Introduction 000	Edit Distance R 0000 C	Regular DOOOO	Repeated Variable	s Overview	Outlook 00
	Marah's Darres		\/		
	Watching Patter	rns with	variables	s under Edi	τ
		Distar	nce		
		Biotai			
	Paweł Gawrychowsk	i Florin	Manea S	stefan Siemer	
	T awer Gawryenowsk			Ceran Orenier	

University of Wrocław, Göttingen University

gawry@cs.uni.wroc.pl, florin.manea@cs.uni-goettingen.de, stefan.siemer@cs.uni-goettingen.de

RP 2022

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
●00	0000	00000		00	00

## Introduction

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
○●○	0000	00000		00	00

## Pattern Matching

## Pattern Matching Problem

Find occurrences of a pattern  $\alpha$  in a word w.

Example:  $\alpha = bab$ , w = aaaababbb

aaaababbb

bab

Paweł Gawrychowski, Florin Manea, Stefan Siemer

## Pattern Matching with Variables

Exact Matching Problem for P: Match<sub>P</sub> Input: A pattern  $\alpha \in P$ , with  $|\alpha| = m$ , a word w, with |w| = n. Question: Is there a substitution h with  $h(\alpha) = w$ ?

Outlook

Example:  $\alpha = x_1 x_1 bab x_2 x_2$ , w = aaaababbb

$\mathtt{x_1} \to \mathtt{aa}$	aaaababbb
$\mathtt{x}_2 \to b$	aaaababbb

Example:  $\alpha = x_1 babx_2$ , w = aaaababbb

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
000	●000	00000		00	00

## Edit Distance

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction 000	Edit Distance 0●00	Regular 00000	Repeated Variables	Overview 00	Outlook 00

## Edit Distance

### String metric

For words  $u, w \in \Sigma^*$ , the *edit distance* between u and w is defined as the minimal number  $d_{ED}(u, w)$  of letter insertions, letter deletions, and letter to letter substitutions which one has to apply to u to obtain w.

Example:  $u = abbab, w = baaba, \Delta = 3$ 

$$abbab \xrightarrow[deletion]{} bbab \xrightarrow[substitution]{} baab \xrightarrow[insertion]{} baaba$$

Paweł Gawrychowski, Florin Manea, Stefan Siemer

# Pattern Matching with Variables and Edit Distance

Approximate Matching Decision Problem for P: MisMatchPInput:A pattern  $\alpha \in P$ , with  $|\alpha| = m$ , a word w, with|w| = n, an integer  $\Delta \leq m$ .Question:Is  $d_{\rm ED}(\alpha, w) \leq \Delta$ ?

Example:  $\alpha = x_1x_1babx_2x_2$ , w = aaababbb,  $\Delta = 1$ 

$\mathtt{x_1} \to \mathtt{aa}$	aaaababbb
$\mathtt{x_2} \to b$	aaababbb

For Hamming Distance results see paper at MFCS 2021.

## Pattern Matching with Variables and Edit Distance

Approximate Matching Minimisation Problem for P: MinMisMatch<br/>PInput:A pattern  $\alpha \in P$ , with  $|\alpha| = m$ , a word w, with<br/>|w| = n.Question:Compute  $d_{\rm ED}(\alpha, w)$ .

Example:  $\alpha = x_1x_1babx_2x_2$ , w = aaababbb

$\mathtt{x_1}  ightarrow \mathtt{aa}$	aaaababbb
$\mathtt{x_2} \to b$	aaababbb

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction 000	Edit Distance 0000	Regular ●0000	Repeated Variables 0000	Overview 00	Outlook 00



Paweł Gawrychowski, Florin Manea, Stefan Siemer

Repeated Variables

Overview 00 Outlook 00

# Regular Pattern Definition

## Definition

$$\alpha \in \text{Reg if } \alpha = w_0 \prod_{i=1}^{M} (x_i w_i), \text{ with } w_i \in \Sigma^{\star}.$$

Example:  $\alpha = abx_1abx_2x_3baab$ .

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction 000	Edit Distance 0000	Regular 00●00	Repeated Variables	Overview 00	Outlook 00
Context					

- MisMatch<sub>Reg</sub> with distance  $\Delta = 0$  in  $\mathcal{O}(n)$ .
- MisMatch<sub>Reg</sub> (HD) in  $\mathcal{O}(n\Delta)$ .
- $x_1 w_1 x_2$  known as Landau and Vishkin Algo in  $\mathcal{O}(n\Delta)$ .
- We extend the idea of Landau and Vishkin to achieve  $\mathcal{O}(n\Delta)$ .

Introduction 000	Edit Distance 0000	Regular 000●0	Repeated Variables	Overview 00	Outlook 00
Algorith	m outline				

- Interpret regular variables as free insertions on that position.
- First dynamic programming in  $\mathcal{O}(nm)$ .
- For increasing distance calculate maximal matching expansion on a diagonal in the DP matrix (Landau Vishkin idea).
- Consider surrounding diagonals and furthest previously reached variable + LCP extensions.
- Compute  $\Delta$  distance increment steps for *n* diagonals, hence  $\mathcal{O}(n\Delta)$ .

## Rectangular lower bound

Hardness follows directly from rectangular lower bound of edit distance (Backurs, Indyk).  $\alpha = xuy$ , where u is a string of terminals and x and y are variables.

