

Kolmogorov complexity, prime numbers, and complexity lower bounds

LMS Computer Science Colloquium

November/2020

Igor Carboni Oliveira

University of Warwick

Research funded by a Royal Society University Research Fellowship

Overview

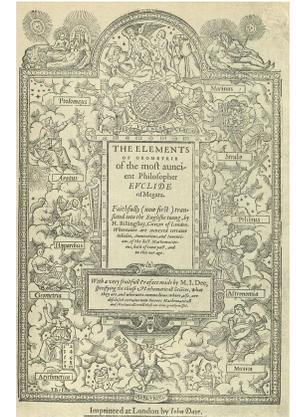
Maths 1. Are there infinitely many prime numbers with “simple” descriptions?

CS 2. Is it hard to detect patterns in data?



Quanta

Maths/CS 3. Is there a fast deterministic algorithm that, given n , outputs an n -bit prime?



Euclid's *Elements*

Wikipedia

This talk:

New insights using a **probabilistic** extension of (time-bounded) Kolmogorov complexity

Background and Motivation

1. Number Theory: Mersenne Primes

Primes of the form $M_n = 2^n - 1$.



51 Mersenne primes are known. The largest known prime is $2^{82589933} - 1$.

In **binary representation**, a Mersenne prime is of the form 11111.....11111.

Mersenne primes admit a short and effective representation.

“Simplest” possible representation of an n-bit prime.

Q. Are there infinitely many **Mersenne primes**?



Q. Are there infinitely many **primes of “minimum description length”**?

Time-bounded Kolmogorov Complexity

Mersenne primes admit a [short](#) and [effective](#) representation.

Kolmogorov Complexity?

001010100100100100100100100100100001001011001010101010101

Levin (1984) proposed the following notion of complexity for strings.

$$Kt(x) \stackrel{\text{def}}{=} \min_{\substack{\text{TM } M, \text{ time } t \\ M \text{ prints } x \text{ in time } t}} |M| + \log t$$

[short](#) [effective](#)



For every n -bit string x : $\log n \leq Kt(x) \leq n + O(\log n)$

111...111

2. Complexity Theory: Intractability

Problems about the **complexity of strings** play a significant role in theory of computing.

(e.g. *learning & cryptography*)

Undecidable

Kolmogorov complexity

Given x , estimate $C(x)$

Exponential Time vs Polynomial Time

Levin's Kt complexity

$$Kt(x) \stackrel{\text{def}}{=} \min_{\substack{\text{TM } M, \text{ time } t \\ M \text{ prints } x \text{ in time } t}} |M| + \log t$$

Given x , estimate $Kt(x)$ **Q.** *Is it in polynomial time?*

e.g. [ABKvMR'06]

$\log t$ function is a threshold

3. Algorithms: Deterministic constructions

POLYMATH 4

$$\underbrace{[100 \dots 000, 111 \dots 111]}_{n \text{ bits}}_2 = [2^{n-1}, 2^n - 1]_{10}$$

Challenge of deterministically generating primes: Given n , output an n -bit prime.

Best known deterministic algorithm runs in time $2^{n/2}$. [Lagarias-Odlyzko'87]

\implies For every large n , there is an n -bit prime p_n with $\text{Kt}(p_n) \leq \frac{n}{2} + o(n)$.



“Simple objects are easier to find”

If there is a sequence p_n of n -bit primes with $\text{Kt}(p_n) \leq \gamma_n$ then primes can be deterministically generated in time $\approx 2^{\gamma_n}$.

Summary

“Simplicity” as bounded Kt complexity (e.g. Mersenne primes).

Connections to basic questions in Maths/CS.

-  Are there n -bit primes of Kt complexity $o(n)$?
-  It is hard to estimate $Kt(x)$ of a given string x ?
-  Deterministic prime generation in time $2^{o(n)}$?

These remain longstanding problems relevant to number theory, algorithms, and complexity.



A theory of probabilistic representations

[O-Santhanam'17] Pseudodeterministic constructions in subexponential time.

[O'19] Randomness and intractability in Kolmogorov complexity.

[Lu-O'20] An efficient coding theorem via probabilistic representations and its applications.

Definition of rKt complexity

[O'19] A randomized analogue of Levin's Kt complexity:



$$\text{rKt}(x) \stackrel{\text{def}}{=} \min_{\substack{\text{randomized TM } M, \text{ time } t \\ \Pr_M[M \text{ prints } x \text{ in time } t] \geq 2/3}} |M| + \log t$$

A **short** and **effective** **probabilistic** procedure that is likely to generate the observed data.

