# Revitalizing Research on Approximate String Matching Algorithms
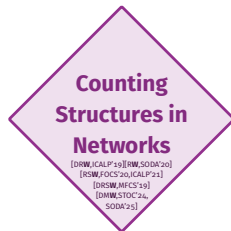
Philip Wellnitz

National Institute of Informatics

**Approximate String Matching**
[BKW,SODA'19]
[CKW,FOCS'20;'22;'25]
[CKW,FOCS'23]
[NKW,STOC'24,
SODA'25]

**Counting Structures in Networks**
[DRW,ICALP'19][RW,SODA'20]
[RSW,FOCS'20,ICALP'21]
[DRSW,MFCS'19]
[DMW,STOC'24,
SODA'25]

**Dense Subset Sum**
[BW,SODA'20]

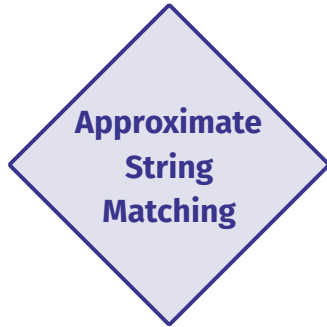**Scheduling Problems**
[BFHSW,ICALP'20,
Algorithmica'22]

**Generalized Dominating Set**
[FMMNSSW,SODA'23,
ToCT'25,TALG'25]
[GSW,STACS'25]

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    2-1

**Approximate String Matching**

...

(TU Kaiserslautern).

(Simons Institute and UC Berkeley), (University of Salzburg), (University of Salzburg).

(University of Pennsylvania), (University of Pennsylvania), (University of Pennsylvania).

(Reichman University, Herzliya, Israel and Birkbeck, University of London), (UC Berkeley and Max Planck Institute for Informatics, SIC, Saarbrücken, Germany), (Max Planck Institute for Informatics, SIC, Saarbrücken, Germany).

(National Institute of Informatics).

(Université Paris Cité, CNRS, IRIF, F-75013, Paris, France); (Universitat Politècnica de Catalunya)

...

...

(ETH Zurich); (Toyota Technological Institute at Chicago)

(Merton College, University of Oxford, United Kingdom); (Mathematical Institute, University of Bonn, Germany); (Max Planck Institute for Informatics, Saarland Informatics Campus (SIC), Saarbrücken, Germany)

(The Academic College of Tel Aviv-Yafo), (UC Berkeley), (Weizmann Institute of Science), (University of California San Diego).

(CISPA Helmholtz Center for Information Security, Saarbrücken); (Simon Fraser University); (Duke University); (CISPA Helmholtz Center for Information Security, Saarbrücken); (Max Planck Institute for Informatics, SIC, Saarbrücken)

...

...

(Alibaba Group); (Microsoft Research); (Washington University in St. Louis); (Columbia University); (University of Chicago)

(Stony Brook University); (Max Planck Institute for Informatics, Saarbruecken)

(Univ. Bordeaux, CNRS, Bordeaux INP, LaBRI, France); (University of Warsaw, Poland); (Saarland University and Max Planck Institute for Informatics, Saarbrücken, Germany); (University of Warsaw, Poland)

(ETH Zürich); (TU Berlin)

(Carnegie Mellon University)

...

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms
4-1

# An Example

...

(TU Kaiserslautern).

(Simons Institute and UC Berkeley), (University of Salzburg), (University of Salzburg).

(University of Pennsylvania), (University of Pennsylvania), (University of Pennsylvania).

(Reichman University, Herzliya, Israel and Birkbeck, University of London), (UC Berkeley and Max Planck Institute for Informatics, SIC, **Saarbrücken**, Germany), (Max Planck Institute for Informatics, SIC, **Saarbrücken**, Germany).

(National Institute of Informatics).

(Université Paris Cité, CNRS, IRIF, F-75013, Paris, France); (Universitat Politècnica de Catalunya)

...

...

(ETH Zurich); (Toyota Technological Institute at Chicago)

(Merton College, University of Oxford, United Kingdom); (Mathematical Institute, University of Bonn, Germany); (Max Planck Institute for Informatics, Saarland Informatics Campus (SIC), Saarbr|cken, Germany)

(The Academic College of Tel Aviv-Yafo), (UC Berkeley), (Weizmann Institute of Science), (University of California San Diego).

(CISPA Helmholtz Center for Information Security, **Saarbrücken**); (Simon Fraser University); (Duke University); (CISPA Helmholtz Center for Information Security, **Saarbrücken**); (Max Planck Institute for Informatics, SIC, **Saarbrücken**)

...

...

(Alibaba Group); (Microsoft Research); (Washington University in St. Louis); (Columbia University); (University of Chicago)

(Stony Brook University); (Max Planck Institute for Informatics, Saarbruecken)

(Univ. Bordeaux, CNRS, Bordeaux INP, LaBRI, France); (University of Warsaw, Poland); (Saarland University and Max Planck Institute for Informatics, Saarbrücken, Germany); (University of Warsaw, Poland)

(ETH Zürich); (TU Berlin)

(Carnegie Mellon University)

...

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 4-2

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

···

(TU Kaiserslautern).

(Simons Institute and UC Berkeley), (University of Salzburg), (University of Salzburg).

(University of Pennsylvania), (University of Pennsylvania), (University of Pennsylvania).

(Reichman University, Herzliya, Israel and Birkbeck, University of London), (UC Berkeley and Max Planck Institute for Informatics, SIC, **Saarbrücken**, Germany), (Max Planck Institute for Informatics, SIC, **Saarbrücken**, Germany).

(National Institute of Informatics).

