr}[[e']] \mid \text{tr}[[K']] \mid w' \rangle$$

and thus the evaluators agree. □

THEOREM 7.6. $\mu_{\mathcal{D}}(E[e], A) = \mu(\text{tr}[\![e]\!], \text{tr}[\![E]\!], A)$

PROOF. By definition of $\mu_{\mathcal{D}}$ and Lemmas 7.3 and 7.5. □

The equational theory for \mathcal{D} is the *pullback* of the \mathcal{L} equational theory over $\text{tr}[\![_]\!]$. Compare with Sabry and Felleisen [1993], which explores the pullback of $\lambda\beta\eta$ and related calculi over the call-by-value CPS transformation. For our language \mathcal{D} , associativity and commutativity combine to yield a generalization of their β_{flat} and β'_{Ω} equations to “single-evaluation” contexts S :

$$S ::= [\] \mid (S e) \mid (e S) \mid \text{let } x = S \text{ in } e \mid \text{let } x = e \text{ in } S \\ \mid \text{op}(e, \dots, S, e, \dots) \mid \text{if } S \text{ then } e \text{ else } e$$

$$\text{let } x = e \text{ in } S[x] =_{\text{ctx}} S[e] \quad \text{when } x \notin FV(S) \cup FV(e) \quad (\text{let}_S)$$

8 INFORMALLY RELATED WORK

Our language and semantics are based on that of Culpepper and Cobb [2017], but unlike that language, which is simply-typed, ours is untyped and thus has recursion and nonterminating programs. Consequently, our logical relation must use step-indexing rather than type-indexing. Using an untyped language instead of a typed one not only introduces recursion; it increases the universe of expressions to which the theory applies, but it also makes the equivalence stricter, since the untyped language admits both more expressions and more contexts.

The construction of our logical relation follows the tutorial of Pitts [2010] on the construction of biorthogonal, step-indexed [Ahmed 2006] logical relations. Instead of termination, we use the program measure as the observable behavior, following Culpepper and Cobb [2017]. But unlike that work, where the meaning of an expression is a measure over arbitrary syntactic values, we define the meaning of an expression and continuation together (representing a whole program) as a measure over the reals. This allows us to avoid the complication of defining a relation on measurable sets of syntactic values [Culpepper and Cobb 2017, the \mathcal{A} relation].

There has been previous work on contextual equivalence for probabilistic languages with only *discrete* random variables. In particular, Bizjak and Birkedal [2015] define a step-indexed, biorthogonal logical relation whose structure is similar to ours, except that they sum where we integrate, and they use the probability of termination as the basic observation whereas we compare measures. Others have applied bisimulation techniques [Crubillé and Lago 2014; Sangiorgi and Vignudelli 2016] to languages with discrete choice; Ehrhard et al. [2014] have constructed fully abstract models for PCF with discrete probabilistic choice using probabilistic coherence spaces.

Staton et al. [2016] gives a denotational semantics for a higher-order, typed language with continuous random variables, scoring, and normalization but without recursion. Using a variant of that denotational semantics, Staton [2017] proves the soundness of the let-reordering transformation for a first-order language.

ACKNOWLEDGMENTS

This material is based upon work sponsored by the Air Force Research Laboratory (AFRL) and the Defense Advanced Research Projects Agency (DARPA) under Contract No. FA8750-14-C-0002. The views expressed are those of the authors and do not reflect the official policy or position of the Department of Defense or the U.S. Government.

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 695412).