Basic properties of rKt

$$\text{rKt}(x) \stackrel{\text{def}}{=} \min_{\substack{\text{randomized TM } M, \text{ time } t \\ \Pr_M[M \text{ prints } x \text{ in time } t] \geq 2/3}} |M| + \log t$$

$$\text{Kt}(x) \stackrel{\text{def}}{=} \min_{\substack{\text{TM } M, \text{ time } t \\ M \text{ prints } x \text{ in time } t}} |M| + \log t$$

For every string x , $\text{rKt}(x) \leq \text{Kt}(x)$.

Q. *Are there strings that admit a more succinct representation using randomness?*

[O'19] If $E \notin \text{i.o.SIZE}[2^{\epsilon n}]$, for every string x we have $\text{rKt}(x) = \Theta(\text{Kt}(x))$.

As far as we know, gap between rKt and Kt could be maximum.

*Proxy measure
to investigate Kt?*

A theory of probabilistic representations



[O-Santhanam'17] Pseudodeterministic constructions in subexponential time.

[O'19] Randomness and intractability in Kolmogorov complexity.

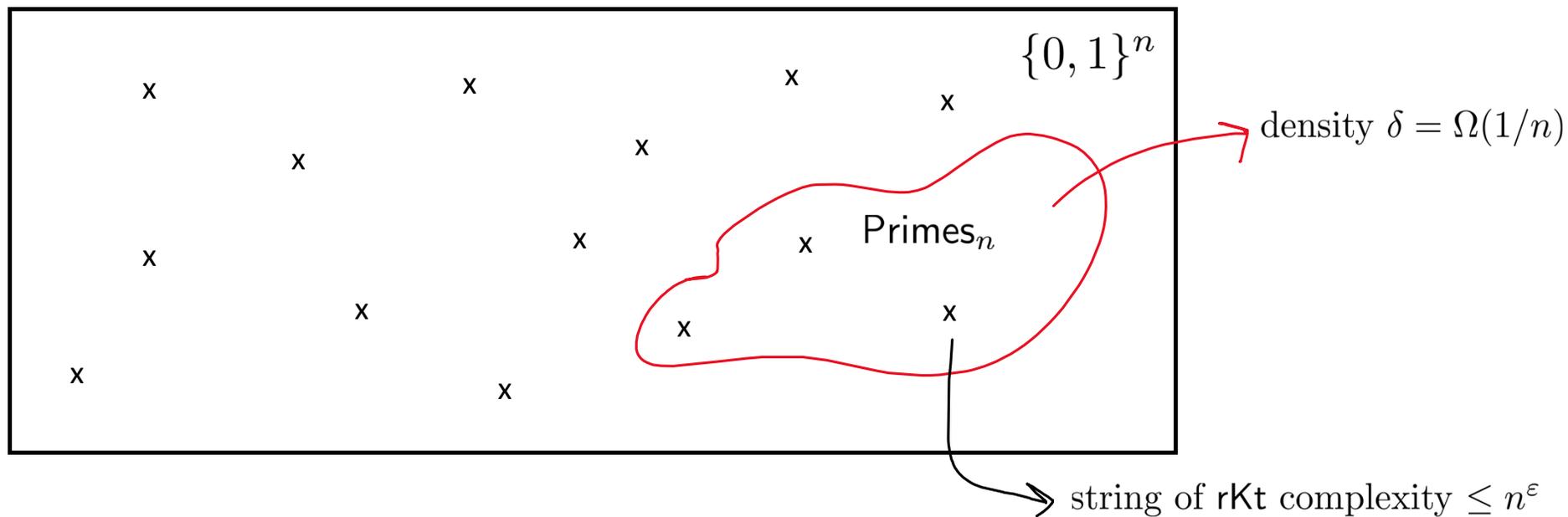
[Lu-O'20] An efficient coding theorem via probabilistic representations and its applications.

Probabilistic representations can see patterns in prime numbers

Recall: Open to show \exists primes of Kt complexity $< n/2$.

Theorem [O-Santhanam'17, O'19]. $\forall \varepsilon > 0$, for infinitely many values of n , \exists n -bit prime p_n such that $\text{rKt}(p_n) \leq n^\varepsilon$.

Informally, some primes are structured enough to admit “short” and “effective” probabilistic representations.

