

Subsequences With Gap Constraints: Complexity Bounds for Matching and Analysis Problems

Joel D. Day ¹ Maria Kosche ² Florin Manea ² Markus Schmid ³

¹Loughborough University, UK

²Göttingen University, Germany

³Humboldt University of Berlin, Germany

16th International Conference on Reachability Problems
17 – 19 October 2022
Kaiserslautern

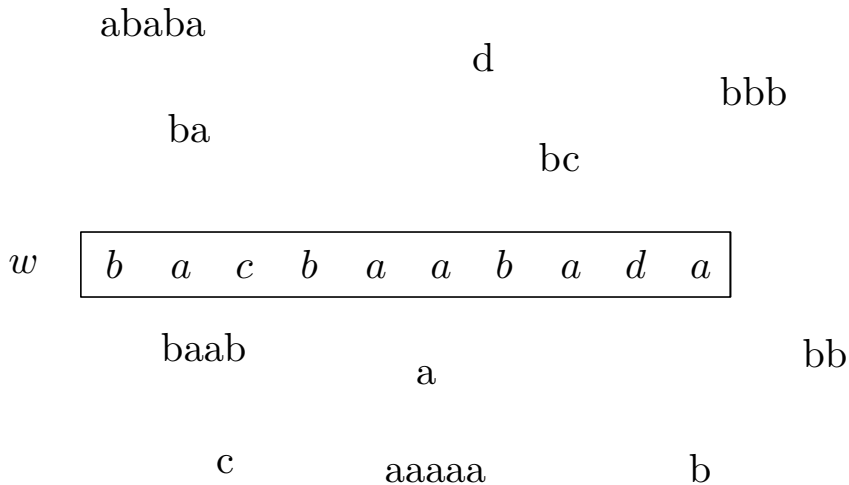
Classical Subsequences

Subsequences

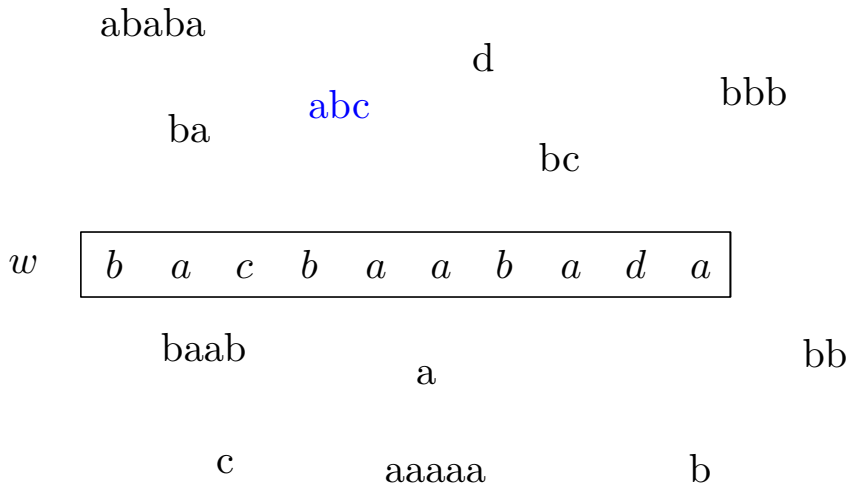
w

b	a	c	b	a	a	b	a	d	a
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

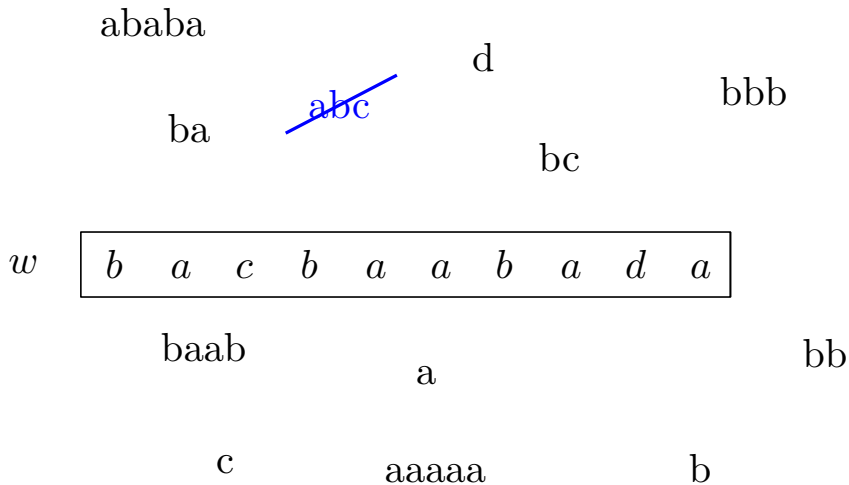
Subsequences



Subsequences



Subsequences



Subsequences

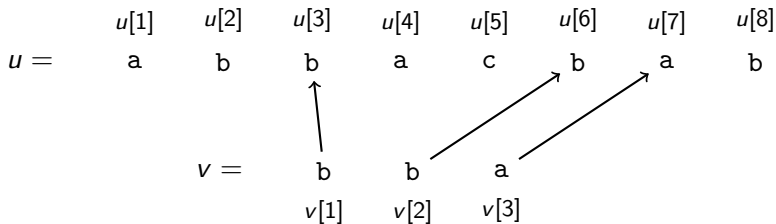
Σ is a finite alphabet, e. g., $\Sigma = \{a, b, c, d\}$.

	$u[1]$	$u[2]$	$u[3]$	$u[4]$	$u[5]$	$u[6]$	$u[7]$	$u[8]$
$u =$	a	b	b	a	c	b	a	b

$v =$	b	b	a
	$v[1]$	$v[2]$	$v[3]$

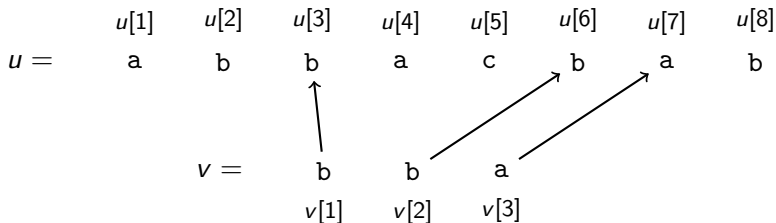
Subsequences

