

Randomness classification

... fun guaranteed

Vladimir Matlach

Diego Krivochen

Jiri Milicka

Lukas Zamecnik

1000010100000010000111000000110110010010100110110010101100011100010001010111011
10111110001000100000000111111110001001010011010111111100101010111010101
100000010101001100010100111000110110010110110111100111111011010000001101100110
10011101101011010001111111101101100101000010111111000001111010110001100010011011
1010011011100111001011000100001100111100010111001110011101000110100110001010
1000101001110111011011110001011100100110011110101100111110110011100000011100
011000011110110001000100101111101010001001100000001101000100111001010001101
101100111001101111110001001001000110111011011011010000011101111011110111100
001101100101100011100111000010111010110011111011110001101100100100001010101
1101001000111000001011101101001010110001001101101100101101111000111111010110001
1111010010101110010011100100111001100100011011001001101010100111110100101011011
000110000111111010011001110000111110110010011101110110011001001001000101111101
000110110100100011011000111011111001101011100011010011001101011001010011010000
1010100111101100010010001001101110110100011110000001000000000100111110100001
111100000110100011001100101101110101111011111110101111011111110111011111101110011111
100100011110010011010111001010001110011001001110010010111110111010000101000001
0010000110000100011011110000000100010011110011010010001000001111100110110000110101
111111111001100001111101000101001011101011000101100111111101010100100011100101
0111101110101010100100101101
1000011101001000111111001110100110001100101001000111010000001110100000011011010001
10011100101110101110100010011001101010101011001000111110111101000100111101111001000101001110
101001101100100100110111110101000101011010111010000111010000010000000100111100111100111
00011111100100100011101010010001111011111010001100000100101011110011111001000

... what if i told you

Randomness

... now we know

knowing anything about sequences that...

should be random

– or –

look random

is important

Randomness

What is randomness?

... lack of patterns and unpredictability

... each token has the same probability to appear next in seq.

... no better prediction of next tokens than guess

... truly random string = no rules, no grammar

Randomness from linguist perspective

Non random strings have rules

... rules create patterns

... natural languages are not random strings

Linguistics has non-random sequences

... we naturally want to know the kind of rules behind it

It shows up that...

... quantification of some sequence properties might answer the questions

... beyond linguistics, various textual sequences exists that we really like to understand

the DNA

```
ACTTGTCA GTGCATATCGTAATCGTGTACGTATTGCAAGCGTTATCGGATATGC  
GATCGGTAGCTGCTATCGATCGCGGTATTATATATCGGTGACTGCGATGATGCG  
CGATTATCGGATCTAGCTAGCTGATTATCGATCAGTCTGATGCATGCCGCGCGT  
ATATCTAGCGTATTTTCGATCGGCGTATGCATGCTGCATCGATCACGCATTAG  
CATGTTGCCGGATTATCGATCGATGCTGCATGTATAAATCGATCGGATATT  
CGCGGCTATATTAGGCTATAGCTATGAGT...
```

... it is said that huge portion of DNA is just a random rubbish...