(Université Paris Cité, CNRS, IRIF, F-75013, Paris, France); (Universitat Politècnica de Catalunya)

···

···

(ETH Zurich); (Toyota Technological Institute at Chicago)

(Merton College, University of Oxford, United Kingdom); (Mathematical Institute, University of Bonn, Germany); (Max Planck Institute for Informatics, Saarland Informatics Campus (SIC), **Saarbr|cken**, Germany)

(The Academic College of Tel Aviv-Yafo), (UC Berkeley), (Weizmann Institute of Science), (University of California San Diego).

(CISPA Helmholtz Center for Information Security, **Saarbrücken**); (Simon Fraser University); (Duke University); (CISPA Helmholtz Center for Information Security, **Saarbrücken**); (Max Planck Institute for Informatics, SIC, **Saarbrücken**)

···

···

(Alibaba Group); (Microsoft Research; St. Louis); (Washington University in (Columbia University); (University of Chicago)

(Stony Brook University); (Max Planck Institute for Informatics, **Saarbruecken**)

(Univ. Bordeaux, CNRS, Bordeaux INP, LaBRI, France); (University of Warsaw, Poland); (Saarland University and Max Planck Institute for Informatics, **Saarbrücken**, Germany); (University of Warsaw, Poland)

(ETH Zürich); (TU Berlin)

(Carnegie Mellon University)

···

Task: Find **Saarbrücken** in a text.

Task: Find **Saarbrücken** in a text.

Or Saarbruecken.

Task: Find **Saarbrücken** in a text.

Or Saarbruecken. Or Sarrebruck.

Task: Find **Saarbrücken** in a text.

Or Saarbruecken. Or Sarrebruck. Or Saarbrucken, Saarbr|cken, SaarbrÃ¼cken, ….

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 5-4

> **Approximate String Matching** >

early 1980's

Given a text *T*, a pattern *P*, and an integer *k*, identify the (starting positions of) substrings of *T* that are at **edit distance** of at most *k* to *P*.

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    6-1

> **Approximate String Matching**

early 1980's

Given a text *T*, a pattern *P*, and an integer *k*, identify the (starting positions of) substrings of *T* that are at **edit distance** of at most *k* to *P*.

**Edit distance:** minimum number of insertions, deletions, or substitutions of single characters to transform one string into another string

Basic Tricks and Tools, Previous Work

New Algorithms via Structural Insights

Spotlight Extension: Quantum Algorithms for ASM

Spotlight Extension: String Matching with Weighted Edits

Open Problems and Future Directions

**Approximate String Matching**

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.



Focus: Obtain starting positions of occurrences

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   8-1

⟨ **Approximate String Matching** ⟩ early 1980's

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences



Philip Wellnitz

Revitalizing Research on Approximate String Matching Algorithms   8-2

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Approximate String Matching**      early 1980's

Given: text *T*, pattern *P*, threshold *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.



Focus: Obtain starting positions of occurrences

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   8-3

⟨ **Approximate String Matching** ⟩ early 1980's

Given: text *T*, pattern *P*, threshold *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   8-4

⟨ **Approximate String Matching** ⟩          early 1980's

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

"Standard Trick":        write $t := |T|$, $p := |P|$
Split $T$ into overlapping fragments of len $\ell + p + k$

⤳ $O(t/\ell)$ instances,
    each "responsible" for its first $\ell$ positions

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms          8-5

**Approximate String Matching** early 1980's

Given: text *T*, pattern *P*, threshold *k*; Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Standard Trick":    write $t := |T|$, $p := |P|$
Split *T* into overlapping fragments of len $\ell + p + k$

⤳ $O(t/\ell)$ instances,
   each "responsible" for its first $\ell$ positions

⤳ Useful special cases: $\ell = k$ and $\ell = 0.5\,p$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    8-6

**Approximate String Matching**

early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.



Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   9-1

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

⟨ **Approximate String Matching** ⟩    early 1980's

Given: text $T$, pattern $P$, integer $k$;    Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:    (typically fast)
Compute (small) superset of starting positions



Filter

**Approximate String Matching**

early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

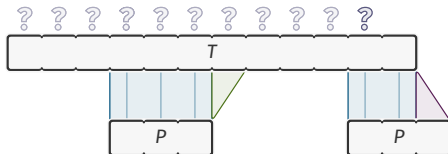Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:     (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:     (typically slow)
Check for occ at each remaining position



Filter

Verify

**Approximate String Matching**

early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:     (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:     (typically slow)
Check for occ at each remaining position

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

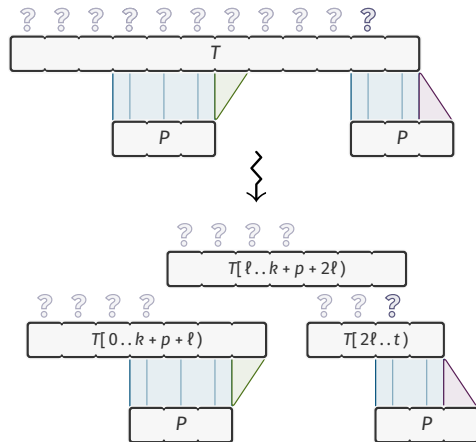**Approximate String Matching**     early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:     (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:     (typically slow)
Check for occ at each remaining position

**Approximate String Matching**

early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:      (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:      (typically slow)
Check for occ at each remaining position

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms      9-6
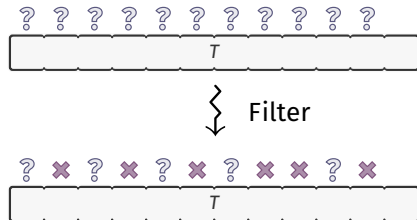
**Approximate String Matching**  early 1980's

Given: text $T$, pattern $P$, integer $k$;  Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:  (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:  (typically slow)
Check for occ at each remaining position

