

Algorithmen auf Sequenzen

Erweiterte Patternklassen

Dominik Kopczynski

Lehrstuhl für Algorithm Engineering (LS11)
Fakultät für Informatik
TU Dortmund

Übersicht

- Bisher wurden nur exakte Pattern besprochen.
- Es ist auch interessant Pattern zu betrachten, bei denen bestimmte Zeichen an einer Position beliebig sein dürfen (verallgemeinerte Strings).
- Eine weitere Erweiterung der Pattern wäre es Gaps beschränkter Länge zu betrachten.
- Auch kann es vorkommen, dass Zeichen im Pattern optional sein können.
- In der Summe könnte so nach "einfachen" regulären Ausdrücken gesucht werden.

Verallgemeinerte Strings

- Verallgemeinerte Strings sind Strings, die *Teilmengen* des Alphabets als Zeichen besitzen.
- Die Patternmenge {Meier, Meyer} könnte somit zu Me[iy]er verallgemeinert werden.
- Wenn ein(!) beliebiges Zeichen erlaubt ist, kann ein Pattern folgendermaßen dargestellt werden: image#.jpg.

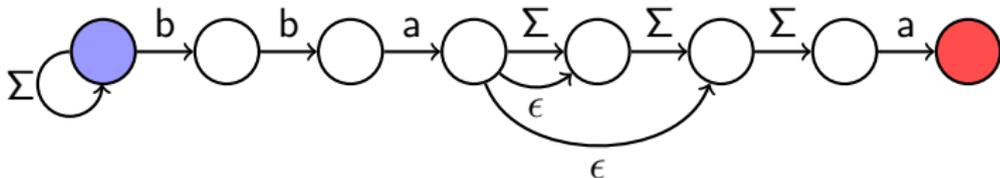
Verallgemeinerte Strings

- Verallgemeinerte Strings sind Strings, die *Teilmengen* des Alphabets als Zeichen besitzen.
- Die Patternmenge {Meier, Meyer} könnte somit zu Me[iy]er verallgemeinert werden.
- Wenn ein(!) beliebiges Zeichen erlaubt ist, kann ein Pattern folgendermaßen dargestellt werden: image#.jpg.
- Bitparallele Algorithmen lassen sich ohne viel Aufwand anpassen, indem beispielsweise beim Shift-And nur die Masken angepasst werden. Beispiel: $P = abba\#b$ mit $\Sigma = \{a, b\}$.

	b#abba
$mask^a$	011001 ₂
$mask^b$	110110 ₂

Gaps beschränkter Länge

- Unter Gaps beschränkter Länge wird eine nichtleere Folge beliebiger Zeichen verstanden.
- Ein Beispieldattern sei: $P = bba\#(1, 3)a$.
- Hierbei muss es 1 beliebiges Zeichen und bis zu $3 - 1 = 2$ optionale beliebige Zeichen geben.
- Es werden also zwei ϵ -Übergänge benötigt.



Gaps beschränkter Länge

Einschränkungen:

- 1 Die Gaps $\#(,)$ sind nicht an erster oder letzter Position im Pattern erlaubt.
- 2 Zwei (oder mehrere) Gaps hintereinander sind nicht erlaubt, denn $\#(u, v)\#(u', v') \hat{=} \#(u + u', v + v')$.
- 3 Für $\#(u, v)$ muss gelten: $1 \leq u \leq v$.

Implementierung

Für die Masken des Alphabets gilt: bei einem Gap $\#(u, v)$ werden für jedes Zeichen aus dem Alphabet genau v 1en hinzugefügt.

Beispiel mit $P = bba\#(1, 3)a$ und $\Sigma = \{a, b, c\}$:

	a###abb
$mask^a$	1111100 ₂
$mask^b$	0111011 ₂
$mask^c$	0111000 ₂

Implementierung

Zur Simulation der ϵ -Übergänge bringt die Subtraktion auf Bit-Ebene den gewünschten Effekt.

- Bitmaske I speichert für jedes Gap-Token die Position, von der ϵ -Kanten ausgehen. Angenommen an Position i im Pattern ist ein Gap-Token, so wird $I[i - 1] := 1$ gesetzt.

Implementierung

Zur Simulation der ϵ -Übergänge bringt die Subtraktion auf Bit-Ebene den gewünschten Effekt.

