

# Grammar compression with probabilistic context-free grammar

Hiroaki Naganuma\*, Diptarama Hendrian\*, Ryo Yoshinaka\*, Ayumi Shinohara\*, and Naoki Kobayashi†

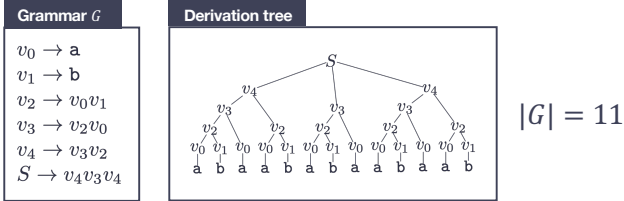
\*Tohoku University, †The University of Tokyo

**Abstract:** We propose a new approach for universal lossless text compression, based on grammar compression. In the literature, a target string  $T$  has been compressed as a context-free grammar  $G$  in Chomsky normal form satisfying  $L(G) = T$ . Such a grammar is often called a straight-line program (SLP). In this work, we consider a probabilistic grammar  $G$  that generates  $T$ , but not necessarily as a unique element of  $L(G)$ . In order to recover the original text  $T$  unambiguously, we keep both the grammar  $G$  and the derivation tree of  $T$  from the start symbol in  $G$ , in compressed form. We show some simple evidence that our proposal is indeed more efficient than SLPs for certain texts, both from theoretical and practical points of view.

## Existing approaches for grammar compression

### Approach 1: Universal

- Given a text  $T$ , an encoder constructs CFG  $G$  such that  $L(G) = \{T\}$ .

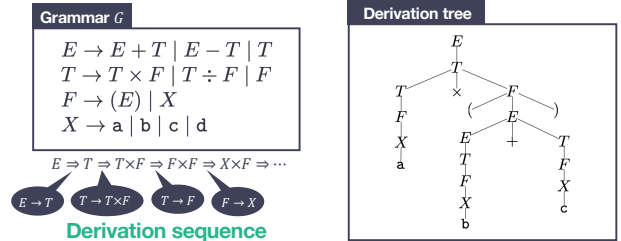


An encoder tries to construct the minimum CFG.  
 (The smallest grammar problem [Charikar+, 2005])

$|G|$ : A grammar size of CFG  $G$ .  
 It is defined as the total length of strings on the right hand sides of all production rules.

### Approach 2: Domain-specific

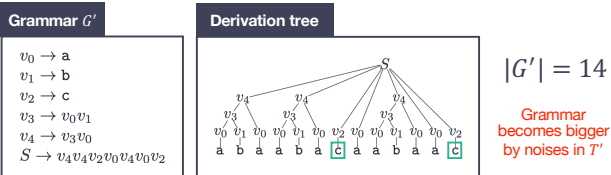
- A fixed CFG  $G$  is used when an encoder compresses an input text  $T$  such that  $T \in L(G)$
- A compressed data is a derivation sequence of  $T$  in  $G$  which is a sequence of production rules in  $G$ .



## Our proposal framework of grammar compression using PCFG

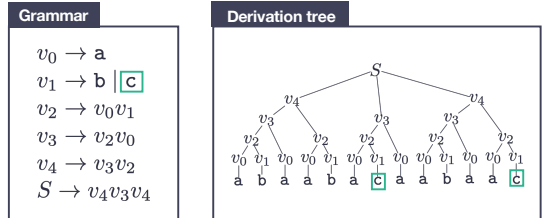
### Disadvantage of approach 1: sensitive to noise

The example of CFG  $G'$  constructed from  $T' = \text{abaaba} \square \text{aaba} \square$



Idea for improvement

- Given a text  $T$ , an encoder constructs CFG  $G$  such that  $T \in L(G)$ .
- A compressed data consists of  $G$  and a derivation sequence of  $T$  in  $G$ .
- A derivation sequence is encoded by **arithmetic coding**.



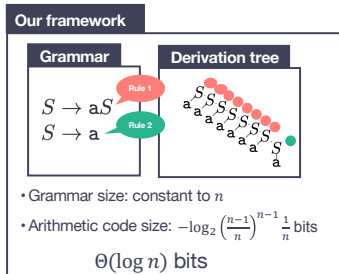
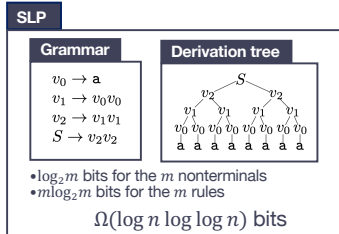
### Disadvantage of approach 2: no universality

We cannot compress a text  $T$  such that  $T \notin L(G)$  effectively.

## Effectiveness of our proposed scheme

### Theoretical result

Comparison of SLP compression and our framework for  $T = a^n$  where  $n = 2^m$



### Experimental result

Compressing "noisy" 20th Fibonacci strings (10946 bytes) by ideal grammars in our framework, grammars output by our prototype compression algorithm, and some existing compressors