Voynich manuscript

... *very famous*

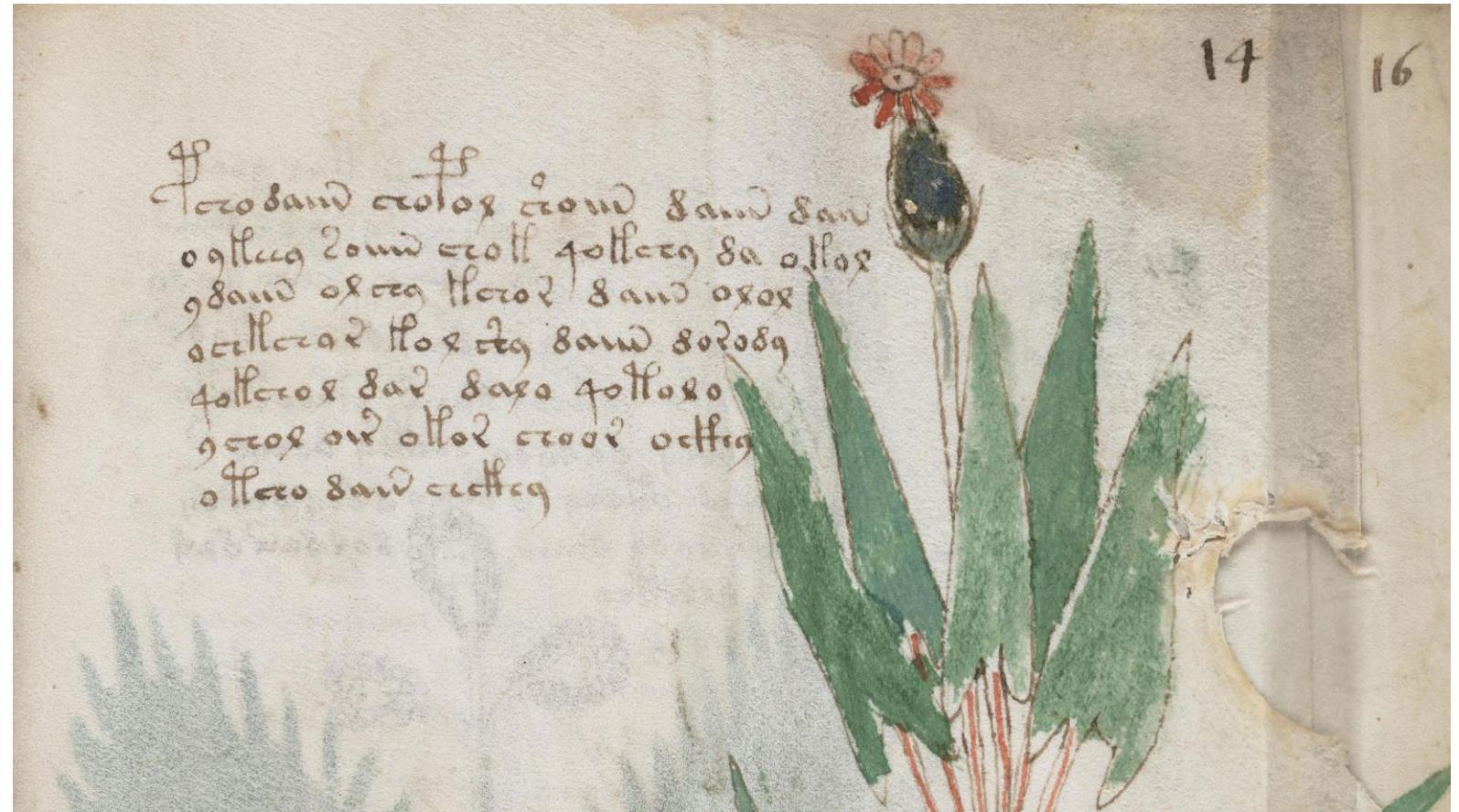
... *early 15th century* (carbon-dated)

... *a lot of \$\$\$ from Roman Emperor Rudolph II.*

... *unknown language*

... *might be encrypted*

... *might be a fraud*



So we live in a world of...

unkwnown sequences

... that we really want to understand better

And also in the world where...

we pay a lot of money for the *real true random*

... b/c of banks, computer security, casinos

What we want?

We want to take this... *(click)*

We want to take this...

```
10000101000000010000111000000110110010010100110110010101100011100010001010111011  
10111110001000100000000011111111000100101001101010111111100101010111010101  
10000000101010011000101001110001110110010110110111100111111011010000001101100110  
10011101101011010001111111101101100101000010111111000001111010110001100010011011  
10100110111001100011000100001100111100010111001011100111010001101001100001010  
1000101001110111011011110001011100100110011110101110011110110011100000011100  
011000011110110001000100101111010100010100110000000110100001001110010100001101  
1011001110011011111100010010011000110111011011010000011101111011110111100111100
```

And be able to answer

*... Are these sequences **random** & how much?*

... Encrypted?

... Monkey-typed?

... Natural Language look like?

... or some trivial repetitions?

And say:

„Yes, it is military hard encryption“

There sure exist

methods detecting *true random*

Like...

Diehard tests battery, Entropy, Chi-squared test, mean-test, Pi constant approximation with a needle etc.

But! ... those methods answer only

„*Is that sequence really perfect random?*“

The rest of our questions are unanswered...

SO here comes our method that can answer them!

And it shows up, that it can be done...

with just the most **basic quantitative linguistics tools...**

... and programming.

All we need...