Σ is a finite alphabet, e. g., $\Sigma = \{a, b, c, d\}$.



Subsequences

Σ is a finite alphabet, e. g., $\Sigma = \{a, b, c, d\}$.

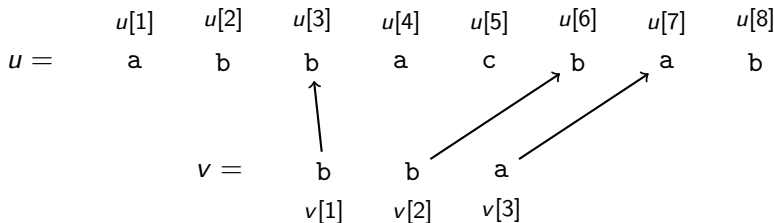


Notation

Embedding: $e : \{1, \dots, |v|\} \rightarrow \{1, \dots, |u|\}$ with $e(1) < \dots < e(|v|)$.

Subsequences

Σ is a finite alphabet, e. g., $\Sigma = \{a, b, c, d\}$.



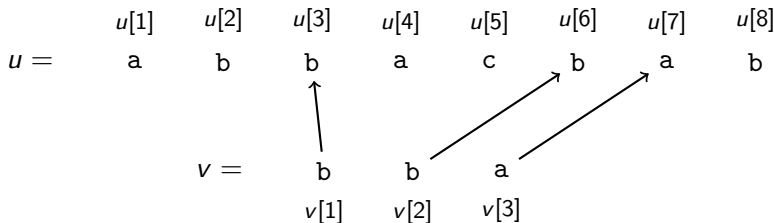
Notation

Embedding: $e : \{1, \dots, |v|\} \rightarrow \{1, \dots, |u|\}$ with $e(1) < \dots < e(|v|)$.

$v \preceq_e u$: e is an embedding with $v[i] = u[e(i)] \ \forall i \in \{1, 2, \dots, |v|\}$.

Subsequences

Σ is a finite alphabet, e. g., $\Sigma = \{a, b, c, d\}$.



Notation

Embedding: $e : \{1, \dots, |v|\} \rightarrow \{1, \dots, |u|\}$ with $e(1) < \dots < e(|v|)$.

$v \preceq_e u$: e is an embedding with $v[i] = u[e(i)] \ \forall i \in \{1, 2, \dots, |v|\}$.

v is subsequence of u ($v \preceq u$) if there is some embedding e with $v \preceq_e u$.