- Bitmaske I speichert für jedes Gap-Token die Position, von der ϵ -Kanten ausgehen. Angenommen an Position i im Pattern ist ein Gap-Token, so wird $I[i - 1] := 1$ gesetzt.
- Bitmaske F speichert für jedes Gap-Token die Position, die nach der letzten ϵ -Kante erreicht wird. Angenommen an Position i im Pattern ist ein Gap-Token, so wird $F[i + v - u + 1] := 1$ gesetzt.

a###abb

F 0100000₂

I 0000100₂

Implementierung

- Um festzustellen, ob sich ein aktiver Zustand am Beginn eines Gaps befindet, wird die aktive Zustandsmenge mit der Startbitmaske verundet: $A \& I$.
- Durch Subtraktion von F und $A \& I$ werden alle Bits zwischen $[u : v]$ auf 1 gesetzt.

$$\begin{array}{r}
 F \quad \quad 0100000_2 \\
 A \& I \quad 0000100_2 \\
 \hline
 - \quad \quad 0011100_2
 \end{array}$$

Implementierung

- Um festzustellen, ob sich ein aktiver Zustand am Beginn eines Gaps befindet, wird die aktive Zustandsmenge mit der Startbitmaske verundet: $A \& I$.
- Durch Subtraktion von F und $A \& I$ werden alle Bits zwischen $[u : v]$ auf 1 gesetzt.

$$\begin{array}{r}
 F \quad 0100000_2 \\
 A \& I \quad 0000100_2 \\
 \hline
 - \quad 0011100_2
 \end{array}$$

- Problem, wenn es mehrere Gaps gibt:

$$\begin{array}{r}
 F \quad 010000100000_2 \\
 A \& I \quad 000000000100_2 \\
 \hline
 - \quad 010000011100_2
 \end{array}$$

Implementierung

- Lösung: Durch die Negation von F wird eine Maske erzeugt, die genutzt werden kann, um ungewünschte 1en zu löschen.

$$\begin{array}{r}
 F \qquad \qquad \qquad 010000100000_2 \\
 A \ \& \ I \qquad \qquad 000000000100_2 \\
 \hline
 F - (A \ \& \ I) \qquad 010000011100_2
 \end{array}$$

Implementierung

- Lösung: Durch die Negation von F wird eine Maske erzeugt, die genutzt werden kann, um ungewünschte 1en zu löschen.

F	010000100000 ₂
$A \& I$	000000000100 ₂
$F - (A \& I)$	010000011100 ₂
$\sim F$	101111011111 ₂
$(F - (A \& I)) \& \sim F$	000000011100 ₂

Implementierung

- Lösung: Durch die Negation von F wird eine Maske erzeugt, die genutzt werden kann, um ungewünschte 1en zu löschen.

$$\begin{array}{r}
 F \qquad \qquad \qquad 010000100000_2 \\
 A \ \& \ I \qquad \qquad \qquad 000000000100_2 \\
 \hline
 F - (A \ \& \ I) \qquad \qquad \qquad 010000011100_2 \\
 \sim F \qquad \qquad \qquad 101111011111_2 \\
 \hline
 (F - (A \ \& \ I)) \ \& \ \sim F \qquad \qquad \qquad 000000011100_2
 \end{array}$$

- Wegen der Einschränkung, dass zwei Gaps nicht unmittelbar aufeinander folgen dürfen, ist kein F -Zustand gleichzeitig ein I -Zustand des nächsten Gaps.

Implementierung

Erweiterte Update-Funktion:

- 1 **Herkömmliche Shift-And-Update-Funktion anwenden:**

$$A = ((A \ll 1) | 1) \& \text{mask}[c]$$

- 2 **Zustandsvektor um neue aktive Zustände erweitern:**

$$A = A | ((F - (A \& I)) \& \sim F)$$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bbacca}$:

	$a\#\#\text{abb}$		
$mask^a$	1111100_2	F	0100000_2
$mask^b$	0111011_2	$(A \& I)$	0000000_2
$mask^c$	0111000_2	$=$	0000000_2

$$A = 0000000$$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba\#(1,3)a}$ und $T = \text{bbacca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit b: $A = 0000001$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba\#(1,3)a}$ und $T = \text{bbacca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 2: $A = 0000001$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{b}^{\text{red}}\text{acca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit b: $A = 0000011$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{b}^{\text{red}}\text{acca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	$=$	0000000 ₂

Update 2: $A = 0000011$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bb}\text{a}\text{cca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit a: $A = 0000100$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bb}\text{a}\text{cca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000100 ₂
$mask^c$	0111000 ₂	=	0011100 ₂

Update 2: $A = 0011100$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bba}\text{c}\text{ca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit c: $A = 0111000$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bba}\text{c}\text{ca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 2: $A = 0111000$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bbac}\text{ca}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit c: $A = 0110000$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bbac}\text{ca}$:

	a###abb		
$mask^a$	1111100_2	F	0100000_2
$mask^b$	0111011_2	$(A \& I)$	0000000_2
$mask^c$	0111000_2	$=$	0000000_2