n-grams

&

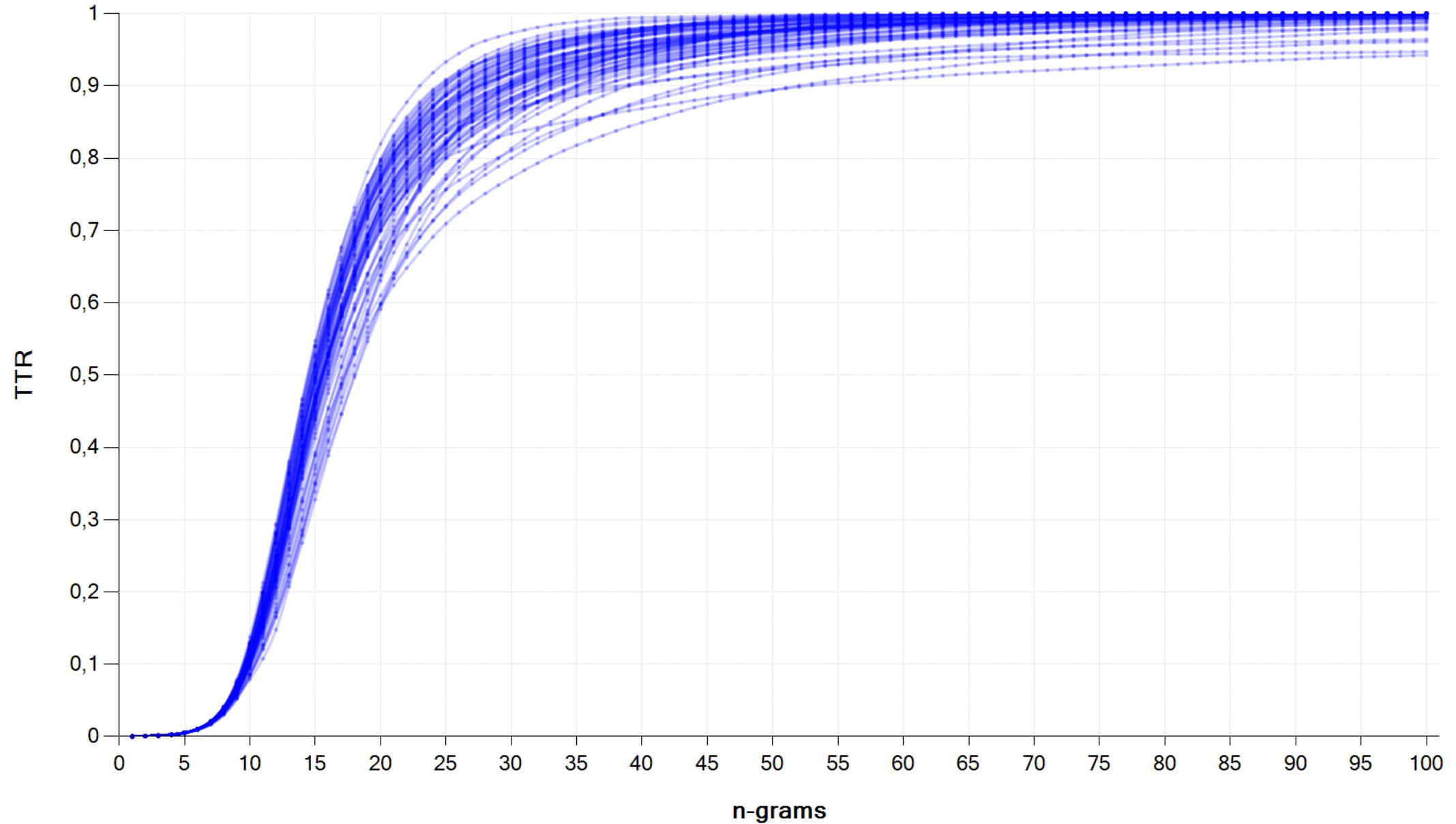
type to token ratio

and programming.