Subsequences

Set of Subsequences of length k

For given integer k and word w , let $Subseq(k, w)$ denote the set of subsequences of length k of w .

$$Subseq(k, w) = \{v \in \Sigma^k \mid v \preceq w\}.$$

Example:

$$Subseq(3, abbacbab) = \{aaa, aab, aac, aba, abb, abc, aca, acb, \\ baa, bab, bac, bba, bbb, bbc, bca, bcb, \\ cab, cba, cbb\}$$

Subsequences are a central concept in many different areas of TCS:

- ▶ Formal languages and logics (piecewise testable languages, subword order and downward closures).
- ▶ Combinatorics on words.
- ▶ Modelling concurrency.
- ▶ Database theory (event stream processing).
- ▶ Algorithms (longest common subsequence, shortest common supersequence).

Computational Problems for Subsequences

Matching

Input: $u, v \in \Sigma^*$

Question: $v \preceq u$?

Computational Problems for Subsequences

Matching

Input: $u, v \in \Sigma^*$

Question: $v \preceq u?$

Analysis Problems

Input: $u, v \in \Sigma^*, k \in \mathbb{N}$

Questions: $Subseq(k, u) = Subseq(k, v)?$ (Equivalence)

$Subseq(k, u) \subseteq Subseq(k, v)?$ (Containment)

$Subseq(k, u) = \Sigma^k?$ (Universality)

