# Statistical properties and evaluation of random $\lambda$-terms – a survey

Maciej Bendkowski[1]

[1]Theoretical Computer Science Departament
Jagiellonian University in Kraków

EUTypes meeting
Kraków 23-24 February 2019

# Why consider random $\lambda$-terms?

▶ Random generation of structures containing *scoped variables* with applications to software testing. For instance, the following *random* $\lambda$-term[1] exhibited a bug in GHC's strictness analizer:

$$(\lambda a.\, \mathsf{seq}\, a\, (\mathsf{seq}\, a\, \bot)\, \mathsf{tail})(\lambda a.\, \mathsf{seq}\, \bot\, (+1))$$

were

$$\bot\ \mathsf{seq}\ a = \bot$$
$$a\ \mathsf{seq}\ b = b.$$

---

[1] M. Pałka. *Random Structured Test Data Generation for Black-Box Testing.* PhD thesis. 2014.

# Why consider random λ-terms?

▶ Random generation of structures containing *scoped variables* with applications to software testing. For instance, the following *random* λ-term[1] exhibited a bug in GHC's strictness analizer:

$$(\lambda a.\, \mathsf{seq}\, a\, (\mathsf{seq}\, a\, \bot)\, \mathsf{tail})(\lambda a.\, \mathsf{seq}\, \bot\, (+1))$$

were

$$\bot\, \mathsf{seq}\, a = \bot$$
$$a\, \mathsf{seq}\, b = b.$$

▶ Such structures are difficult to analyse and inspire the development of new combinatorial techniques.

[1] M. Pałka. *Random Structured Test Data Generation for Black-Box Testing*. PhD thesis. 2014.

# How to represent $\lambda$-terms?

Let us focus on $\lambda$-terms with de Bruijn indices.

$T ::= \lambda T \mid T\,T \mid N$
$N ::= 0 \mid S\,N.$

$$\lambda xyz.xz(yz) \equiv \lambda\lambda\lambda\underline{2}\underline{0}(\underline{1}\underline{0})$$

Benefits:

---

# How to represent $\lambda$-terms?

Let us focus on $\lambda$-terms with de Bruijn indices.

$T ::= \lambda T \mid T\,T \mid N$
$N ::= 0 \mid S\,N.$

$$\lambda xyz.xz(yz) \equiv \lambda\lambda\lambda\underline{2}0(\underline{1}0)$$

Benefits:

▶ We do not have to worry about $\alpha$-equivalence.

---

[2]Requires some non-trivial extensions of well-known tools.

# How to represent $\lambda$-terms?

Let us focus on $\lambda$-terms with de Bruijn indices.

$T ::= \lambda T \mid T\, T \mid N$
$N ::= 0 \mid S\, N.$

$$\lambda xyz.xz(yz) \equiv \lambda\lambda\lambda\underline{2}\underline{0}(\underline{1}\underline{0})$$

Benefits:

▶ We do not have to worry about $\alpha$-equivalence.

▶ It is perhaps the simplest model we can start to analyse.

---

[2] Requires some non-trivial extensions of well-known tools.

# How to represent $\lambda$-terms?

Let us focus on $\lambda$-terms with de Bruijn indices.

$T ::= \lambda T \mid T\,T \mid N$
$N ::= 0 \mid S\,N.$

$$\lambda xyz.xz(yz) \equiv \lambda\lambda\lambda\underline{2}\,\underline{0}(\underline{1}\,\underline{0})$$
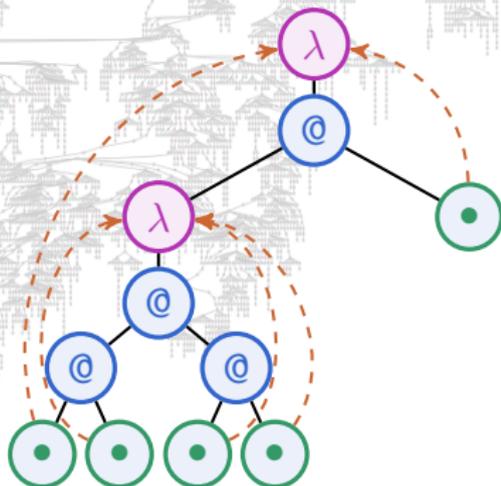
Benefits:

- ▶ We do not have to worry about $\alpha$-equivalence.
- ▶ It is perhaps the simplest model we can start to analyse.
- ▶ (Some) techniques of *analytic combinatorics* are applicable[2].

---

[2]Requires some non-trivial extensions of well-known tools.

# Alternative representations?

What about a representation were only *closed* terms are allowed?
Also, how should we measure the term size?