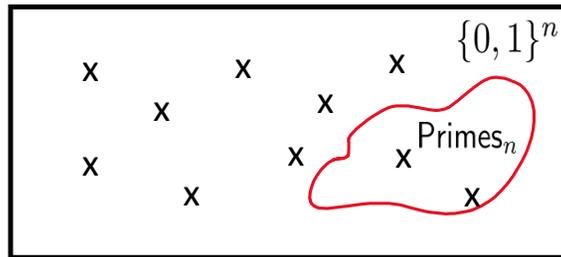


Intuition: A random n -bit string “hits” Primes_n with probability δ .

Construct a pseudorandom distribution \mathcal{D} supported over $\mathcal{R}_{\leq n^\epsilon}^{\text{rKt}}$.

└─→ “Tests”: $\text{Primes}_n \in \text{DTIME}[n^k]$.

Pseudorandomness



Example. $\mathcal{T} = \{x_1, \dots, x_n\}$

$G: \{0, 1\}^1 \rightarrow \{0, 1\}^n$

$G(0) = 000 \dots 000$

$G(1) = 111 \dots 111$

Fact: G ε -fools \mathcal{T} for $\varepsilon = 0$.

\mathcal{D} ε -fools a class of tests \mathcal{T} if:

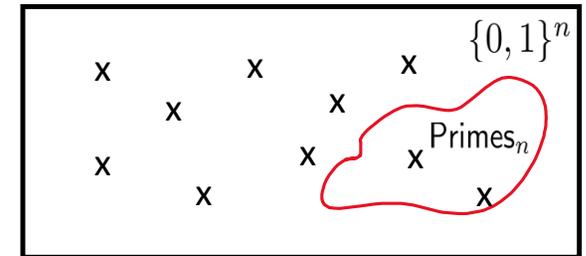
For every function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ in \mathcal{T} ,

$$\left| \Pr_{x \sim \{0, 1\}^n} [f(x) = 1] - \Pr_{y \sim \mathcal{D}} [f(y) = 1] \right| \leq \varepsilon.$$

$G: \{0, 1\}^s \rightarrow \{0, 1\}^n$, $s \ll n$, generates a distribution $\mathcal{D} \equiv G(\mathcal{U}_s)$.

$G: \{0, 1\}^s \rightarrow \{0, 1\}^n$, $s \ll n$, generates a distribution $\mathcal{D} \equiv G(\mathcal{U}_s)$.

“Fools” a class of tests \mathcal{T} .



Crucial: n -bit outputs of G have low rKt complexity: “explained” by a seed of length s .

Long list of works on PRGs in TCS:
 → **Unconditional** PRGs against “**weak**” tests.
 → **Conditional** PRGs against “**expressive**” tests.

[IW’97] $E \not\subseteq \text{i.o.SIZE}[2^{\lambda n}] \rightarrow \exists G: \{0, 1\}^{O(\log n)} \rightarrow \{0, 1\}^n$ that fools $\text{SIZE}[n^k]$.



(★)



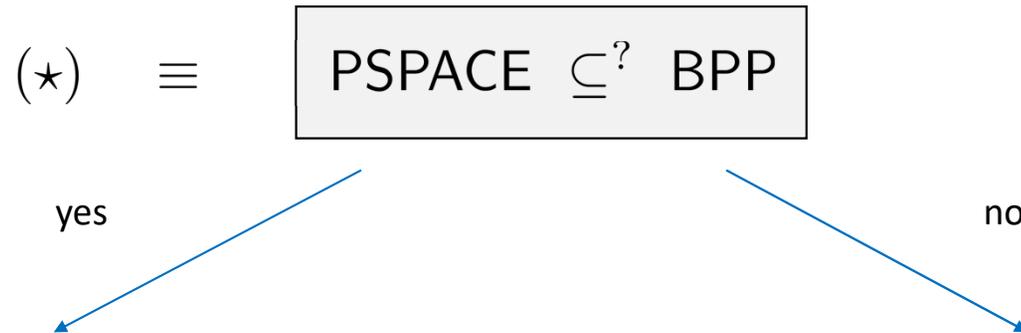
NOT (★) \implies “Easiness”

Time-bounded computations are more powerful.

Certain “patterns” in primes might become evident.

Example: The lexicographic first n -bit prime.