Subsequence Problems are Simple

a b c b c a b c a b a c

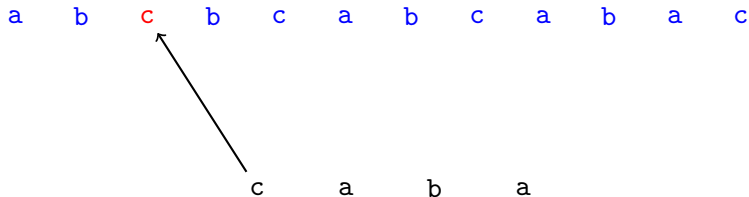
c a b a

Subsequence Problems are Simple

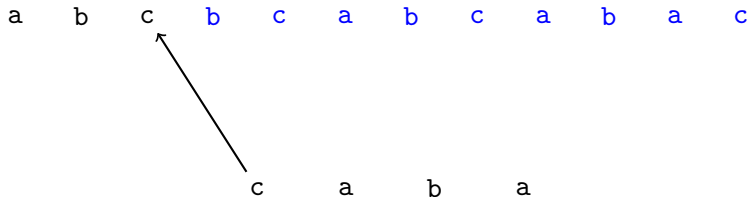
a b c b c a b c a b a c

c a b a

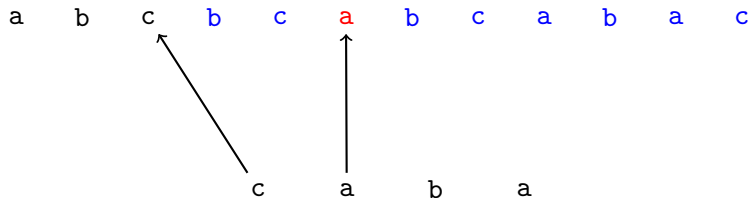
Subsequence Problems are Simple



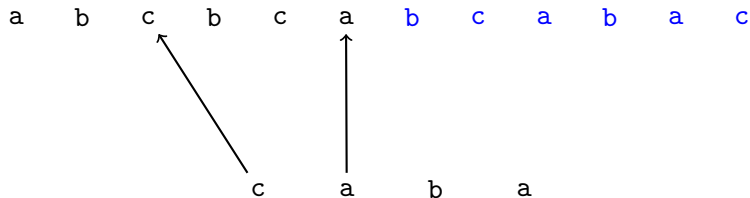
Subsequence Problems are Simple



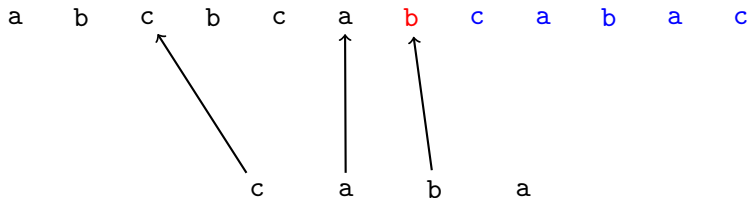
Subsequence Problems are Simple



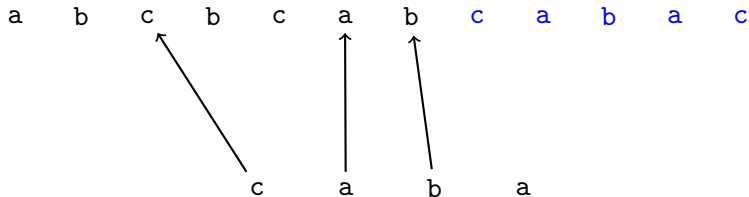
Subsequence Problems are Simple



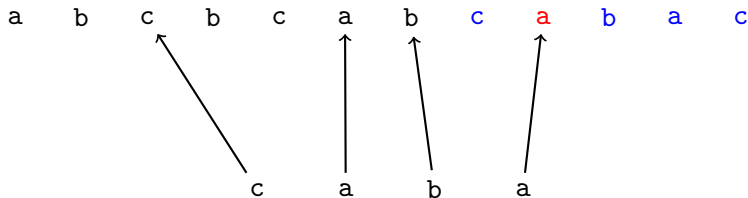
Subsequence Problems are Simple



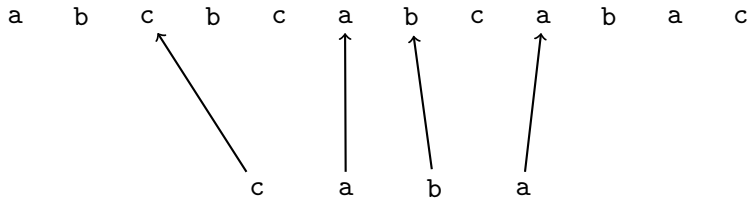
Subsequence Problems are Simple



Subsequence Problems are Simple



Subsequence Problems are Simple



Subsequence Problems are Simple

Observation

Matching is trivially solvable in linear time.

Subsequence Problems are Simple

Observation

Matching is trivially solvable in linear time.

Observation

We can construct a DFA for $\text{Subseq}(k, u)$ of size $O(k|u|)$.
 \Rightarrow the analysis problems can be solved in polynomial time.

Subsequence Problems are Simple

Observation

Matching is trivially solvable in linear time.

Observation

We can construct a DFA for $\text{Subseq}(k, u)$ of size $O(k|u|)$.
 \Rightarrow the analysis problems can be solved in polynomial time.

Theorem (Gawrychowski et al., STACS 2021)

Equivalence can be decided in linear time.

Subsequences with Gap Constraints

Shift from Classical Scenario to Gap Constraints

Classical subsequences ...

... are usually considered with arbitrary embeddings.

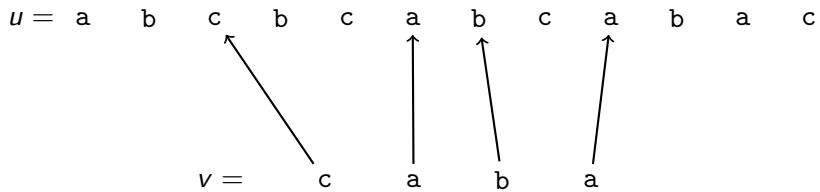
For practical scenarios, it is reasonable to introduce *gap constraints*.

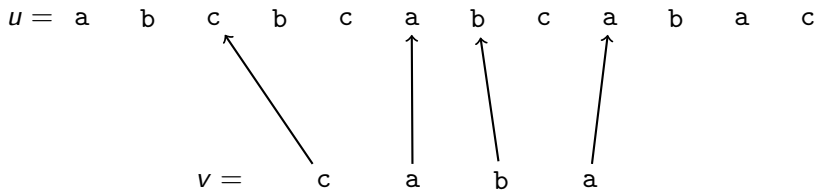
We restrict the length of gaps.

- ▶ Alignments of bio-sequences.
- ▶ Modelling single processor scheduling with fairness properties.

And we restrict allowed symbols that can occur in gaps.