Assume that abstractions, applications and variables contribute
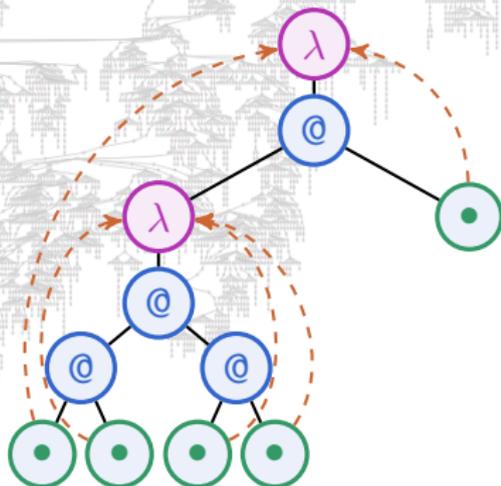some *weight* to the term size.



Variable weight options:

▶ No weight.

# Alternative representations?

What about a representation were only *closed* terms are allowed?
Also, how should we measure the term size?

Assume that abstractions, applications and variables contribute
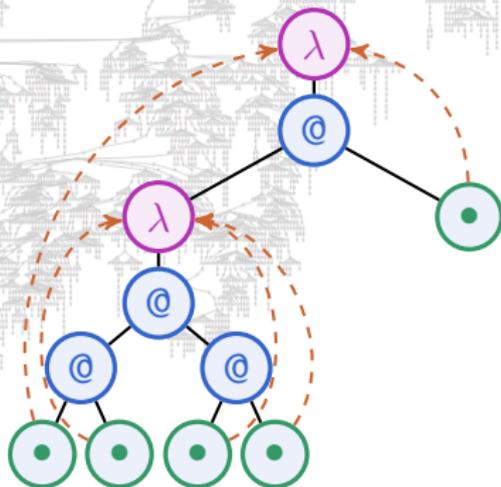some *weight* to the term size.



Variable weight options:
- ▶ No weight.
- ▶ Constant weight.

# Alternative representations?

What about a representation were only *closed* terms are allowed?
Also, how should we measure the term size?

Assume that abstractions, applications and variables contribute
some *weight* to the term size.



Variable weight options:

▶ No weight.

▶ Constant weight.

▶ Weight proportional
to the binder distance
(cf. de Bruijn indices).

# Counting $\lambda$-terms – how hard can it be?

Let's start with the following *counting problem*. Given $n$, what is the number $\Lambda_n$ of $\lambda$-terms of size $n$? How does $\Lambda_n$ change with $n \to \infty$?

# Counting $\lambda$-terms – how hard can it be?

Let's start with the following *counting problem*. Given $n$, what is the number $\Lambda_n$ of $\lambda$-terms of size $n$? How does $\Lambda_n$ change with $n \to \infty$?

Bodini, Gardy, Gittenberger, Jacquot '13

Assume that variables contribute constant weight one. For special cases, such as linear or affine terms, the counting problem admits asymptotic estimates. In general, the problem remains <u>unsolved</u>.

# Counting $\lambda$-terms – how hard can it be?

Let's start with the following *counting problem*. Given $n$, what is the number $\Lambda_n$ of $\lambda$-terms of size $n$? How does $\Lambda_n$ change with $n \to \infty$?

## Bodini, Gardy, Gittenberger, Jacquot '13

Assume that variables contribute constant weight one. For special cases, such as linear or affine terms, the counting problem admits asymptotic estimates. In general, the problem remains <u>unsolved</u>.

## David, Grygiel, Kozik, Raffalli, Theyssier, Zaionc '13

Assume that variables contribute no weight. The counting problem is still unsolved, however asymptotically almost all $\lambda$-terms are shown to be <u>strongly normalizing</u>.

# Back to the de Bruijn representation

What if we assume that $\underline{n} \equiv S^{(n)}0$ has *natural* weight $n + 1$?
Recall that in the de Bruijn representation we have

$T ::= \lambda T \mid T\,T \mid N$
$N ::= 0 \mid S\,N.$

# Back to the de Bruijn representation

What if we assume that $\underline{n} \equiv S^{(n)}0$ has *natural* weight $n + 1$?
Recall that in the de Bruijn representation we have

$T ::= \lambda T \mid T\,T \mid N$
$N ::= 0 \mid S\,N.$

### B., Grygiel, Lescanne, Zaionc '16

In the de Bruijn representation, the counting problem admits asymptotic estimates. More importantly, asymptotically almost all $\lambda$-terms are not strongly normalising.