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

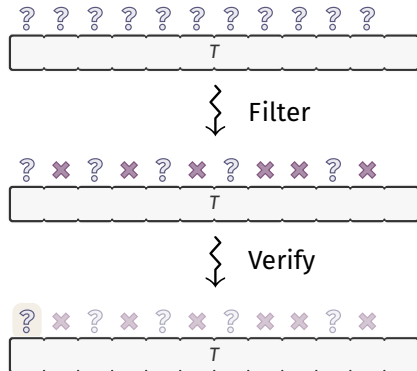**Approximate String Matching**                    early 1980's

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:     (typically fast)
Compute (small) superset of starting positions

Step 2, Verify:     (typically slow)
Check for occ at each remaining position

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    9-8

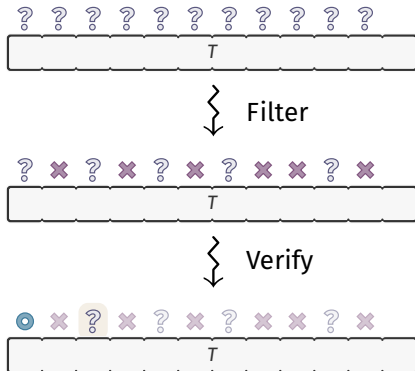**Approximate String Matching**                    early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:      (typically fast)
Compute (small) superset of starting positions

Step 2', Multi-pos Verify:      (typically slow)
Check for occ at each remaining position,
multiple positions at once

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics
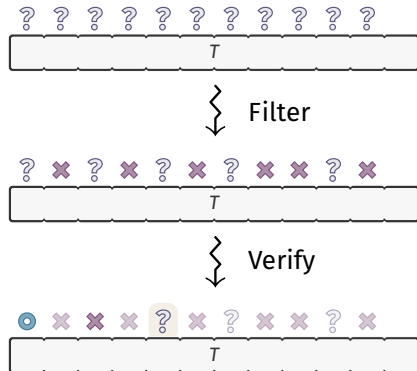
**Approximate String Matching**　　　early 1980's

Given: text *T*, pattern *P*, integer *k*;　Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:　(typically fast)
Compute (small) superset of starting positions

Step 2', Multi-pos Verify:　(typically slow)
Check for occ at each remaining position,
multiple positions at once

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms　9-10

◆◆◆◆◆◆◆◆◆◆◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇

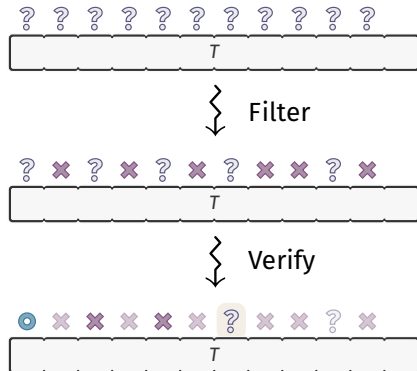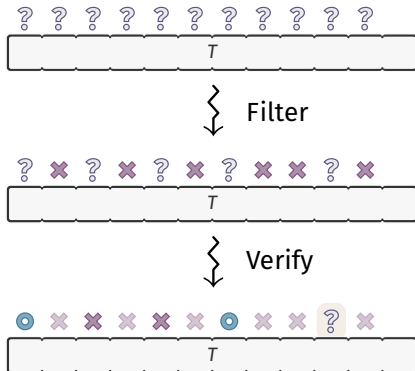**Approximate String Matching**                    early 1980's

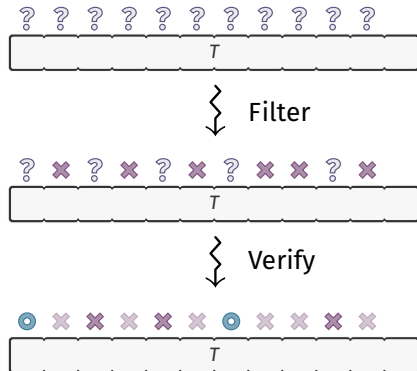Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Focus: Obtain starting positions of occurrences

"Filter and Verify Paradigm"

Step 1, Filter:        (typically fast)
Compute (small) superset of starting positions

Step 2', Multi-pos Verify:        (typically slow)
Check for occ at each remaining position,
multiple positions at once

## Approximate String Matching

early 1980's

Given: text *T*, pattern *P*, threshold *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

| Edit dist/Verify Algorithm | ASM Algorithm | ($t := |T|, p := |P|$) |
|---|---|---|
| Textbook DP ⇝ $O(tp)$ <br> [Sellers, 1980] and others | Verify pos 1 by 1 ⇝ $O(t^2p)$ | substr |
| | | |
| | | |



$e_{i,0} = i$ (delete $T[0..i]$)

$e_{0,j} = j$ (insert $P[0..j]$)

$$e_{i,j} = \min \begin{cases} e_{i-1,j} + 1 & \text{(del from } T\text{)} \\ e_{i,j-1} + 1 & \text{(ins in } T\text{)} \\ e_{i-1,j-1} \\ \quad +[T[i] \neq P[j]] & \text{(match/} \\ & \text{subst)} \end{cases}$$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   10-1

## Approximate String Matching

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

| Edit dist/Verify Algorithm | ASM Algorithm | ($t := |T|, p := |P|$) |
|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ | Verify pos 1 by 1 $\rightsquigarrow O(t^2 p)$ | substr |
| [Sellers, 1980] and others | Verify all pos at once $\rightsquigarrow O(tp)$ | start pos |