REFERENCES

- Amal J. Ahmed. 2006. Step-Indexed Syntactic Logical Relations for Recursive and Quantified Types. In *Proc. 15th European Symposium on Programming (ESOP '06)*. 69–83. https://doi.org/10.1007/11693024_6
- Ales Bizjak and Lars Birkedal. 2015. Step-Indexed Logical Relations for Probability. In *Proc. 18th International Conference on Foundations of Software Science and Computation Structures (FoSSaCS '15)*. 279–294. https://doi.org/10.1007/978-3-662-46678-0_18
- Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szyczak. 2017. A Lambda-calculus Foundation for Universal Probabilistic Programming (long version). <https://arxiv.org/abs/1512.08990>
- Johannes Borgström, Ugo Dal Lago, Andrew D. Gordon, and Marcin Szyczak. 2016. A lambda-calculus foundation for universal probabilistic programming. In *Proc. 21st ACM SIGPLAN International Conference on Functional Programming (ICFP '16)*. 33–46. <https://doi.org/10.1145/2951913.2951942>
- Raphaëlle Crubillé and Ugo Dal Lago. 2014. On Probabilistic Applicative Bisimulation and Call-by-Value λ -Calculi. In *Proc. 23rd European Symposium on Programming (ESOP '14)*. 209–228. https://doi.org/10.1007/978-3-642-54833-8_12
- Ryan Culpepper and Andrew Cobb. 2017. Contextual Equivalence for Probabilistic Programs with Continuous Random Variables and Scoring. In *Proc. 26th European Symposium on Programming (ESOP '17)*. 368–392. https://doi.org/10.1007/978-3-662-54434-1_14
- Thomas Ehrhard, Christine Tasson, and Michele Pagani. 2014. Probabilistic coherence spaces are fully abstract for probabilistic PCF. In *Proc. 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '14)*. 309–320. <https://doi.org/10.1145/2535838.2535865>
- Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. 2009. *Semantics Engineering with PLT Redex* (1st ed.). The MIT Press.
- Noah D. Goodman, Vikash K. Mansinghka, Daniel M. Roy, Keith Bonawitz, and Joshua B. Tenenbaum. 2008. Church: a language for generative models. In *Proc. 24th Conference in Uncertainty in Artificial Intelligence (UAI '08)*. 220–229.
- Daniel Huang and Greg Morrisett. 2016. An Application of Computable Distributions to the Semantics of Probabilistic Programming Languages. In *Proc. 25th European Symposium on Programming (ESOP '16)*. 337–363. https://doi.org/10.1007/978-3-662-49498-1_14
- Oleg Kiselyov and Chung-chieh Shan. 2009. Embedded Probabilistic Programming. In *Proc. IFIP TC 2 Working Conference on Domain-Specific Languages (DSL '09)*. 360–384. https://doi.org/10.1007/978-3-642-03034-5_17
- Vikash Mansinghka, Daniel Selsam, and Yura Perov. 2014. Venture: a higher-order probabilistic programming platform with programmable inference. <http://arxiv.org/abs/1404.0099>
- Ian Mason and Carolyn Talcott. 1991. Equivalence in functional languages with effects. *J. Funct. Program.* 1, 3 (1991), 287–327. <https://doi.org/10.1017/S095679680000125>
- Kevin P. Murphy. 2007. Conjugate Bayesian analysis of the Gaussian distribution. <https://www.cs.ubc.ca/~murphyk/Papers/bayesGauss.pdf>
- Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. 2016. Probabilistic Inference by Program Transformation in Haku (System Description). In *Proc. 13th International Symposium on Functional and Logic Programming (FLOPS '16)*. 62–79. https://doi.org/10.1007/978-3-319-29604-3_5
- Brooks Paige and Frank Wood. 2014. A Compilation Target for Probabilistic Programming Languages. In *Proc. 31th International Conference on Machine Learning (ICML '14)*. 1935–1943.
- Sungwoo Park, Frank Pfenning, and Sebastian Thrun. 2008. A Probabilistic Language Based on Sampling Functions. *ACM Trans. Program. Lang. Syst.* 31, 1, Article 4 (Dec. 2008), 4:1–4:46 pages. <https://doi.org/10.1145/1452044.1452048>
- Andrew M. Pitts. 2010. Step-Indexed Biorthogonality: a Tutorial Example. In *Modelling, Controlling and Reasoning About State (Dagstuhl Seminar Proceedings)*, Amal Ahmed, Nick Benton, Lars Birkedal, and Martin Hofmann (Eds.). <http://drops.dagstuhl.de/opus/volltexte/2010/2806/>
- Amr Sabry and Matthias Felleisen. 1993. Reasoning about programs in continuation-passing style. *LISP and Symbolic Computation* 6, 3 (01 Nov 1993), 289–360.
- Davide Sangiorgi and Valeria Vignudelli. 2016. Environmental bisimulations for probabilistic higher-order languages. In *Proc. 43rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '16)*. 595–607. <https://doi.org/10.1145/2837614.2837651>
- Sam Staton. 2017. Commutative Semantics for Probabilistic Programming. In *Proc. 26th European Symposium on Programming (ESOP '17)*. 855–879. https://doi.org/10.1007/978-3-662-54434-1_32
- Sam Staton, Hongseok Yang, Frank Wood, Chris Heunen, and Ohad Kammar. 2016. Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints. In *Proc. 31st IEEE Symposium on Logic in Computer Science (LICS '16)*. 525–534. <https://doi.org/10.1145/2933575.2935313>
- Mitchell Wand, Ryan Culpepper, Theophilos Giannakopoulos, and Andrew Cobb. 2018. Contextual Equivalence for a Probabilistic Language with Continuous Random Variables and Recursion. <https://arxiv.org/abs/1807.02809>

Frank Wood, Jan-Willem van de Meent, and Vikash Mansinghka. 2014. A New Approach to Probabilistic Programming Inference. In *Proc. 17th International Conference on Artificial Intelligence and Statistics (AISTATS '14)*. 1024–1032.