# What about closed $\lambda$-terms?

Let $T_m$ denote the class of *m-open $\lambda$-terms*, i.e. terms which after prepending $m$ head abstractions become closed. Then, the set $T_0$ of closed terms satisfies the following *infinite* specification:

$$\begin{cases} T_0 ::= \lambda T_1 \mid T_0 T_0 \\ T_1 ::= \lambda T_2 \mid T_1 T_1 \mid \underline{0} \\ T_2 ::= \lambda T_3 \mid T_2 T_2 \mid \underline{0} \mid \underline{1} \\ \qquad \cdots \\ T_m ::= \lambda T_{m+1} \mid T_m T_m \mid \underline{0} \mid \underline{1} \mid \cdots \mid \underline{m-1} \\ \qquad \cdots \end{cases}$$

### Bodini, Gittenberger, Gołębiewski '18

The counting problem for closed $\lambda$-terms admits asymptotic estimates. Moreover, it is possible to *sample* (uniformly) random, closed $\lambda$-terms of large, target size *n*.

# Sampling closed, simply-typed $\lambda$-terms?

For some applications, sampling closed terms can be insufficient. Sometimes, we need stronger guarantees about the *properties* of generated terms, such as *termination*.

## B., Grygiel, Tarau '17

Combining samplers for closed $\lambda$-terms and a careful use of rejection, it is possible design practical samplers for closed, simply-typed $\lambda$-terms for *moderate* target sizes.

Benchmark term sizes:[3]

---

[3] Using a standard PC as a reference point.

# Sampling closed, simply-typed $\lambda$-terms?

For some applications, sampling closed terms can be insufficient. Sometimes, we need stronger guarantees about the *properties* of generated terms, such as *termination*.

### B., Grygiel, Tarau '17

Combining samplers for closed $\lambda$-terms and a careful use of rejection, it is possible design practical samplers for closed, simply-typed $\lambda$-terms for *moderate* target sizes.

Benchmark term sizes:[3]

▶ Closed $\lambda$-terms: $n \leq 1\,000\,000\,000$.

---

[3]Using a standard PC as a reference point.

# Sampling closed, simply-typed λ-terms?

For some applications, sampling closed terms can be insufficient. Sometimes, we need stronger guarantees about the *properties* of generated terms, such as *termination*.

## B., Grygiel, Tarau '17

Combining samplers for closed λ-terms and a careful use of rejection, it is possible design practical samplers for closed, simply-typed λ-terms for *moderate* target sizes.

Benchmark term sizes:[3]

- ▶ Closed λ-terms: $n \leq 1\,000\,000\,000$.
- ▶ Closed, simply-typed λ-terms: $n \leq 250$.

---

[3]Using a standard PC as a reference point.

# Typical properties of $\lambda$-terms

### B., Bodini, Dovgal '18

Suppose that we sample a large, random (unrestricted) $\lambda$-term of size *n*. What shape and properties should we expect?

| Parameter | Mean, $\sim$ | | Distribution | |
|---|---|---|---|---|
| | plain | closed | plain | closed |
| Variables | $0.307n$ | | Normal | |
| Abstractions | $0.258n$ | | Normal | |
| Successors | $0.129n$ | | Normal | |
| Redexes | $0.091n$ | | Normal | |
| Index value | $0.420$ | | Geometric | |
| Head abstractions | 0.420 | 1.447 | Geometric | Discrete |
| *m*-openness | 2.019 | 0 | Discrete | trivial |
| Free variables | 5.722 | 0 | Discrete | trivial |

# Even better control over generated $\lambda$-terms?

## B., Bodini, Dovgal '18

Using quite general tools[4], it is possible to *skew* the uniform distribution of generated $\lambda$-terms, and gain additional control over some of their parameters. For instance, request more abstractions or favour larger de Bruijn indices.



---
[4] See also https://github.com/maciej-bendkowski/boltzmann-brain

# Typical properties of λ-terms (II)

What about evaluation of large, random terms? For instance, how long does it take *on average* to find the leftmost-outermost redex in a random term?



It is always a constant time overhead:

- ▶ Plain terms ≈ 6.222,
- ▶ Closed terms ≈ 6.054.

What can be said about the typical cost of substitution?

# Substitution resolution and explicit substitutions

It is possible to carry out substitutions in various ways,
e.g. substitutions can be carried out strictly or suspended
until truely required (i.e. resolved non-strictly).