$e_{i,0} = 0$ (delete $T[\,0..i\,]$)
$e_{0,j} = j$ (insert $P[\,0..j\,]$)

$$e_{i,j} = \min \begin{cases} e_{i-1,j} + 1 & \text{(del from } T\text{)} \\ e_{i,j-1} + 1 & \text{(ins in } T\text{)} \\ e_{i-1,j-1} & \\ \quad +[T[\,i\,] \neq P[\,j\,]] & \text{(match/} \\ & \text{subst)} \end{cases}$$

**Approximate String Matching**

Given: text *T*, pattern *P*, threshold *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

| Edit dist/Verify Algorithm | ASM Algorithm | ($t := |T|, p := |P|$) |
|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ <br> [Sellers, 1980] and others | Verify pos 1 by 1 $\rightsquigarrow O(t^2p)$ <br> Verify all pos at once $\rightsquigarrow O(tp)$ | substr <br> start pos |
| DP, diagonally compute only entries ≤ *k* <br> $\rightsquigarrow O(t + kp)$ | Verify pos 1 by 1 $\rightsquigarrow O(tkp)$ <br> Verify $\ell \leq t$ pos at once <br> $\rightsquigarrow O(t/\ell \cdot (\ell + k)p) = O(tp)$ | substr <br> start pos |



Prune whenever $e_{i,j} > k$.

Philip Wellnitz

Revitalizing Research on Approximate String Matching Algorithms   10-3

Inter-University Research Institute Corporation
Research Organization of Information and Systems
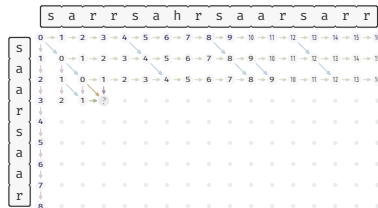**National Institute of Informatics**

## Approximate String Matching

early 1980's

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

| Edit dist/Verify Algorithm | ASM Algorithm | $(t := |T|, p := |P|)$ |
|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ <br> [Sellers, 1980] and others | Verify pos 1 by 1 $\rightsquigarrow O(t^2 p)$ <br> Verify all pos at once $\rightsquigarrow O(tp)$ | substr <br> start pos |
| DP, diagonally compute only entries $\leq k$ <br> $\rightsquigarrow O(t + kp)$ | Verify pos 1 by 1 $\rightsquigarrow O(tkp)$ <br> Verify $\ell \leq t$ pos at once <br> $\rightsquigarrow O(t/\ell \cdot (\ell + k)p) = O(tp)$ | substr <br> start pos |
| compute DP diagonals, jump over common substr in $O(1)$ <br> $\rightsquigarrow O(t + k^2)$ [Landau, Vishkin'89]. | Verify pos 1 by 1 $\rightsquigarrow O(tk^2)$ <br> Verify $\ell \leq t$ pos at once <br> $\rightsquigarrow O(t/\ell \cdot (\ell + k)k) = O(tk)$ | substr <br> start pos |



Compute furthest pos on diag $d$ reachable w/ $\ell = 0, \ldots, k$ edits.
$\rightsquigarrow$ Jump over equal substrings in $O(1)$ time

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
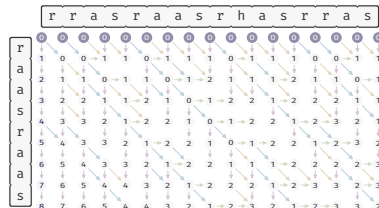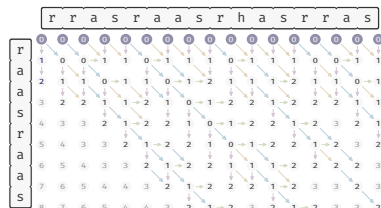Revitalizing Research on Approximate String Matching Algorithms   10-4

# Classical Algorithms

## **Approximate String Matching**

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

| Edit dist/Verify Algorithm | ASM Algorithm | $(t := |T|, p := |P|)$ | |
|---|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ | Verify pos 1 by 1 $\rightsquigarrow O(t^2p)$ | | substr |
| [Sellers, 1980] and others | Verify all pos at once $\rightsquigarrow O(tp)$ | | start pos |
| DP, diagonally compute only entries $\leq k$ | Verify pos 1 by 1 $\rightsquigarrow O(tkp)$ | | substr |
| | Verify $\ell \leq t$ pos at once | | start pos |
| $\rightsquigarrow O(t + kp)$ | $\rightsquigarrow O(t/\ell \cdot (\ell + k)p) = O(tp)$ | | |
| compute DP diagonals, jump over common substr in $O(1)$ | Verify pos 1 by 1 $\rightsquigarrow O(tk^2)$ | | substr |
| | Verify $\ell \leq t$ pos at once | | start pos |
| $\rightsquigarrow O(t + k^2)$ [Landau, Vishkin'89]. | $\rightsquigarrow O(t/\ell \cdot (\ell + k)k) = O(tk)$ | | |



Compute furthest pos on diag $d$ reachable w/ $\ell = 0, \dots, k$ edits.
$\rightsquigarrow$ Jump over equal substrings in $O(1)$ time

## Approximate String Matching

Given: text $T$, pattern $P$, threshold $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