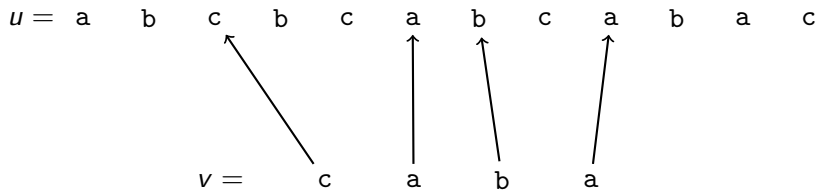
- ▶ Complex event processing → forbidding events in specific positions of a subsequence.





$$\text{gap}_e(u, i) = u[e(i) + 1..e(i + 1) - 1]$$

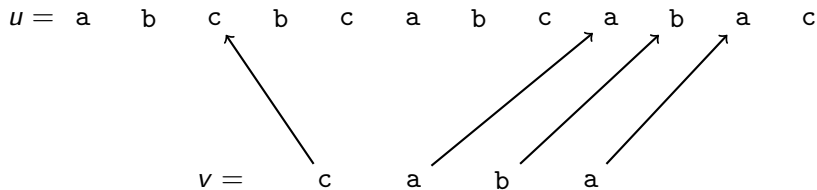
$$i \in \{1, \dots, |v| - 1\}$$



$$\text{gap}_e(u, 1) = bc$$

$$\text{gap}_e(u, 2) = \varepsilon$$

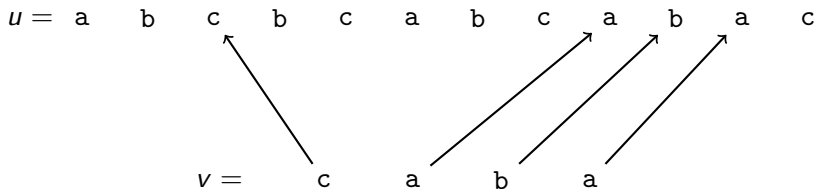
$$\text{gap}_e(u, 3) = c$$



$$\text{gap}_e(u, 1) = bcabc$$

$$\text{gap}_e(u, 2) = \varepsilon$$

$$\text{gap}_e(u, 3) = \varepsilon$$



$$\text{gap}_e(u, 1) = bcabc$$

$$\text{gap}_e(u, 2) = \varepsilon$$

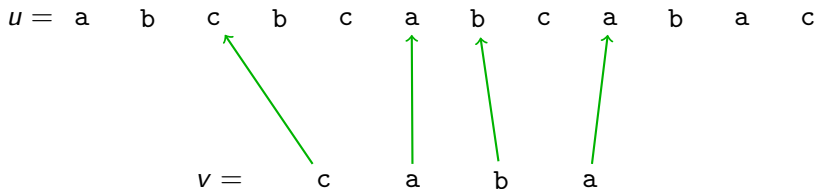
$$\text{gap}_e(u, 3) = \varepsilon$$

Gap constraints

$gc = (C_1, \dots, C_{|v|-1})$, where $C_i \subseteq \Sigma^*$ for every $i \in \{1, \dots, |v| - 1\}$.

(tuple of gap constraints)

The embedding e satisfies gc w.r.t. v , if, for every $i \in [|v| - 1]$, $\text{gap}_e(u, i) \in C_i$.



$$\text{gap}_e(u, 1) = bc$$

$$\text{gap}_e(u, 2) = \varepsilon$$

$$\text{gap}_e(u, 3) = c$$

$$C_1 = \{ab, ac, bc, bcabc\}$$

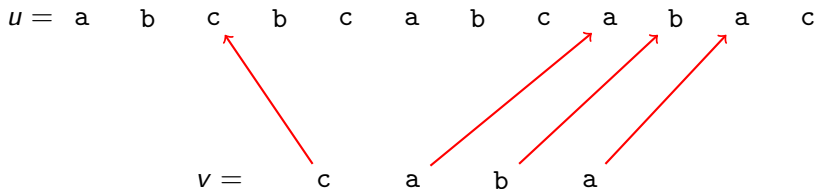
$$C_2 = \{\varepsilon, a, b\}$$

$$C_3 = \{a, b, c\}$$

Gap constraints

$gc = (C_1, \dots, C_{|v|-1})$, where $C_i \subseteq \Sigma^*$ for every $i \in \{1, \dots, |v| - 1\}$.
 (tuple of gap constraints)

The embedding e satisfies gc w.r.t. v , if, for every $i \in [|v| - 1]$, $\text{gap}_e(u, i) \in C_i$.



$$\text{gap}_e(u, 1) = bcabc$$

$$\text{gap}_e(u, 2) = \varepsilon$$

$$\text{gap}_e(u, 3) = \varepsilon$$

$$C_1 = \{ab, ac, bc, bcabc\}$$

$$C_2 = \{\varepsilon, a, b\}$$

$$C_3 = \{a, b, c\}$$

Gap constraints

$gc = (C_1, \dots, C_{|v|-1})$, where $C_i \subseteq \Sigma^*$ for every $i \in \{1, \dots, |v| - 1\}$.
(tuple of gap constraints)

The embedding e satisfies gc w.r.t. v , if, for every $i \in [|v| - 1]$, $\text{gap}_e(u, i) \in C_i$.

