# **ON THE NATURALNESS OF SOFTWARE**

-- Dr. Abram Hindle et al.

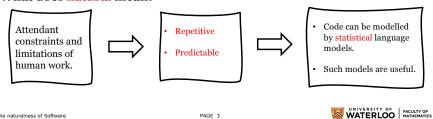
Felix Wang, David R. Cheriton School of Computer Science



# Conjecture

Most software is also natural (like other natural languages).

What does natural mean?



# **PROBLEM TO BE SOLVED**

On the naturalness of Software

PAGE 2

#### Problem to be solved

N-gram: A statistical Language Model

- Code can be modelled by statistical language
- Code (in Java, C, Python) models.

Eclipse Plug-in

PAGE 4

Such models are useful.

# **NEW IDEA**

On the naturalness of Software

PAGE 5

# N-gram

• Example paragraph:

I like pink dresses. I like pink. I do not like dresses. I like green.

### **Repetitiveness & Predictability**

• Why are language models useful for this scenario?

Because Repetitiveness & Predictability appear on:

 $\textbf{1. Lexical level} \qquad \quad \textbf{->} \qquad \textbf{Token sequences}$ 

2. Syntactic level -> Grammatical structures

3. Semantic level -> Probabilistic

On the naturalness of Software PAGE 6

Statistical Language Model

WINIVERSITY OF MACHINITY OF MATHEMATICS

## 3-gram example

PAGE 8

Tokenize:

On the naturalness of Software

<s> <s> I like pink dresses </s> </s>

<s> <s> I like pink </s> </s>

 $<\!\!\mathrm{s}\!\!><\!\!\mathrm{s}\!\!>\mathrm{I}$ do not like dresses  $<\!\!/\mathrm{s}\!\!><\!\!/\mathrm{s}\!\!>$ 

<s> <s> I like green </s> </s>

<s>: start symbol of a sentence. </s>: end symbol of a sentence.



## 3-gram example

Count all possible trigrams:

<s> <s> I: 4

 $<\!\!\mathrm{s}\!\!><\!\!\mathrm{s}\!\!>\mathrm{I}$ like pink dresses  $<\!\!/\mathrm{s}\!\!><\!\!/\mathrm{s}\!\!>$ 

<s> I like: 3

<s> <s> I like pink </s> </s> <s> <s> I do not like dresses </s> </s>

I like pink: 2

<s> <s> I like green </s> </s>

I like green: 1

I do not: 1

<s>: start symbol of a sentence.

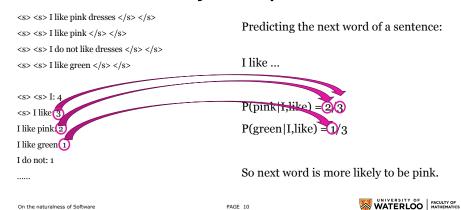
.....

</s>: end symbol of a sentence

On the naturalness of Software PAGE 9



#### 3-gram example



## 3-gram example

<s> <s> I like pink dresses </s> </s> <s> I like pink </s>	• Sparsity:	
<s> <s> I do not like dresses </s> </s>	I like dresses.	
<s> <s> I like green </s> </s>		
<s> <s> I: 4</s></s>	It's a perfectly valid sentence, but:	
<s> I like: 3</s>	P(dresses I,like) = o	
I like pink: 2		
I like green: 1		
I do not: 1	Since there was not enough learning	
	data.	
On the naturalness of Software	PAGE 11 WATERLOO	FACULTY O

#### N-gram

#### Solution:

- 1. Fall back to (N-1)-gram if the probability of N-gram is o.
- 2. Kneser-Ney Smoothing:

$$P_{\mathrm{KN}}(w_i|w_{i-n+1}^{i-1}) = \frac{\max\left(c_{KN}(w_{i-n+1}^i) - d, 0\right)}{c_{KN}(w_{i-n+1}^{i-1})} + \lambda\left(w_{i-n+1}^{i-1}\right)P_{KN}(w_i|w_{i-n+2}^{i-1})$$