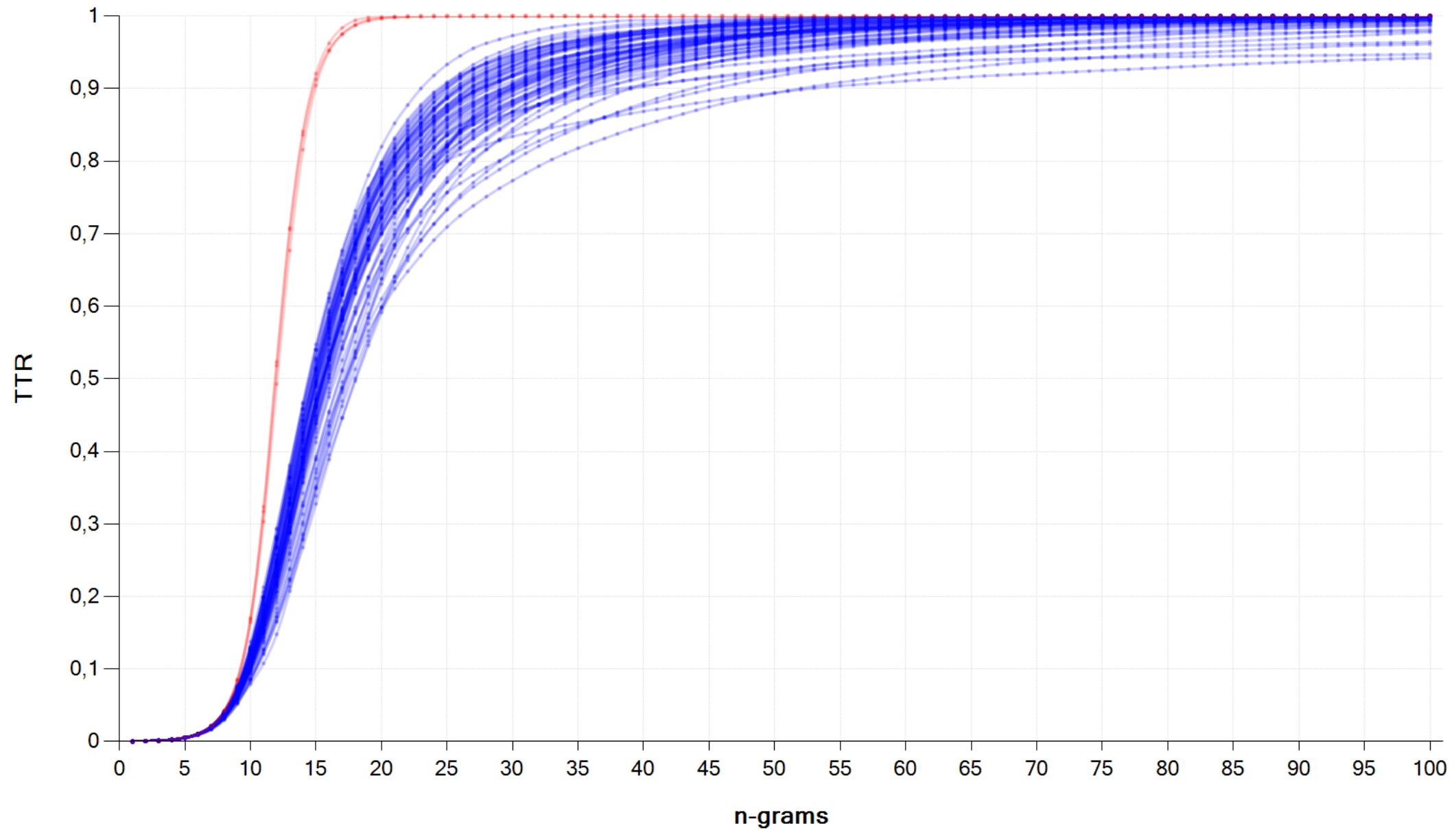
How does the method work?

Method: *Simple but effective*

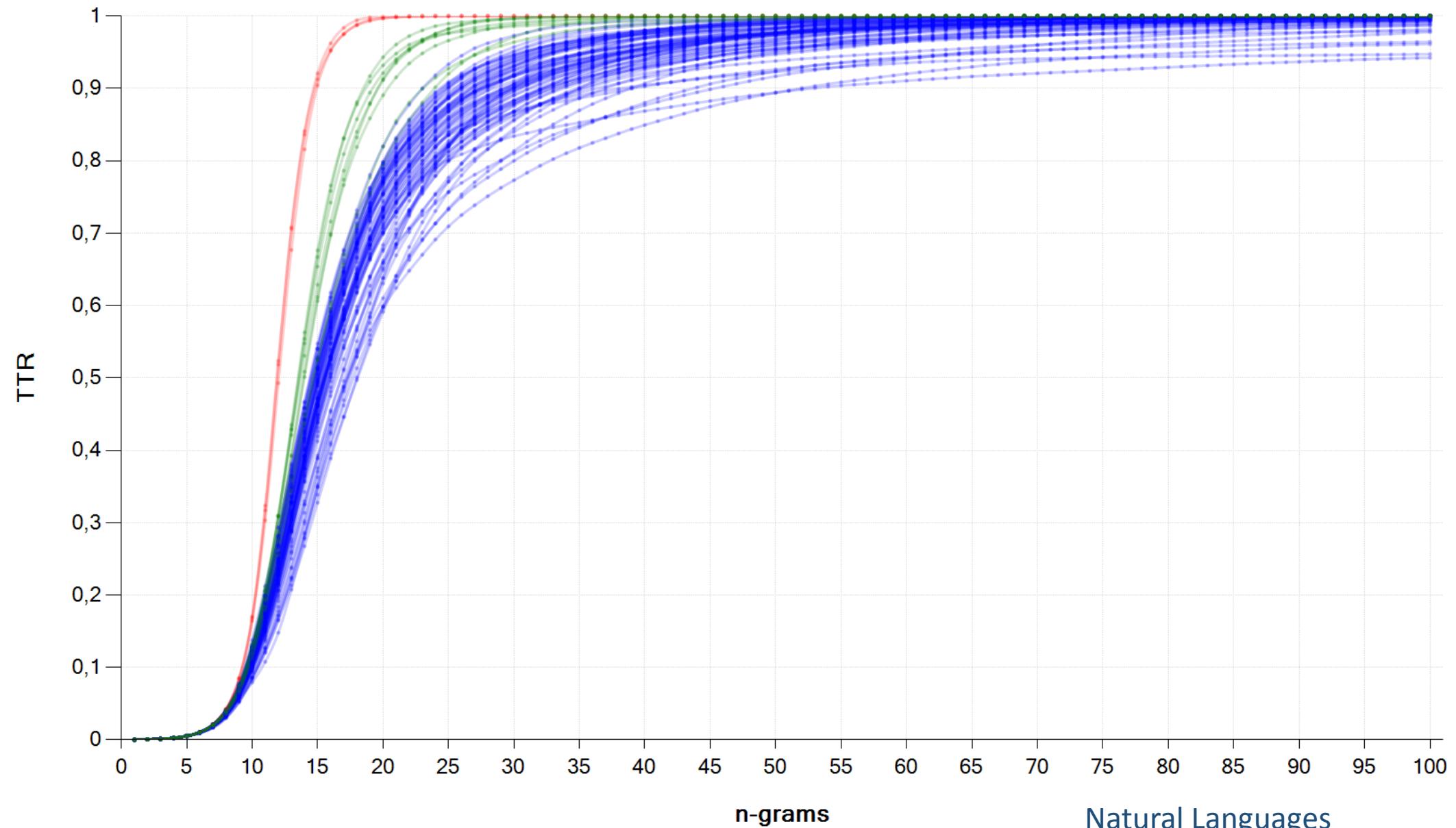
1. Take a **sequence**, in any language or of any type.
2. Convert all distinguishable symbols into **binary**
... new only binary sequence emerges
3. Take **n-grams** of the sequence:
Do for $n=\{2, k\}$, where k is $\frac{1}{2}$ of the sequence length:
Calculate Type to Token Ratio (TTR)
4. Plot **n-gram** size and calculated TTR values to chart