Notations

$v \preceq_{gc} u$: $v \preceq_e u$ for some embedding e satisfying gc .
(v is a gc -subsequence of u)

Notations

$v \preceq_{gc} u$: $v \preceq_e u$ for some embedding e satisfying gc .
(v is a gc -subsequence of u)

$$Subseq(gc, u) = \{v \in \Sigma^{|gc|+1} \mid v \preceq_{gc} u\}.$$

Gap Constraints

Notations

$v \preceq_{gc} u$: $v \preceq_e u$ for some embedding e satisfying gc .
(v is a gc -subsequence of u)

$$Subseq(gc, u) = \{v \in \Sigma^{|gc|+1} \mid v \preceq_{gc} u\}.$$

Example

$gc = (C_1, C_2)$ with $C_1 = C_2 = \{a, b, c\}$

$$Subseq(gc, abbacbab) = \{abc, bab, bca, abb\}.$$

Problems for Gap Constrained Subsequences

Matching

Input: $u, v \in \Sigma^*$, $k := |v|$, $(k - 1)$ -tuple gc of gap constraints.

Question: $v \preceq_{gc} u$?

Problems for Gap Constrained Subsequences

Matching

Input: $u, v \in \Sigma^*$, $k := |v|$, $(k - 1)$ -tuple gc of gap constraints.

Question: $v \preceq_{gc} u$?

Analysis Problems

Input: $u, v \in \Sigma^*$, $(k - 1)$ -tuple gc of gap constraints.

Questions: $Subseq(gc, u) = Subseq(gc, v)$? (Equivalence)

$Subseq(gc, u) \subseteq Subseq(gc, v)$? (Containment)

$Subseq(gc, u) = \Sigma^k$? (Universality)

Problems for Gap Constrained Subsequences

Matching

Input: $u, v \in \Sigma^*$, $k := |v|$, $(k - 1)$ -tuple gc of gap constraints.

Question: $v \preceq_{gc} u$?

Analysis Problems

Input: $u, v \in \Sigma^*$, $(k - 1)$ -tuple gc of gap constraints.

Questions: $Subseq(gc, u) = Subseq(gc, v)$? (Equivalence)

$Subseq(gc, u) \subseteq Subseq(gc, v)$? (Containment)

$Subseq(gc, u) = \Sigma^k$? (Universality)

Remark

k always means the length of subsequences!

(I. e., $gc = (C_1, \dots, C_{k-1})$.)

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.

... *length constraints* if every $C_i = \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$ for some $L^-(i), L^+(i) \in \mathbb{N} \cup \{0, +\infty\}$.

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.

... *length constraints* if every $C_i = \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$ for some $L^-(i), L^+(i) \in \mathbb{N} \cup \{0, +\infty\}$.

... *reg-len constraints* if every $C_i = C'_i \cap \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$.

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.
(Represented as: DFAs.)

... *length constraints* if every $C_i = \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$ for some $L^-(i), L^+(i) \in \mathbb{N} \cup \{0, +\infty\}$.

... *reg-len constraints* if every $C_i = C'_i \cap \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$.

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.

(Represented as: DFAs.)

... *length constraints* if every $C_i = \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$ for some $L^-(i), L^+(i) \in \mathbb{N} \cup \{0, +\infty\}$.

(Represented as: $(L^-(i), L^+(i)) \in \mathbb{N}^2$ in binary encoding.)

... *reg-len constraints* if every $C_i = C'_i \cap \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$.

Types of Gap-Constraints

Gap constraints $gc = (C_1, \dots, C_{k-1})$ are ...

... *regular constraints* if every C_i is a regular language.

(Represented as: DFAs.)

... *length constraints* if every $C_i = \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$ for some $L^-(i), L^+(i) \in \mathbb{N} \cup \{0, +\infty\}$.