#### Theorem

 $MisMatch_{Reg}$  can not be solved in time  $\mathcal{O}(|w|^h \Delta^g)$  (or  $\mathcal{O}(|w|^h |\alpha|^g)$ ) where  $h + g = 2 - \epsilon$  with  $\epsilon > 0$ , unless the Orthogonal Vectors Conjecture fails

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
000	0000	00000	●000	00	00

## **Repeated Variables**

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
000	0000	00000	0●00	00	00

## 1Var Definition

#### Definition

# $\alpha \in 1$ Var if there exists exactly one variable $x_1$ with several occurences.

Example:  $\alpha = abx_1x_1abx_1x_1x_1baabx_1 \in 1$  Var.

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
000	0000	00000	00●0	00	00

## Algorithm outline

#### Theorem

MisMatch<sub>1Var</sub> and MinMisMatch<sub>1Var</sub> can be solved in  $O(n^{3|\alpha|_{x_1}})$  time, where  $x_1$  is the single variable occurring in  $\alpha$ .

- Brute force possible intervals for all occurrences of x<sub>1</sub>.
- Edit Distance Median-String problem for selected intervals (Sankoff 75).

Introduction 000	Edit Distance 0000	Regular 00000	Repeated Variables 000●	Overview 00	Outlook 00
Deduct	ian				

#### Theorem

Vennenon

 $MisMatch_{1Var}$  is W[1]-hard w.r.t. the number of occurrences of the single variable x of the input pattern  $\alpha$ .

 $\begin{array}{lll} \mbox{Median String: MS} \\ \mbox{Input:} & k \mbox{ strings } w_1, \ldots, w_k \in \sigma^* \mbox{ an integer } \Delta. \\ \mbox{Question:} & \mbox{Does there exist a string } s \mbox{ such that} \\ & \sum_{i=1}^k d_{\rm ED}(w_i,s) \leq \Delta? \\ & (\mbox{The string } s \mbox{ for which } \sum_{i=1}^k d_{\rm ED}(w_i,s) \mbox{ is minimum is} \\ & \mbox{ called the median string of the strings } \{w_1, \ldots, w_k\}. \\ \end{array}$ 

**PM instance:** w encodes the k strings seperated by long borders of two fresh symbols.  $\alpha$  encodes the same long borders with  $x_1$  in between.

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction 000	Edit Distance 0000	Regular 00000	Repeated Variables 0000	Overview ●0	Outlook 00



Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction	Edit Distance	Regular	Repeated Variables	Overview	Outlook
000	0000	00000		⊙●	00

## Result Overview

Class	$Match(w, \alpha)$	MisMatch( $w, \alpha, \Delta$ )	$MisMatch(w, \alpha, \Delta)$
		for $d_{\text{HAM}}(\cdot, \cdot)$	for $d_{\text{ED}}(\cdot, \cdot)$
Reg	O(n) [folklore]	$O(n\Delta)$ , matching	$O(n\Delta)$ , matching
		cond. lower bound	cond. lower bound
1Var	O(n) [folklore]	<i>O</i> ( <i>n</i> )	$O(n^{3 \alpha _{\times}})$
$(\operatorname{var}(\alpha) = \{x\})$			W[1]-hard w.r.t. $ \alpha _x$
NonCross	$O(nm \log n)$ [1]	$O(n^3p)$	NP-hard
1RepVar	$O(n^2)$ [1]	$O(n^{k+2}m)$	NP-hard for $k \ge 1$
k = # x-blocks		W[1]-hard w.r.t. <i>k</i>	
kLOC	$O(mkn^{2k+1})$ [2]	$O(n^{2k+2}m)$	NP-hard for $k \ge 1$
	W[1]-hard w.r.t. <i>k</i>	W[1]-hard w.r.t. <i>k</i>	
kSCD	$O(m^2 n^{2k})$ [1]	NP-hard for $k \ge 2$	NP-hard for $k \ge 1$
	W[1]-hard w.r.t. k		
kRepVar	$O(n^{2k})$ [1]	NP-hard for $k \ge 1$	NP-hard for $k \ge 1$
	W[1]-hard w.r.t. <i>k</i>		

<sup>1</sup>Henning Fernau, Florin Manea, Robert Mercas, and Markus L. Schmid. Pattern matching with variables...

<sup>2</sup> Joel D. Day, Pamela Fleischmann, Florin Manea, and Dirk Nowotka. Local patterns

<sup>3</sup>Daniel Reidenbach and Markus L. Schmid. Patterns with bounded treewidth.

Paweł Gawrychowski, Florin Manea, Stefan Siemer

Introduction 000	Edit Distance 0000	Regular 00000	Repeated Variables 0000	Overview 00	Outlook ●0



Paweł Gawrychowski, Florin Manea, Stefan Siemer



## Extensions & Outlook

- Matching lower and upper bounds for non-cross (HD).
- Enumeration algorithms
- Bounding the number of variables.
- Other metrics (Damerau-Levenshtein).
- Combining exact and approximate matching.
- Restrictions on Variables (e.g. RegEx membership).
- Analysing use cases (e.g. database theory, learning theory).



## Extensions & Outlook

- Matching lower and upper bounds for non-cross (HD).
- Enumeration algorithms
- Bounding the number of variables.
- Other metrics (Damerau-Levenshtein).
- Combining exact and approximate matching.
- Restrictions on Variables (e.g. RegEx membership).
- Analysing use cases (e.g. database theory, learning theory).

Thank you.