An easy smoothing example (Laplace):

$$p = \frac{numerator + 1}{denominator + size \ of \ vocabulary}$$

<s> <s> I like pink dresses </s> </s>

<s> <s> I do not like dresses </s> </s>

<s> <s> I like pink </s> </s>

<s> <s> I like green </s> </s>

<s> <s> I: 4

<s> I like: 3 I like pink: 2

I like green: 1

I do not: 1

#### N-gram

<s> <s> I like pink dresses </s> </s>

<s> <s> I like pink </s> </s>

<s> <s> I do not like dresses </s> </s>

<s> <s> I like green </s> </s>

<s> <s> I: 4

<s> I like: 3

I like pink: 2

I like green: 1

I do not: 1

On the naturalness of Software

Metrics:

Cross-entropy Loss/Log-transformed perplexity:

$$H_{\mathcal{M}}(s) = -\frac{1}{n} \sum_{i=1}^{n} \log p_{\mathcal{M}}(a_i | a_1 \dots a_{i-1})$$

Rationale: How "confident" the language model is able to guess the next word.

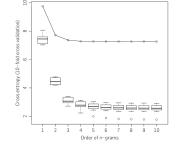
PAGE 13

WATERLOO FACULTY OF MATHEMATICS

### RQ1:Do n-gram models capture regularities in software?

PAGE 15

- Data: Java projects
- Findings:
- 1. Software unigram entropy is much lower than uniform distributions.
- 2. Cross-entropy declines rapidly with n-gram order.
- Claim:
- 1. Java contains a lot of local repetitiveness.
- 2. Software is far more regular than English.



#### **Experiment**

• Data:

Natural languages: Two famous corpora (Brown corpus, Gutenberg corpus).

Coding languages: Java projects, Ubuntu applications in C, after removing comments.

Randomly select 90% for training

 Training technique (10-fold cross-validation) WATERLOO | FACULTY OF MATHEMATICS

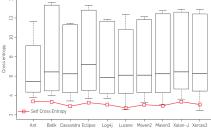
On the naturalness of Software

and 10% for validating.

PAGE 14

# RQ2:Is the local regularity language-specific or project specific?

- Data: 10 different Java projects
- Train a trigram model on each project.
- Findings:
- Self cross entropy is lower.
- Cross-project entropy is higher.



• Claim:

Each project has its own type of local, non-Java-specific regularity.

On the naturalness of Software

PAGE 16



#### RQ3: Similarities within app domain, and differences between domains?

- Data: 10 Ubuntu application domains
- Findings:
- 1. Cross-domain entropy is higher.
- 2. Within-domain entropy is lower.

Ubuntu Domain	Version		Tokens		
		Lines	Total	Unique	
Admin	10.10	9092325	41208531	1140555	
Doc	10.10	87192	362501	15373	
Graphics	10.10	1422514	7453031	188792	
Interpreters	10.10	1416361	6388351	201538	
Mail	10.10	1049136	4408776	137324	
Net	10.10	5012473	20666917	541896	
Sound	10.10	1698584	29310969	436377	
Tex	10.10	1405674	14342943	375845	
Text	10.10	1325700	6291804	155177	
Web	10.10	1743376	11361332	216474	

• Claim:

A lot of local regularity is repeated within application domains, and much less so across domains.

PAGE 17

On the naturalness of Software



# Eclipse Plug-in's achievement & findings

- 1. NGSE works best with shorter tokens (probably because coders choose shorter names for frequently used entities).
- 2. MSE saved more keystrokes compared to ECSE.

	Top 2	Top 6	Top 10
ECSE	42743	77245	95318
MSE	68798	103100	120750
Increase	61%	33%	27%

#### WATERLOO | FACULTY OF MATHEMATICS

### **Eclipse Suggestion Plug-in**

• Eclipse's built-in suggestion engine (ECSE) vs N-gram model suggestion engine (NGSE)

• NGSE: good at short tokens.

• ECSE: good at long tokens.