(Represented as: $(L^-(i), L^+(i)) \in \mathbb{N}^2$ in binary encoding.)

... *reg-len constraints* if every $C_i = C'_i \cap \{v \in \Sigma^* \mid L^-(i) \leq |v| \leq L^+(i)\}$.

(Represented as: $((L^-(i), L^+(i)), A'_i)$, where A'_i is DFA for C'_i .)

Research Questions

Main Research Question

Investigate the complexity of matching and analysis problems in the presence of gap constraints.

Research Questions

Main Research Question

Investigate the complexity of matching and analysis problems in the presence of gap constraints.

Conditional Lower Bound Hypotheses

- ▶ *Exponential Time Hypothesis* (ETH):
No $2^{o(n)}$ $\text{poly}(n + m)$ algo. for 3-Satisfiability.
- ▶ *Strong Exponential Time Hypothesis* (SETH):
 $\forall \epsilon > 0 \exists k$: No $O(2^{n(1-\epsilon)}) \text{poly}(n)$ algo. for k -Satisfiability.
- ▶ *Orthogonal Vectors Hypothesis* (OVH):
 $\forall \epsilon > 0$: No $O(n^{2-\epsilon} \text{poly}(d))$ algo. for OV.

Complexity Bounds for Matching

Complexity Bounds for Matching

Upper Bound

The matching problem can be solved in time $O(|u||gc|)$.

Complexity Bounds for Matching

Upper Bound

The matching problem can be solved in time $O(|u||gc|)$.

Remark

For gap-constraints gc , let $|gc|$ be

- ▶ the total number of the DFAs states, if gc are regular constraints or reg-len constraints.
- ▶ the number of constraints with $L^+(i) - L^-(i) > 0$, if gc are length constraints.

Complexity Bounds for Matching

Conditional Lower Bound

The matching problem cannot be solved in time $O(|u|^h |gc|^g)$ with $g + h = 2 - \epsilon$ for some $\epsilon > 0$ unless OVH fails.

Complexity Bounds for Matching

Conditional Lower Bound

The matching problem cannot be solved in time $O(|u|^h |gc|^g)$ with $g + h = 2 - \epsilon$ for some $\epsilon > 0$ unless OVH fails.

For length constraints: this even holds for $|\Sigma| = 4$ and length constraints $(0, \ell)$ with $\ell \leq 6$.

For regular constraints: this even holds for $|\Sigma| = 4$ and all regular constraints are expressed by constant size DFAs.

Complexity Bounds for Matching

Conditional Lower Bound

The matching problem cannot be solved in time $O(|u|^h |gc|^g)$ with $g + h = 2 - \epsilon$ for some $\epsilon > 0$ unless OVH fails.

For length constraints: this even holds for $|\Sigma| = 4$ and length constraints $(0, \ell)$ with $\ell \leq 6$.

For regular constraints: this even holds for $|\Sigma| = 4$ and all regular constraints are expressed by constant size DFAs.

Proof Sketch

Fine-grained reduction from the *orthogonal vectors* problem (OV).

Complexity Bounds for Analysis Problems

Complexity Bounds for Analysis Problems

Remark

We only discuss the **non-universality** problem with length constraints.

→ We are checking $\text{Subseq}(gc, u) \neq \Sigma^k$.

All results hold analogously for non-equivalence and non-containment.

Complexity Bounds for Analysis Problems

Remark

We only discuss the **non-universality** problem with length constraints.

→ We are checking $\text{Subseq}(gc, u) \neq \Sigma^k$.

All results hold analogously for non-equivalence and non-containment.

Upper Bound

The non-universality problem can be solved in time $O(|\Sigma|^k |gc| |u|)$.

(Achieved by a brute force algorithm.)

Complexity Bounds for Analysis Problems

Conditional Lower Bounds

For every fixed alphabet Σ with $|\Sigma| \geq 3$, the non-universality problem with length constraints can be solved in time $2^{O(k)}|gc||u|$. Moreover, it cannot be solved

- ▶ in subexponential time $2^{o(k)} \text{poly}(|u|, k)$ (unless ETH fails),
- ▶ in time $O(2^{k(1-\epsilon)} \text{poly}(|u|, k))$ (unless SETH fails).

These lower bounds hold even if all length constraints are $(1, 5)$.