Update 2: $A = 0110000$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bbacc}\text{a}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	=	0000000 ₂

Update 1 mit a: $A = 1100000$

Beispiel für den Gap-Shift-And

Sei $P = \text{bba}\#(1, 3)\text{a}$ und $T = \text{bbacc}\text{a}$:

	a###abb		
$mask^a$	1111100 ₂	F	0100000 ₂
$mask^b$	0111011 ₂	$(A \& I)$	0000000 ₂
$mask^c$	0111000 ₂	$=$	0000000 ₂

Update 2: $A = 1100000$ Treffer

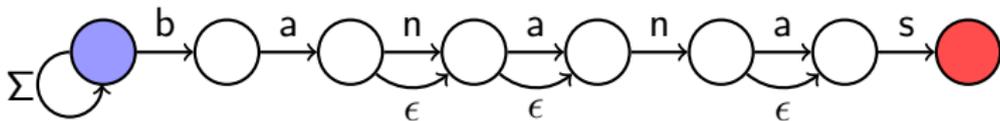
Optionale Zeichen

- Mit einer zusätzlichen Erweiterung können Zeichen verarbeitet werden, die bis zu einmal vorkommen dürfen.
- Eine Notation für die Patternmenge $\{\text{color, colour}\}$ wäre beispielsweise: $P = \text{colou?r}$.
- Konsekutive ϵ -Transitionen (Blöcke) sind erlaubt.
- Auch diese Erweiterung lässt sich mit dem Shift-And-Algorithmus realisieren.

Optionale Zeichen

Implementierung:

Für die folgenden Berechnungen sind drei zusätzliche Masken nötig. Am Beispiel-NFA für das Pattern $P = \text{ban? a? na? s}$ werden die Masken verdeutlicht:



$I :$	0	1	0	0	1	0	0
$F :$	0	0	0	1	0	1	0
$O :$	0	0	1	1	0	1	0

Optionale Zeichen

Implementierung:

- Sobald ein Zustand innerhalb eines Blocks aktiv ist, sind seine Folgezustände (innerhalb des Blocks) auch aktiv.
- Wie findet man den ersten aktiven Zustand?

Optionale Zeichen

Implementierung:

- Sobald ein Zustand innerhalb eines Blocks aktiv ist, sind seine Folgezustände (innerhalb des Blocks) auch aktiv.
- Wie findet man den ersten aktiven Zustand?
- Subtraktion wieder ausnutzen: wenn bei einer Binärzahl um 1 subtrahiert wird, werden alle niederwertigen Bits bis einschließlich der ersten 1 geflippt.
- Alle höherwertigen Bits bleiben erhalten.

Optionale Zeichen

Implementierung:

- Sobald ein Zustand innerhalb eines Blocks aktiv ist, sind seine Folgezustände (innerhalb des Blocks) auch aktiv.
- Wie findet man den ersten aktiven Zustand?
- Subtraktion wieder ausnutzen: wenn bei einer Binärzahl um 1 subtrahiert wird, werden alle niederwertigen Bits bis einschließlich der ersten 1 geflippt.
- Alle höherwertigen Bits bleiben erhalten.

$$\begin{array}{r}
 1101010000 \\
 - \qquad \qquad \qquad 1 \\
 \hline
 1101001111
 \end{array}$$

$$\begin{array}{r}
 1101011000 \\
 - \qquad \qquad \qquad 100 \\
 \hline
 1101010100
 \end{array}$$

Optionale Zeichen

Implementierung:

- Damit der Minuend immer größer als der Subtrahend ist, wird die aktive Zustandsmenge mit den Endpositionen der Blöcke verodert, sei $A_f := A \mid F$.

Optionale Zeichen

Implementierung:

- Damit der Minuend immer größer als der Subtrahend ist, wird die aktive Zustandsmenge mit den Endpositionen der Blöcke verodert, sei $A_f := A \mid F$.
- Mit der Subtraktion werden die niederwertigsten Bits geflippt: $A_f - I$.

Optionale Zeichen

Implementierung:

- Damit der Minuend immer größer als der Subtrahend ist, wird die aktive Zustandsmenge mit den Endpositionen der Blöcke verodert, sei $A_f := A \mid F$.
- Mit der Subtraktion werden die niederwertigsten Bits geflippt: $A_f - I$.
- Da es keine Maschinenoperation für bitweise Äquivalenz gibt, wird die exklusive Veroderung genutzt: $a \equiv b \hat{=} \sim a \oplus b$.