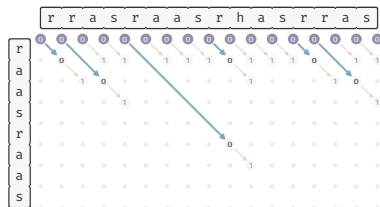| Edit dist/Verify Algorithm | ASM Algorithm | ($t := |T|, p := |P|$) |
|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ [Sellers, 1980] and others | Verify pos 1 by 1 $\rightsquigarrow O(t^2p)$ <br> Verify all pos at once $\rightsquigarrow O(tp)$ | substr <br> start pos |
| DP, diagonally compute only entries $\leq k$ $\rightsquigarrow O(t + kp)$ | Verify pos 1 by 1 $\rightsquigarrow O(tkp)$ <br> Verify $\ell \leq t$ pos at once $\rightsquigarrow O(t/\ell \cdot (\ell + k)p) = O(tp)$ | substr <br> start pos |
| compute DP diagonals, jump over common substr in $O(1)$ $\rightsquigarrow O(t + k^2)$ [Landau, Vishkin'89]. | Verify pos 1 by 1 $\rightsquigarrow O(tk^2)$ <br> Verify $\ell \leq t$ pos at once $\rightsquigarrow O(t/\ell \cdot (\ell + k)k) = O(tk)$ | substr <br> start pos |



Compute furthest pos on diag $d$ reachable w/ $\ell = 0, \dots, k$ edits. $\rightsquigarrow$ Jump over equal substrings in $O(1)$ time
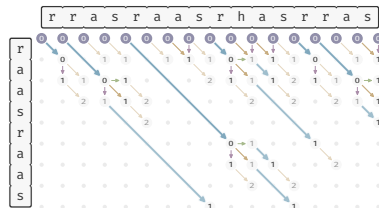
Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   10-6

**Approximate String Matching** early 1980's

Given: text *T*, pattern *P*, threshold *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

| Edit dist/Verify Algorithm | ASM Algorithm | ($t := |T|, p := |P|$) |
|---|---|---|
| Textbook DP $\rightsquigarrow O(tp)$ <br><br> [Sellers, 1980] and others | Verify pos 1 by 1 $\rightsquigarrow O(t^2 p)$ <br> Verify all pos at once $\rightsquigarrow O(tp)$ | substr <br> start pos |
| DP, diagonally compute only entries ≤ *k* <br> $\rightsquigarrow O(t + kp)$ | Verify pos 1 by 1 $\rightsquigarrow O(tkp)$ <br> Verify $\ell \leq t$ pos at once <br> $\rightsquigarrow O(t/\ell \cdot (\ell + k)p) = O(tp)$ | substr <br> start pos |
| compute DP diagonals, jump over common substr in $O(1)$ <br> $\rightsquigarrow O(t + k^2)$ [Landau, Vishkin'89]. | Verify pos 1 by 1 $\rightsquigarrow O(tk^2)$ <br> Verify $\ell \leq t$ pos at once <br> $\rightsquigarrow O(t/\ell \cdot (\ell + k)k) = O(tk)$ | substr <br> start pos |

$\rightsquigarrow$ Didn't use filter yet...

NII Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   10-7

**Approximate String Matching**

early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Idea: Can we filter out diagonals that will never lead to an occurrence?

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
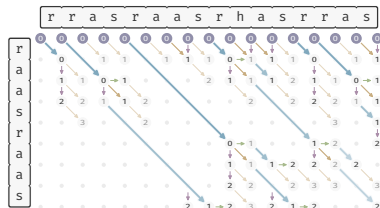Revitalizing Research on Approximate String Matching Algorithms    11-1

**Approximate String Matching** early 1980's

Given: text *T*, pattern *P*, integer *k*;  Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Idea: Can we filter out diagonals that will never lead to an occurrence?

**Approximate String Matching**

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance ≤ $k$ to $P$.

Idea: Can we filter out diagonals that will never lead to an occurrence?
Sometimes we can; and if we cannot, there might be structure to exploit!



Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   11-3

**Approximate String Matching**    early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

Idea: Can we filter out diagonals that will never lead to an occurrence?
Sometimes we can; and if we cannot, there might be structure to exploit!

> **Approximate String Matching** $\rangle$     early 1980's

Given: text *T*, pattern *P*, integer *k*;     Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

◆ Idea: Find parts of *P* that are rare in *T*

⤳ Can filter out candidates for occurrences

**Approximate String Matching**　　　　early 1980's

Given: text *T*, pattern *P*, integer *k*;　Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ Idea: Find parts of *P* that are rare in *T*
  ⤳ Can filter out candidates for occurrences
- ◆ Seen: Highly repetitive parts *R* are bad…
  ⤳ *R* may appear ≈ *t* times in *T*

**Approximate String Matching**    early 1980's

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ Idea: Find parts of *P* that are rare in *T*
  ⤳ Can filter out candidates for occurrences

- ◆ Seen: Highly repetitive parts *R* are bad...
  ⤳ *R* may appear ≈ *t* times in *T*

- ◆ Non-repetitive parts *B* are good!
  ⤳ *B* may appear only ≈ *t*/*b* times in *T*

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    12-3

Inter-University Research Institute Corporation
Research Organization of Information and Systems
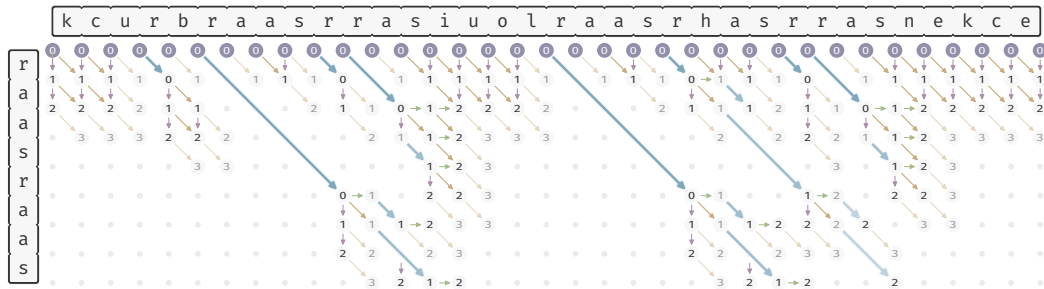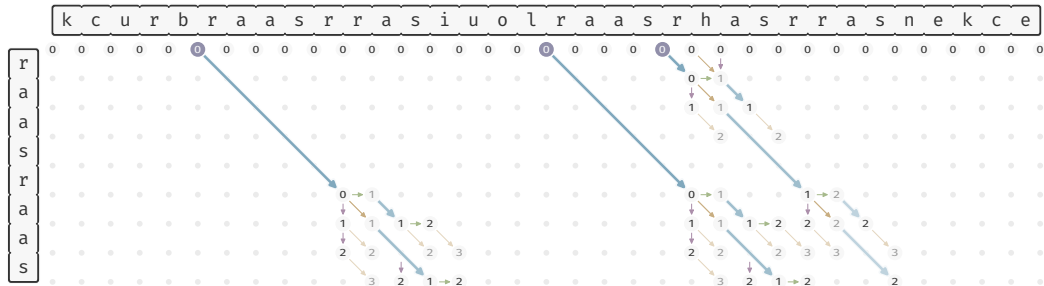National Institute of Informatics