Complexity Bounds for Analysis Problems

Conditional Lower Bounds

For every fixed alphabet Σ with $|\Sigma| \geq 3$, the non-universality problem with length constraints can be solved in time $2^{O(k)}|gc||u|$. Moreover, it cannot be solved

- ▶ in subexponential time $2^{o(k)} \text{poly}(|u|, k)$ (unless ETH fails),
- ▶ in time $O(2^{k(1-\epsilon)} \text{poly}(|u|, k))$ (unless SETH fails).

These lower bounds hold even if all length constraints are $(1, 5)$.

Theorem

For every fixed alphabet Σ with $|\Sigma| = 2$, non-universality with length constraints is NP-complete even if each length constraint is $(0, 0)$ or $(3, 9)$.

Proof

Reduction from CNF-Sat.

Quick Overview

	Classical	Gap Constraints
Matching	$O(n)$	$O(u gc)$
Equivalence	$O(n)$	
Containment	$O(nm \Sigma)$	$O(\Sigma ^k gc u)$
Universality	$O(n)$	

Additional Results

Gap Length Equalities

In addition to gap constraints, we consider gap length equalities of the form “ $|\text{gap}_i| = |\text{gap}_j|$ ” meaning $|\text{gap}_e(w, i)| = |\text{gap}_e(w, j)|$.

Gap Length Equalities

In addition to gap constraints, we consider gap length equalities of the form “ $|\text{gap}_i| = |\text{gap}_j|$ ” meaning $|\text{gap}_e(w, i)| = |\text{gap}_e(w, j)|$.

(Or even more complex variant: $2|\text{gap}_7| + |\text{gap}_3| \leq |\text{gap}_2|$.)

Additional Results

Gap Length Equalities

In addition to gap constraints, we consider gap length equalities of the form “ $|\text{gap}_i| = |\text{gap}_j|$ ” meaning $|\text{gap}_e(w, i)| = |\text{gap}_e(w, j)|$.

(Or even more complex variant: $2|\text{gap}_7| + |\text{gap}_3| \leq |\text{gap}_2|$.)

Theorem

Matching with gap constraints and gap length equalities is NP-complete even for $|\Sigma| = 2$ and gap constraints $C_i = \Sigma^*$.

Proof

Reduction from 3-CNF-SAT.

Sets of Gap Constrained Subsequences with Multiplicities:

$$\textit{Subseq}(2, \text{abba}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

$$\textit{Subseq}(2, \text{abab}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

Sets of Gap Constrained Subsequences with Multiplicities:

$$\text{Subseq}(2, \text{abba}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

$$\text{Subseq}(2, \text{abab}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

$$\text{Subseq}^m(2, \text{abba}) = \{(\text{aa}, 1), (\text{ab}, 2), (\text{ba}, 2), (\text{bb}, 1)\}$$

$$\text{Subseq}^m(2, \text{abab}) = \{(\text{aa}, 1), (\text{ab}, 3), (\text{ba}, 1), (\text{bb}, 1)\}$$

Multiplicities

Sets of Gap Constrained Subsequences with Multiplicities:

$$\text{Subseq}(2, \text{abba}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

$$\text{Subseq}(2, \text{abab}) = \{\text{aa}, \text{ab}, \text{ba}, \text{bb}\}$$

$$\text{Subseq}^m(2, \text{abba}) = \{(\text{aa}, 1), (\text{ab}, 2), (\text{ba}, 2), (\text{bb}, 1)\}$$

$$\text{Subseq}^m(2, \text{abab}) = \{(\text{aa}, 1), (\text{ab}, 3), (\text{ba}, 1), (\text{bb}, 1)\}$$

Theorem

The equivalence problem for gap-constrained subsequences with multiplicities can be decided in polynomial time.

The presented results and topics are also summarized in the survey paper presented at NCMA, 2022:

Combinatorial Algorithms for Subsequence Matching: A Survey
(MK, Tore Koß, Florin Manea, Stefan Siemer)

The presented results and topics are also summarized in the survey paper presented at NCMA, 2022:

Combinatorial Algorithms for Subsequence Matching: A Survey
(MK, Tore Koß, Florin Manea, Stefan Siemer)

Thank you!

Theorem

The non-universality problem with length constraints cannot be solved in running time $O(f(k) \text{poly}(|w|, k))$ for any computable function f (unless $\text{FPT} = \text{W}[1]$).