• MSE:

If ECSE has long tokens within the top n, take ECSE's top n offers. If not, half NGSE, half ECSE.

On the naturalness of Software PAGE 18



WATERLOO FACULTY OF MATHEMATICS

# **POSITIVES**

## Positives #1: Deterministic & Explainable

- The choice of the N-gram model is good.
- The N-gram model is deterministic and explainable. But many language models based on optimization algorithms are not!
- The scope of this paper is to prove: Software is repetitive & predictable.
- Good for a study to "prove" a hypothesis.

On the naturalness of Software

PAGE 21



My critiques of this paper could be fundamentally wrong due to my shallow research experience, please DO correct me during the discussion section.

#### Positives #2: Strong Empirical Evidence

- · Large, real-world dataset
- 10 major Java projects
- 10 Ubuntu domains of C applications.
- Large-scale natural language corpus
- Brown: ~1 million words
- Gutenberg: ~2.5 million words
- Rigorous procedure: 10-fold cross-validation.
- Robust and strong empirical evidence.

On the naturalness of Software

PAGE 22



# Negative #1: Is repeating a good sign?

- Repetition does mean regularity, but it is not always a good sign in Software Engineering.
- Any preprocessing on the codebase that the n-gram model is trained on?
  - removed comments
- tokenized codes
- Any legacy code? Poor practices that caused redundancy (that caused repetition)?
- Code suggestions to the programmer
   As good as possible? Or as repetitive as possible? As a statistical model's natural rationale is to suggest something as repetitive as possible.



#### Negative #2: The MSE algorithm is too superficial

- Eclipse's built-in suggestion engine (ECSE): "are typically based on type information available in context".
  - ->Miss statistical patterns
- N-gram model suggestion engine (NGSE): a statistical model that considers partial semantic, lexical, and syntactic information. ->Miss type information.
- Can we build something by combining the useful information from 2, not by simply choosing one or another?

On the naturalness of Software

PAGE 25



## Negative #4: The scope of this paper

- Most software is natural -> It is also likely to be repetitive and predictable
- Most software is repetitive and predictable -> Most software is natural?
- Not a sufficient and necessary condition!

We begin with the conjecture that most software is also natural, in the sense that it is created by humans at work, with all the attendant constraints and limitations—and thus, like natural language, it is also likely to be repetitive and predictable. We then proceed to ask whether (a)

WATERLOO | FACULTY OF MATHEMATICS

#### Negative #3: Limitations of N-gram

- N-gram's ability to capture grammatical structure is extremely limited.
- Able to: catch local grammatical structures like: ``public static void``, ``System.out.println(``, ``for (i = 0; i < n; i++) {``
- Unable to catch: ``try {//~100 words } catch {}``.
- It will need a 100-gram model! Which will be extremely sparse
- Obviously, there's some grammatical structure between ``try{`` and ``} catch{``.
- N-grams cannot catch some long grammatical structures.

On the naturalness of Software

PAGE 26



## Negative #4: The scope of this paper

- The author did not formally state that "software is natural". It instead said: "(Software) it's regular", " (Software) it's not random", or this experiment supports the conjecture.
- Because this evaluation cannot prove the conjecture, and software really isn't fully a natural language by any means (tolerance of errors? one keyword can have multiple meanings?).
- Do we really need to prove this catchy conjecture?
  - Proving that software is repetitive and predictable is already a huge contribution!