**Approximate String Matching**  early 1980's

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ Idea: Find parts of *P* that are rare in *T*
  ⤳ Can filter out candidates for occurrences
- ◆ Seen: Highly repetitive parts *R* are bad...
  ⤳ *R* may appear ≈ *t* times in *T*
- ◆ Non-repetitive parts *B* are good!
  ⤳ *B* may appear only ≈ *t*/*b* times in *T*
  ⤳ Problem: *B* may appear approximately in *T*

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  12-4

**NI** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Approximate String Matching**

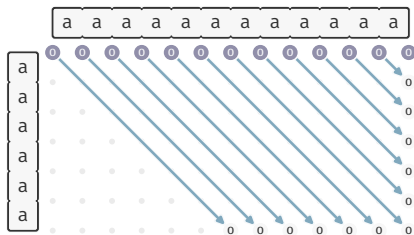Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

◆ Idea: Find parts of *P* that are rare in *T*
  ⤳ Can filter out candidates for occurrences

◆ Seen: Highly repetitive parts *R* are bad…
  ⤳ *R* may appear ≈ *t* times in *T*

◆ Non-repetitive parts *B* are good!
  ⤳ *B* may appear only ≈ *t*/*b* times in *T*
  ⤳ Problem: *B* may appear approximately in *T*
  ⤳ Use > *k* disjoint breaks

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   12-5

Suppose we have: $2k$ disjoint breaks $B_1, \ldots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$    and   ◆ $B_i$ is not periodic   (no long prefix is also a suffix)

$P$

Suppose we have: $2k$ disjoint breaks $B_1, \ldots, B_{2k}$ in $P$ such that

♦ $|B_i| = \Theta(p/k)$   and   ♦ $B_i$ is not periodic   (no long prefix is also a suffix)

♦ In any $k$-edit occ, at least 1 break is matched exactly. (budget for edits is $k$)

Suppose we have: $2k$ disjoint breaks $B_1, \ldots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$    and    ◆ $B_i$ is not periodic   (no long prefix is also a suffix)

◆ In any $k$-edit occ, at least 1 break is matched exactly. (budget for edits is $k$)

◆ Algorithm idea: For each $B_i$, find exact matches of $B_i$ in $T$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   13-3

Suppose we have: $2k$ disjoint breaks $B_1, \ldots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$    and    ◆ $B_i$ is not periodic    (no long prefix is also a suffix)

◆ In any $k$-edit occ, at least 1 break is matched exactly. (budget for edits is $k$)

◆ Algorithm idea: For each $B_i$, find exact matches of $B_i$ in $T$,
Try to extend each exact match into an occ using $O(n + k^2)$-time Verify [LV'89]

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Suppose we have: $2k$ disjoint breaks $B_1, \dots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$  and  ◆ $B_i$ is not periodic  (no long prefix is also a suffix)

◆ In any $k$-edit occ, at least 1 break is matched exactly. (budget for edits is $k$)

◆ Algorithm idea: For each $B_i$, find exact matches of $B_i$ in $T$,
  Try to extend each exact match into an occ using $O(n + k^2)$-time Verify [LV'89]

◆ Have at most $t/(|B_i|/2) = \Theta(kt/p)$ exact occ's of $B_i$

Suppose we have: $2k$ disjoint breaks $B_1, \dots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$   and   ◆ $B_i$ is not periodic   (no long prefix is also a suffix)

◆ In any $k$-edit occ, at least 1 break is matched exactly. (budget for edits is $k$)

◆ Algorithm idea: For each $B_i$, find exact matches of $B_i$ in $T$,
   Try to extend each exact match into an occ using $O(n + k^2)$-time Verify [LV'89]

◆ Have at most $t/(|B_i|/2) = \Theta(kt/p)$ exact occ's of $B_i$
   ⤳ $O(k^2 t/p)$ calls to [LV'89]; $O(k^4 t/p)$ time in total

Suppose we have: $2k$ disjoint breaks $B_1, \ldots, B_{2k}$ in $P$ such that

◆ $|B_i| = \Theta(p/k)$   and   ◆ $B_i$ is not periodic   (no long prefix is also a suffix)

◆ In any $k$-edit occ, at least $k$ breaks are matched exactly. (budget for edits is $k$)

◆ Algorithm idea: For each $B_i$, find exact matches of $B_i$ in $T$,
  Try to extend each exact match into an occ using $O(n + k^2)$-time Verify [LV'89]

◆ Have at most $t/(|B_i|/2) = \Theta(kt/p)$ exact occ's of $B_i$
  $\rightsquigarrow O(k^2t/p)$ calls to [LV'89]; $O(k^4t/p)$ time in total

◆ Can even be improved to $O(k^3t/p)$; but not-enough breaks case is also $O(k^4t/p)$.