Let us therefore consider a *concrete implementation* of substitution
resolution, say $\lambda \upsilon^5$ — a $\lambda$-calculus with explicit substitutions — and
analyse substitution resolution therin, instead.

$$t ::= \underline{n} \mid \lambda t \mid tt \mid t[s]$$
$$s ::= t/ \mid \Uparrow (s) \mid \uparrow$$
$$\underline{n} ::= \underline{0} \mid S\underline{n}.$$

$$
\begin{array}{rll}
(\lambda a)b & \to a[b/] & \text{(Beta)} \\
(ab)[s] & \to a[s](b[s]) & \text{(App)} \\
(\lambda a)[s] & \to \lambda(a[\Uparrow (s)]) & \text{(Lambda)} \\
\underline{0}[a/] & \to a & \text{(FVar)} \\
(S\underline{n})[a/] & \to \underline{n} & \text{(RVar)} \\
\underline{0}[\Uparrow (s)] & \to \underline{0} & \text{(FVarLift)} \\
(S\underline{n})[\Uparrow (s)] & \to \underline{n}[s][\uparrow] & \text{(RVarLift)} \\
\underline{n}[\uparrow] & \to S\underline{n}. & \text{(VarShift)}
\end{array}
$$

---

[5] P. Lescanne. *From $\lambda\sigma$ to $\lambda\upsilon$ – a journey through calculi of explicit substitutions.* 1994.

# $\upsilon$-reduction grammars

Instead of a single $\beta$ rewriting rule, $\lambda\upsilon$ consists of a (Beta) rule and seven auxiliary $\upsilon$ rules governing the execution of substitutions. Luckily, compared with classic $\lambda$-calculus they are much easier to analyse in quantitative terms.
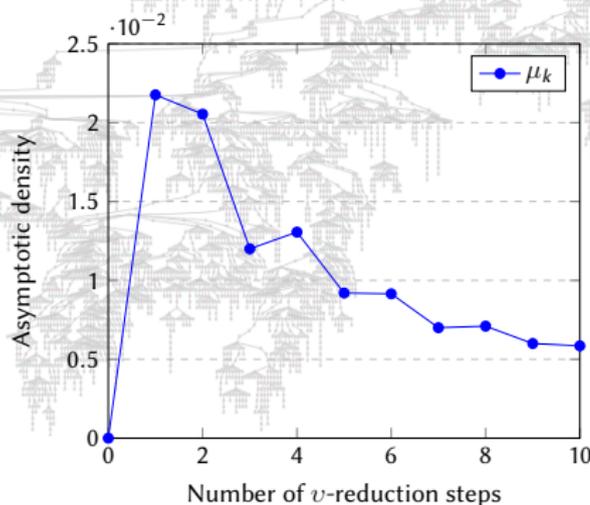
### B. '19

For all $k \geq 0$, the set $G_k$ of $\lambda\upsilon$-terms which reduce to their $\upsilon$-normal forms (i.e. pure forms without explicit substitutions) in $k$ leftmost-outermost $\upsilon$ rewriting steps forms a *regular tree language*.

$$G_0 \to \lambda G_0 \mid G_0 G_0 \mid \underline{n}$$
$$G_1 \to \lambda G_1 \mid G_0 G_1 \mid G_1 G_0$$
$$\mid \underline{0}[(G_0 G_0)/] \mid \underline{0}[\lambda G_0/] \mid \underline{0}[\underline{n}/]$$
$$\mid (S\underline{n})[t/] \mid \underline{0}[\Uparrow (s)] \mid \underline{n}[\uparrow]$$

$$\vdots$$

# $\upsilon$-reduction grammars (II)

$(G_k)_k$ admits a neat hierarchical structure and can be analysed using standard techniques of *analytic combinatorics*. In particular, for any fixed $k \geq 0$, the fraction of terms $\upsilon$-normalising in $k$ steps tends to a computable limit $\mu_k$ as the term size tends to infinity.



Number of $\upsilon$-reduction steps (x-axis), Asymptotic density (y-axis)

| $k$ | $\mu_k$ |
|-----|---------|
| 0   | 0.      |
| 1   | 0.02176 |
| 2   | 0.02054 |
| 3   | 0.01200 |
| 4   | 0.01306 |
| 5   | 0.00920 |
| 6   | 0.00915 |
| 7   | 0.00700 |
| 8   | 0.00710 |
| 9   | 0.00600 |
| 10  | 0.00585 |

# Conclusions

▶ Depending on the representation, typical $\lambda$-terms admit contrasting properties, e.g. strong normalisation or lack thereof.

# Conclusions

▶ Depending on the representation, typical $\lambda$-terms admit contrasting properties, e.g. strong normalisation or lack thereof.

▶ Representations using de Bruijn indices in unary notation are now well-understood and admit effective tools for random generation (at least in the untyped universe).

# Conclusions

▶ Depending on the representation, typical $\lambda$-terms admit contrasting properties, e.g. strong normalisation or lack thereof.

▶ Representations using de Bruijn indices in unary notation are now well-understood and admit effective tools for random generation (at least in the untyped universe).

▶ Although the techniques of analytic combinatoric are quite daunting, it is possible to use them in order to analyse the operational costs of substitution in $\lambda$-calculus.