# **FUTURE WORK**

On the naturalness of Software PAGE 29

#### **Probabilistic Context Free Grammar (PCFG)**

PAGE 31

$S \rightarrow NP VP$	[.80]	$  Det \rightarrow that [.05]   the [.80]  a$	[.15]
$S \rightarrow Aux NP VP$	[.15]	Noun → book	[.10]
$S \rightarrow VP$	[.05]	Noun → flights	[.50]
NP → Det Nom	[.20]	Noun → meal	[.40]
NP → Proper-Noun	[.35]	Verb → book	[.30]
$NP \rightarrow Nom$	[.05]	Verb → include	[.30]
NP → Pronoun	[.40]	Verb → want	[.40]
Nom → Noun	[.75]	$Aux \rightarrow can$	[.40]
Nom → Noun Nom	[.20]	Aux → does	[.30]
Nom → Proper-Noun Nom	[.05]	$Aux \rightarrow do$	[.30]
$VP \rightarrow Verb$	[.55]	Proper-Noun → TWA	[.40]
$VP \rightarrow Verb NP$	[.40]	Proper-Noun → Denver	[.40]
$VP \rightarrow Verb NP NP$	[.05]	Pronoun $\rightarrow you[.40] \mid I[.60]$	

References: Johnson, David, and May Young. "Course: CPSC522/PCFG - UBC Wiki." Wiki.ubc.ca, 2017, wiki.ubc.ca/Course:CPSC522/PCFG. Accessed 20 Sept.

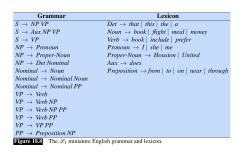
On the naturalness of Software

 Construct an abstract syntax tree of what grammar is allowed in the training data, along with the training process.

 Construct an abstract syntax tree, but this time with a probability.



#### **Context Free Grammar (CFG)**



· CFG: A concept related to traditional NLP aiming to catch the grammatical structure of natural languages.

NP: Noun Phrase VP: Verb Phrase PP: Prepositional phrases Det: Determiner Aux: Auxiliary

· Construct an abstract syntax tree of what grammar is allowed in the training data, along with the training process.

References: Jurafsky, Daniel, and James H. Martin. Speech and Language I rocessing, Computational Linguistics, and Speech Recognition with Language Models. 3rd edn. 2025, Language Processing, Computational Linguistics, and Speech Recognition with Language Models. 3rd edn. 2025, WATERLOO RECOGNITION OF LANGUAGE TO A CONTROL OF LANGUAGE TO A CONT

#### **Combine PCFG and N-gram**

- The N-gram model can work with PCFG, and it's a well-known technique in the field of traditional NLP [1].
- Incorporate these two together, and let the N-gram model be capable of catching complicated grammar in code, and repeat the procedure of this paper again, to see if the model catches something new.
- References
  - [1] Pauls, Adam, and Dan Klein. "Large-scale syntactic language modeling with treelets." Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2012.

PAGE 32

# **RATING**

On the naturalness of Software PAGE 33

# **DISCUSSION POINTS**

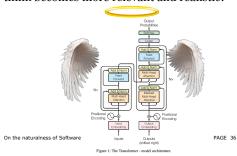
# 3.5/5 (GOOD PAPER BUT NOT LIFE-CHANGING)

The idea of naturalness of this paper is novel, the work done is rigorous, but its reliance on the N-gram model limits its ability to catch something more in-depth.

PAGE 34

On the naturalness of Software

1. One and a half years after this paper was posted on Communications of the ACM, a paper called "Attention is all you need" came out, and everything about language models, computer vision, speech recognition, and machine translation changed fundamentally. For the future directions section envisioned by Dr. Hindle (section 6), which direction do you think becomes irrelevant, and which direction do you think becomes more relevant and realistic?



References: Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).



On the naturalness of Software PAGE 35

2. The N-gram-based Eclipse plug-in showed 27–61% improvement in code suggest the measure used is the number of keystroke savings. Do you think this result convidemonstrated practical benefits, or is it more of a proof-of-concept?

What makes you think it indicates practical benefits? Or what follow-up study coul to make it more convincing?

- 3. Should "Quality of Code" play into our proof for naturalness? If yes, how can we incorporate it? If no, please justify it.
- 4. If we now have Eclipse's built-in suggestion engine and N-gram model suggestio (NGSE) both in hand, how could we build something potentially better than MSE?
- 5. Any disagreement or agreement with my critiques before? Do you have other posnegatives that you want to share?

On the naturalness of Software

PAGE 37



stion. And rincingly			
ld be done			
on engine			
sitives or			
Y OF FACULTY OF MATHEMATICS			