Existing algorithms for the case $n \approx t \approx 2p$

Existing algorithms for the case $n \approx t \approx 2p$

Existing algorithms for the case $n \approx t \approx 2p$

Existing algorithms for the case $n \approx t \approx 2p$

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 14-4

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Existing algorithms for the case $n \approx t \approx 2p$

Why did a new improvement take 24 years?

◆ Simpler problems were not fully understood!

- ◆ Simpler problems were not fully understood!
- ◆ Example 1: The complexity of Edit Distance was settled only in 2014

**Hardness of Edit Distance**                    [BI18] (ann. 2014)

Computing the Edit Distance of to strings of length $n$ is not possible in $O(n^{2-\varepsilon})$ time.

(Unless there is a major breakthrough for the Satisfiability Problem.)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   16-2

◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇

Detour: Why did a new improvement for Approximate String Matching take 24 years?

- ◆ Simpler problems were not fully understood!
- ◆ Example 1: The complexity of Edit Distance was settled only in 2014

> **Hardness of Edit Distance** [BI18] (ann. 2014)
>
> Computing the Edit Distance of to strings of length $n$ is not possible in $O(n^{2-\varepsilon})$ time.
> (Unless there is a major breakthrough for the Satisfiability Problem.)

⇝ Yields lower bound of $O(t + k^2 \cdot t/p)$ for Approximate String Matching.

Philip Wellnitz

Revitalizing Research on Approximate String Matching Algorithms 16-3

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

- ◆ Simpler problems were not fully understood!
- ◆ Example 2: String Matching with Mismatches (no insertions or deletions of chars)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 17-1

◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◆◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇

Detour: Why did a new improvement for Approximate String Matching take 24 years?

- ◆ Simpler problems were not fully understood!
- ◆ Example 2: String Matching with Mismatches (no insertions or deletions of chars)
  - ⤳ $O(tk)$ algorithm [Landau, Vishkin'**86**]
  - ⤳ $\tilde{O}(t\sqrt{k})$ algorithm [Amir, Lewenstein, Porat'**04**]
  - ⤳ $\tilde{O}(t + t/p \cdot k^2)$ algorithm [CFPSS'**16**]

Philip Wellnitz

Revitalizing Research on Approximate String Matching Algorithms   17-2

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇

Detour: Why did a new improvement for Approximate String Matching take 24 years?

◆ Simpler problems were not fully understood!

◆ Example 2: String Matching with Mismatches (no insertions or deletions of chars)
  ⇝ $O(tk)$ algorithm [Landau, Vishkin'**86**]
  ⇝ $\tilde{O}(t\sqrt{k})$ algorithm [Amir, Lewenstein, Porat'**04**]
  ⇝ $\tilde{O}(t + t/p \cdot k^2)$ algorithm [CFPSS'**16**]

⟨ **Optimal String Matching with Mismatches** ⟩     [GU18] (ann. 2017)

There is a $\tilde{O}(t + kt/\sqrt{p})$-time algorithm for String Matching with Mismatches
and no significantly faster algorithm exits.

(Unless there is a major breakthrough for combinatorial Boolean Matrix Multiplication.)

✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦✦◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇◇

Detour: Why did a new improvement for Approximate String Matching take 24 years?

◆ Simpler problems were not fully understood!

◆ Example 2: String Matching with Mismatches (no insertions or deletions of chars)

⤳ $O(tk)$ algorithm [Landau, Vishkin'**86**]

⤳ $\tilde{O}(t\sqrt{k})$ algorithm [Amir, Lewenstein, Porat'**04**]

⤳ $\tilde{O}(t + t/p \cdot k^2)$ algorithm [CFPSS'**16**]

⟨**Optimal String Matching with Mismatches**⟩     [GU18] (ann. 2017)

There is a $\tilde{O}(t + kt/\sqrt{p})$-time algorithm for String Matching with Mismatches
and no significantly faster algorithm exits.

(Unless there is a major breakthrough for combinatorial Boolean Matrix Multiplication.)

⤳ String Matching with Mismatches was "fully" understood only in 2017.

How do we obtain faster algorithms?

How do we obtain faster algorithms?

New insights into the solution structure of
Approximate String Matching

Think: Better filter + more structure when filter fails

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 18-2

Step 0:

What is the solution structure of
Exact String Matching?

NII  Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  19

## (The) Period of a String

| | | |
|---|---|---|
| $p > 0$ is a period of $S$ | $:\iff$ | $S[i] = S[i+p]$ for all $i = 1, \dots, |S| - p$ |
| The period of $S$ | $:\iff$ | smallest period of $S$, write $\mathrm{per}(S)$ |
| $S$ is periodic | $:\iff$ | $\mathrm{per}(S) \leq |S|/2$ |



$S$  a b c a b c a b c a b

$\mathrm{per}(S) = 3$

**(The) Period of a String**

$p > 0$ is a period of $S$ :$\iff$ $S[i] = S[i + p]$ for all $i = 1, \ldots, |S| - p$

The period of $S$ :$\iff$ smallest period of $S$, write $\mathrm{per}(S)$

$S$ is periodic :$\iff$ $\mathrm{per}(S) \leq |S|/2$



$\mathrm{per}(S) = 3$

6 and 9 also periods of $S$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 20-2

**"Periodicity Lemma"**

(Folklore)

Text *T*, pattern *P* with $t \leq \frac{3}{2} p$; one of the following holds:

◆ *P* appears ≤ 1 times in *T*.