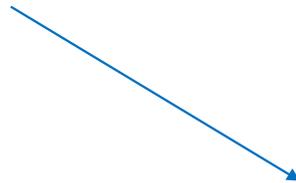
Need to find a “sweet spot” (★) that is useful.



Fact: PSPACE computations can detect first n-bit prime.



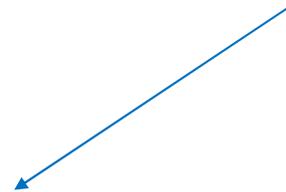
(probabilistic representation)



[IW98, TV07] PRG G under hardness assumption that fools $\text{DTIME}[n^k]$ infinitely often.

$\text{Support}(G) \subseteq \mathcal{R}_{n^\epsilon}^{\text{rKt}}$

$\text{Support}(G)$ intersects Primes_n



Infinitely many primes of bounded rKt complexity

A theory of probabilistic representations

[O-Santhanam'17] Pseudodeterministic constructions in subexponential time.



[O'19] Randomness and intractability in Kolmogorov complexity.

[Lu-O'20] An efficient coding theorem via probabilistic representations and its applications.

Is it hard to detect patterns?

$$\mathcal{R}_{\leq n^\varepsilon}^{\text{rKt}}$$

$$\mathcal{R}_{\geq .99n}^{\text{rKt}}$$

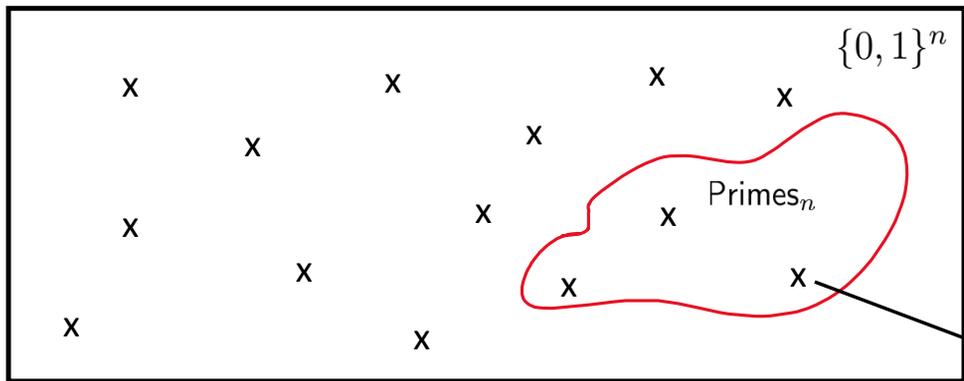
“structured”

“random”

Theorem [O’19]. $\forall \varepsilon > 0$, there is no randomised algorithm running in quasi-polynomial time that accepts strings in $\mathcal{R}_{\leq n^\varepsilon}^{\text{rKt}}$ and rejects strings in $\mathcal{R}_{\geq .99n}^{\text{rKt}}$

We cannot feasibly distinguish “*structured*” strings from “*random*” strings.

Proof of a weaker result: The set $\mathcal{R}_{>n^\epsilon}^{\text{rKt}}$ is not in P.



Info about Primes_n needed in previous proof:

- **Density:** $\delta \geq 1/\text{poly}(n)$
- **Easiness:** Primes_n \in P

low rKt

Lemma. Any **dense** and **easy** set contains, infinitely often, strings x with $\text{rKt}(x) \leq n^\epsilon$.

Proof of a weaker result: The set $\mathcal{R}_{>n^\epsilon}^{\text{rKt}}$ is not in P.

Lemma. Any **dense** and **easy** set contains, infinitely often, strings x with $\text{rKt}(x) \leq n^\epsilon$.

$\mathcal{R}_{>n^\epsilon}^{\text{rKt}}$ is **dense**.

If $\mathcal{R}_{>n^\epsilon}^{\text{rKt}}$ is also **easy** (in P), then we contradict the lemma.

A more delicate argument is used for the **gap problem** and against **BPTIME[quasi-poly]**.

$\mathcal{R}_{\leq n^\epsilon}^{\text{rKt}}$

$\mathcal{R}_{\geq .99n}^{\text{rKt}}$

“structured”

“random”

A theory of probabilistic representations

[O-Santhanam'17] Pseudodeterministic constructions in subexponential time.

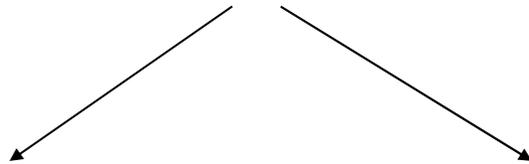
[O'19] Randomness and intractability in Kolmogorov complexity.