70 books, various languages

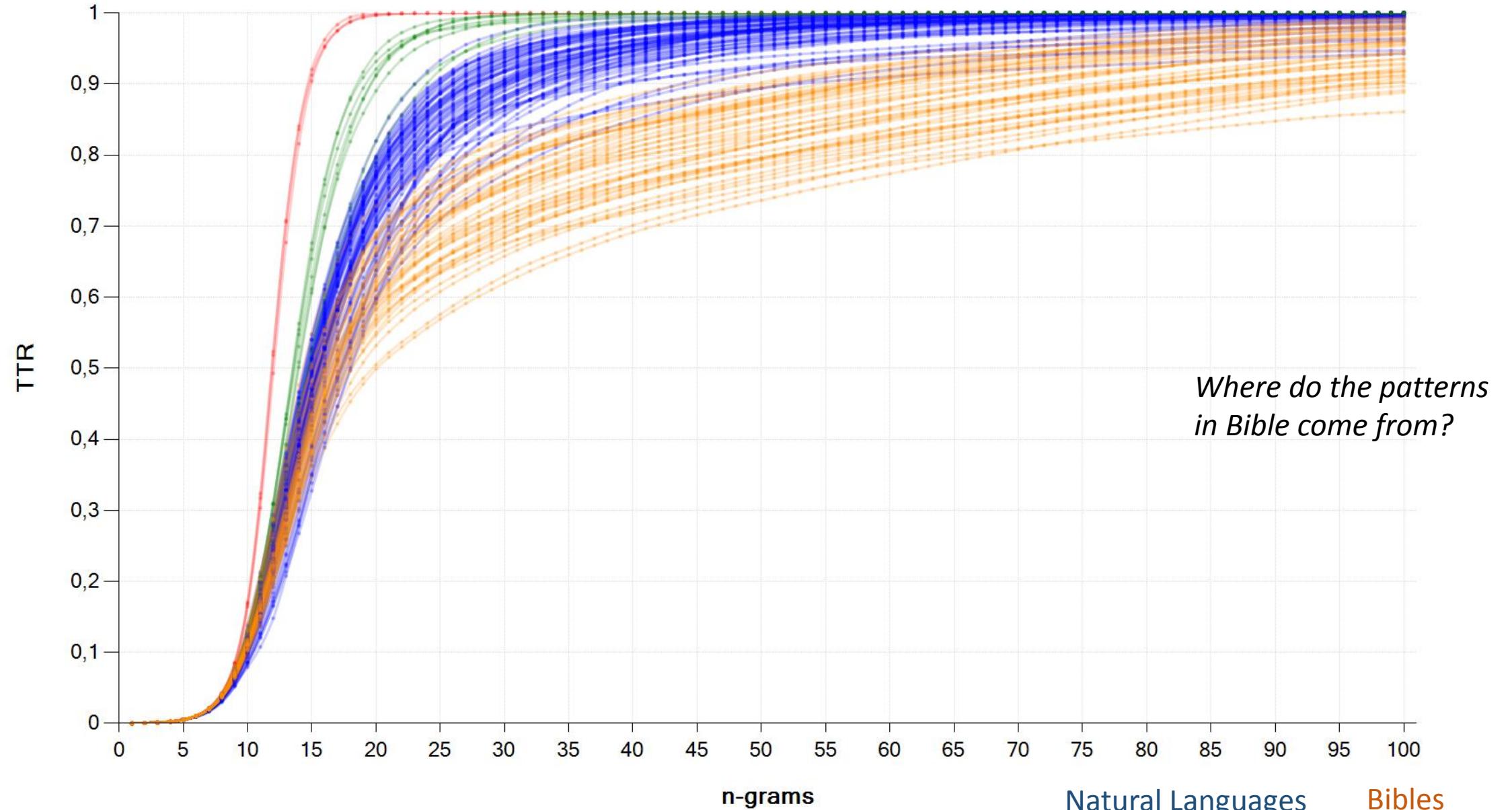


Random sequences from atmospheric noise



Monkey-typed texts

Natural Languages
True Random
Monkey-typed



Bibles in 43 languages

n-grams

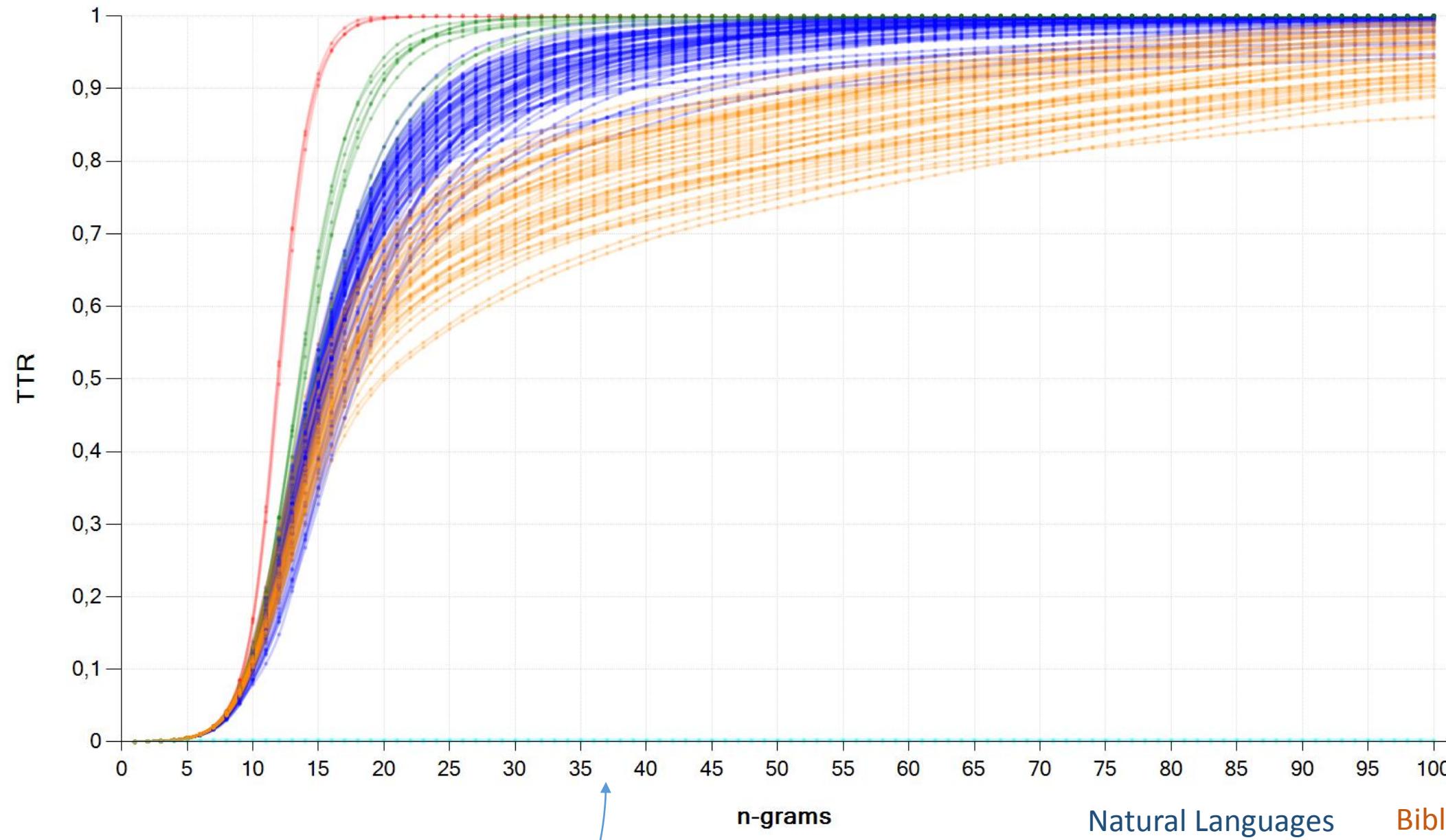
Natural Languages

Bibles

True Random

Monkey-typed