Optionale Zeichen

Implementierung:

- Damit der Minuend immer größer als der Subtrahend ist, wird die aktive Zustandsmenge mit den Endpositionen der Blöcke verodert, sei $A_f := A \mid F$.
- Mit der Subtraktion werden die niederwertigsten Bits geflippt: $A_f - I$.
- Da es keine Maschinenoperation für bitweise Äquivalenz gibt, wird die exklusive Verodung genutzt: $a \equiv b \hat{=} \sim a \oplus b$.
- Seien also $\sim(A_f - I) \oplus A_f$ die äquivalenten Bits.

Optionale Zeichen

Implementierung:

- Damit der Minuend immer größer als der Subtrahend ist, wird die aktive Zustandsmenge mit den Endpositionen der Blöcke verodert, sei $A_f := A \mid F$.
- Mit der Subtraktion werden die niederwertigsten Bits geflippt: $A_f - I$.
- Da es keine Maschinenoperation für bitweise Äquivalenz gibt, wird die exklusive Veroderung genutzt: $a \equiv b \hat{=} \sim a \oplus b$.
- Seien also $\sim(A_f - I) \oplus A_f$ die äquivalenten Bits.
- Zuletzt sollen noch alle Bits auf 0 gesetzt werden, die sich nicht im Block befinden, mit $0 \& (\sim(A_f - I) \oplus A_f)$.

Optionale Zeichen

Implementierung:

- 1 Masken für alle Zeichen aus Σ erstellen, dabei die optionalen Zeichen wie reguläre Zeichen behandeln.
- 2 Herkömmliche Shift-And-Update-Funktion anwenden:
 $A = ((A \ll 1) | 1) \& \text{mask}[c]$
- 3 Zustandsvektor um neue aktive Zustände erweitern:

$$A_f = A | F$$

$$A = A | (0 \& (\sim(A_f - I) \wedge A_f))$$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sanab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂	<hr/>	
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

$$A = 0000000$$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		<hr/>
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 1 mit b: $A = 0000001$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		0000000 ₂
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 2: $A = 0000001$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sanab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		0000000 ₂
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 1 mit a: $A = 0000010$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		<hr/>
		$O \& (\sim(A_f - I) \oplus A_f)$	0001100 ₂

Update 2: $A = 0001110$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sanab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂	<hr/>	
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 1 mit n : $A = 0010100$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		0101100 ₂
		$O \& (\sim(A_f - I) \oplus A_f)$	0101100 ₂

Update 2: $A = 0111100$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂	<hr/>	
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 1 mit n: $A = 0010000$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{banns}$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂		0100000 ₂
		$O \& (\sim(A_f - I) \oplus A_f)$	

Update 2: $A = 0110000$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{bann}s$:

	sananab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂	<hr/>	
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 1 mit s : $A = 1000000$

Beispiel für den Optional-Shift-And

Sei $P = \text{ban?a?na?s}$ und $T = \text{bann}s$:

	sanab		
$mask^a$	0101010 ₂	I	0010010 ₂
$mask^b$	0000001 ₂	F	0101000 ₂
$mask^n$	0010100 ₂	O	0101100 ₂
$mask^s$	1000000 ₂	<hr/>	
		$O \& (\sim(A_f - I) \oplus A_f)$	0000000 ₂

Update 2: $A = 1000000$ Treffer

Zusammenfassung

- Durch einfache Tricks lässt sich der Shift-And-Algorithmus auf erweiterte Patternklassen ausbauen.
- Verallgemeinerte Strings, Gaps beschränkter Länge und optionale Zeichen sind möglich.
- Die Laufzeit für die Erstellung der Masken für Wildcards und Gaps beschränkter Länge beträgt $\mathcal{O}(m|\Sigma|)$.

Zusammenfassung

- Durch einfache Tricks lässt sich der Shift-And-Algorithmus auf erweiterte Patternklassen ausbauen.
- Verallgemeinerte Strings, Gaps beschränkter Länge und optionale Zeichen sind möglich.
- Die Laufzeit für die Erstellung der Masken für Wildcards und Gaps beschränkter Länge beträgt $\mathcal{O}(m|\Sigma|)$.
- Die Erstellung der Masken für optionale Zeichen bleibt unverändert, die Laufzeit beträgt also $\mathcal{O}(m)$.
- Der Speicherplatz beträgt weiterhin $\mathcal{O}(|\Sigma|\lceil m/W \rceil)$.
- Die Laufzeit für die Mustersuche bleibt unverändert bei $\mathcal{O}(n\lceil m/W \rceil)$.