[Lu-O'20] An efficient coding theorem via probabilistic representations and its applications.

Perspective

rKt has enabled results that remain intriguing questions for Levin's **Kt** complexity.



existence of shorter representations

intractability of estimating rKt

Q. Can we further advance **time-bounded Kolmogorov complexity** using probabilistic representations?

Pillars of Kolmogorov Complexity

Three essential results in Kolmogorov complexity:

Language Compression Theorem

$$C(x) \leq |L^n| + O(\log n)$$

Hardness Assumption

Symmetry of Information

$$C(xy) \approx C(x) + C(y|x) \approx C(y) + C(x|y)$$

Hardness Assumption

Source Coding Theorem

$$C(x) \leq \log(1/\delta(x)) + O(1)$$

Time-bounded version?



See e.g. [Troy Lee, PhD Thesis]



Coding Theorem

Shannon's Information Theory

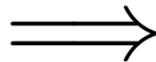
Distributions, entropy, compression, etc.

Coding Theorem in
Kolmogorov Complexity

Kolmogorov Complexity

Individual strings and their complexities.

An object x can be sampled
with probability δ



x admits a representation
of length $\approx \log(1/\delta)$

Interested in establishing an **unconditional time-bounded** version of the **Coding Theorem**.

An *Efficient* Coding Theorem for rKt

[Zhenjian Lu-O'20] “**Samplable** objects admit **short** and **effective** representations.”

Randomized Algorithm $A(1^m)$

Runs in time $T(m)$

Outputs string x with probability $\geq \delta$

$$\implies \text{rKt}(x) = O_A(\log(1/\delta) + \log(T) + \log m)$$

Efficient generation of representation: Given x , A , m , and δ , we can compute in time $\text{poly}(|x|, |A|, \log m, \log(1/\delta))$ and with probability $\geq .99$ a valid rKt representation.

└─ “**Magic**”: Running time has no dependence on T .

**Extremely useful
in applications!**

Application: Efficient universal compression

$x = 0010101000101111010010010111100010101010001010100101111010010100$

“There is a way of compressing it to k bits” (in the sense of **rKt**).

$$\text{rKt}(x) \leq k$$

\implies We can output in polynomial time with probability $\geq .99$ a valid **rKt** encoding of x of complexity $O(k)$.

Open Problem

*Is the existence of succinct representations for primes a **rare** phenomenon?*

Prove that for every large n , there is an n -bit prime of rKt complexity $\leq \varepsilon n$.

*By the **Coding Theorem for rKt**, enough to show that:*

Problem. Is there a probabilistic algorithm that, given n as input, runs in time $2^{\varepsilon n}$ and outputs some fixed n -bit prime with probability at least $1/2^{\varepsilon n}$?

This is a relaxation of the Polymath problem of deterministically generating primes.

Summary: Probabilistic Data Representations

$$K_t \longrightarrow rK_t$$



(under assumptions)

$$K_t \approx rK_t$$

Succinct Descriptions:

Infinitely many primes have rK_t complexity $\leq n^\epsilon$.

Computational Hardness:

It is intractable to estimate the rK_t complexity of an input string.

Coding Theorem:

Samplable objects admit short and effective representations.

Main References

[O-Santhanam'17] Pseudodeterministic constructions in subexponential time (STOC'2017).

[O'19] Randomness and intractability in Kolmogorov complexity (ICALP'19).

[Lu-O'20] An efficient coding theorem via probabilistic representations and its applications (preprint).

[Levin'84] Randomness conservation inequalities; information and independence in mathematical theories. *Information and Control*, 61(1):15-37, 1984.

[ABKvMR'06] Eric Allender, Harry Buhrman, Michal Koucky, Dieter van Melkebeek, and Detlef Ronneburger. Power from random strings. *SIAM J. Comput.*, 35(6):1467-1493, 2006.

[Polymath 4 - TCH12] Terence Tao, Ernest Croot, III, and Harald Helfgott. Deterministic methods to find primes. *Math. Comp.*, 81(278):1233-1246, 2012.

[Troy Lee, PhD Thesis] Kolmogorov complexity and formula lower bounds. University of Amsterdam, 2006.



<https://www.dcs.warwick.ac.uk/~igorcarb/complexity-meetings.html>