*Where do the patterns
in Bible come from?*



Only repetitions

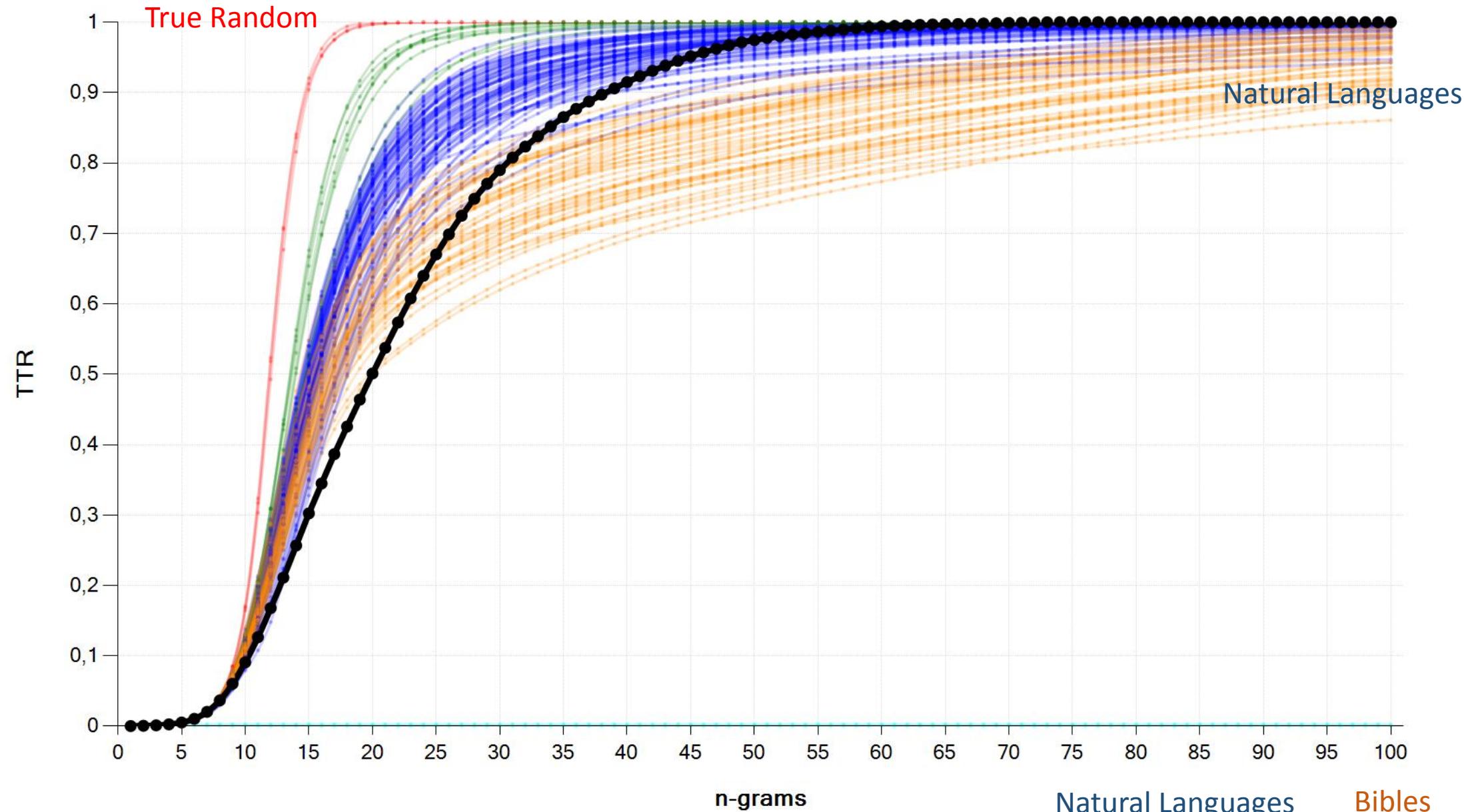
n-grams

Natural Languages

Bibles

True Random
Monkey-typed

Only repetitions



Voynich manuscript

Natural Languages
True Random
Monkey-typed

Bibles
Only repetitions
Voynich

That means:

**We can, with some degree of probability...
determine the pattern nature of sequences.**

From

- ... perfectly random sequences*
- ... to encrypted data,*
- ... monkey-typed,*
- ... natural language texts and*
- ... trivial repetitions*

We have also the next step method

for random sequences

we can distinguish the complexity of generator used

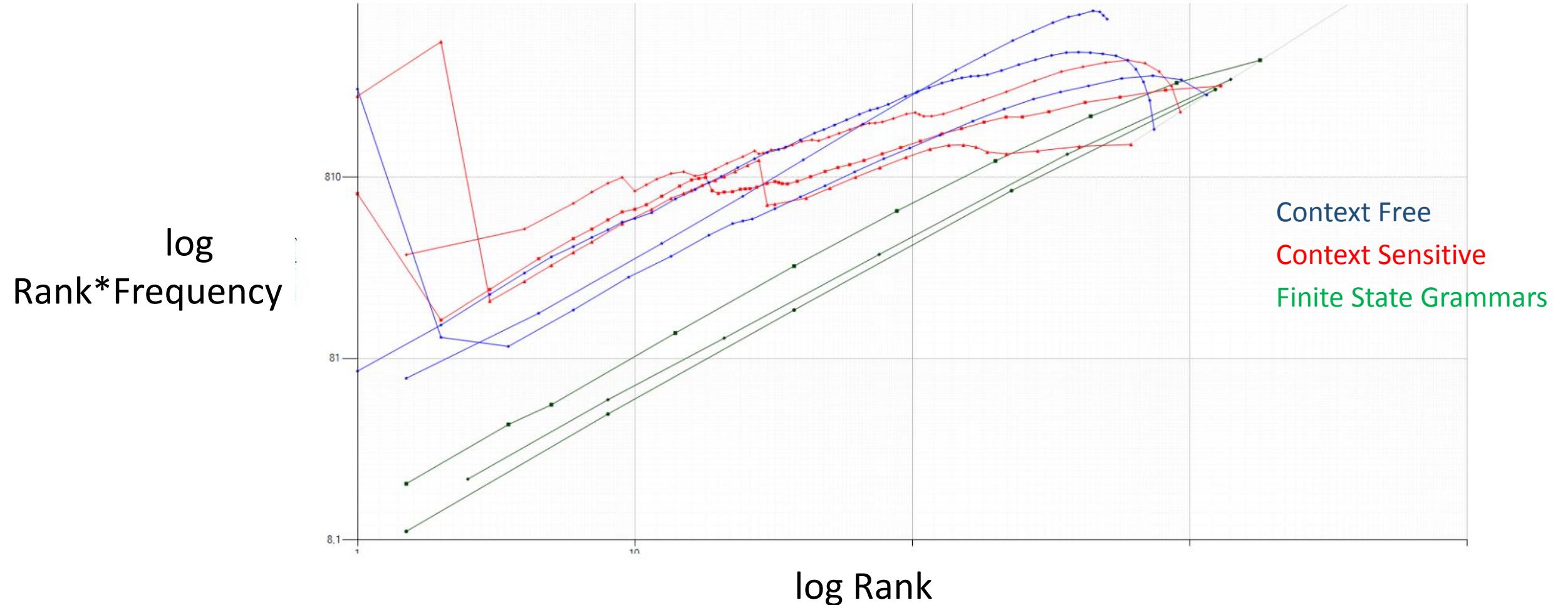
... in terms of **Chomsky grammar hierarchy**
if the generator was like...

Final state machine

Context free machine

Context sensitive machine

... the method is based on Zipfian behavior of n-grams & needs more research



A preview: Stochastic generated sequences @ 14 n-gram

Plotted Zipfian behavior can be used for stochastic generator classification

Summary: What we have?

Possible decision tree distinguishing

*True Random, Pseudo Random, Encryption,
Monkey Typing, Natural Languages , Natural Languages with added
structures, Trivial repetitions*

With a branch distinguishing

... complexity of pseudo-random generator

Summary: What's the difference with the other methods?

Other methods tend to **care only** for
randomness

&

not studying inner structure

Exceptions exist,
but no direct trend to understand „what's inside“

Applications

Linguistics

... analysing unknown sequences & languages

Computer security

... testing entropy sources in computer & classifying the generator complexity

The end

Thank you for attention

Questions?