**"Periodicity Lemma"** (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;     one of the following holds:

◆ $P$ appears ≤ 1 times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; one of the following holds:

◆ $P$ appears $\leq 1$ times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;     one of the following holds:

◆ $P$ appears ≤ 1 times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \leq {}^3\!/_2\, p$; one of the following holds:

◆ $P$ appears $\leq 1$ times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

**"Periodicity Lemma"** (Folklore)

Text $T$, pattern $P$ with $t \le \frac{3}{2} p$; one of the following holds:

◆ $P$ appears ≤ 1 times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \le \frac{3}{2} p$; one of the following holds:

◆ $P$ appears ≤ 1 times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    one of the following holds:

◆ $P$ appears $\leq 1$ times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

## "Periodicity Lemma" (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;     one of the following holds:

◆ $P$ appears ≤ 1 times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   21-9

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**"Periodicity Lemma"** (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$; one of the following holds:
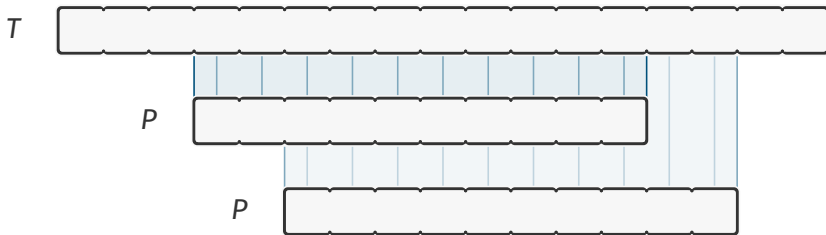
◆ $P$ appears $\leq 1$ times in $T$.

◆ $P$ (and the relevant part of $T$) are periodic with some period $Q$.



### Standard Trick is useful here:

For $t \gg p$ consider separately $O(t/p)$ fragments of $T$; then Periodicity Lemma applies

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 21-10

Step 1:

What is the solution structure of
String Matching with Mismatches?

**Periodicity Lemma** (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   23-1

**Periodicity Lemma**  (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

**Main Result (Mismatches)**  [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

## Periodicity Lemma
(Folklore)

Text $T$, pattern $P$ with $t \le \frac{3}{2}p$;    $P$ appears $\le 1$ times in $T$    or    $P$ is periodic with some period $Q$.

## Main Result (Mismatches)
[BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \le \frac{3}{2}p$; threshold $k$, one of the following holds

$P$ appears $\le O(k)$ times in $T$



*P* appears $2k$ times w/ $\le k$ mism.
but *P* far from periodic

## Periodicity Lemma

(Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

## Main Result (Mismatches)

[BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

> **Periodicity Lemma** > (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

> **Main Result (Mismatches)** > [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   23-5

**Periodicity Lemma** (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$    or    $P$ is almost periodic with some period $Q$

$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

$P$ and $T$ have approximate period a
(are at HD $\leq 2k$ from aa...aa)

## Periodicity Lemma (Folklore)

Text $T$, pattern $P$ with $t \le \frac{3}{2}\,p$;   $P$ appears $\le 1$ times in $T$   or   $P$ is periodic with some period $Q$.

## Main Result (Mismatches) [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \le \frac{3}{2}\,p$; threshold $k$, one of the following holds

$P$ appears $\le O(k)$ times in $T$   or   $P$ is almost periodic with some period $Q$



*P* appears $2k$ times w/ $\le k$ mism.
but *P* far from periodic



*P* and *T* have approximate period a
(are at HD $\le 2k$ from aa…aa)
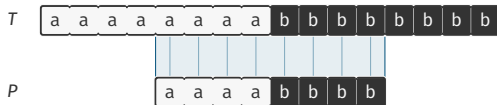
**Periodicity Lemma** (Folklore)

Text *T*, pattern *P* with $t \leq \frac{3}{2} p$;  *P* appears ≤ 1 times in *T*  or  *P* is periodic with some period *Q*.

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

Text *T*, pattern *P* with $t \leq \frac{3}{2} p$; threshold *k*, one of the following holds

*P* appears ≤ $O(k)$ times in *T*  or  *P* is almost periodic with some period *Q*



*P* appears 2*k* times w/ ≤ *k* mism.
but *P* far from periodic

*P* and *T* have approximate period a
(are at HD ≤ 2*k* from aa…aa)

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 23-8

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

## Periodicity Lemma

(Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$;    $P$ appears $\leq 1$ times in $T$    or    $P$ is periodic with some period $Q$.

## Main Result (Mismatches)

[BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$    or    $P$ is almost periodic with some period $Q$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

$P$ and $T$ have approximate period a
(are at HD $\leq 2k$ from aa...aa)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    23-9

**Periodicity Lemma** (Folklore)

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$;  $P$ appears $\leq 1$ times in $T$  or  $P$ is periodic with some period $Q$.

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k)$ times in $T$  or  $P$ is almost periodic with some period $Q$



$P$ appears $2k$ times w/ $\leq k$ mism.
but $P$ far from periodic

$P$ and $T$ have approximate period a
(are at HD $\leq 2k$ from aa…aa)

Step 1.5:

The solution structure of
Approximate String Matching.

**Main Result (Mismatches)** [BKü**W**'19; ChKo**W**'20]

$T$, $P$ with $t \le \frac{3}{2} p$; threshold $k$,  $P$ appears $O(k)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Main Result (Edits)** [ChKo**W**'20]

Text $T$, pattern $P$ with $t \le \frac{3}{2} p$; threshold $k$, one of the following holds

or     $P$ is almost periodic with some period $Q$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  25-1

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

$T$, $P$ with $t \leq \frac{3}{2}\,p$; threshold $k$, $P$ appears $O(k)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Main Result (Edits)** [ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2}\,p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$

$P$ appears $\approx k^2$ times with $\leq k$ edits but $P$ and $T$ far from (approx.) periodic

**Main Result (Mismatches)**  [BKüW'19; ChKoW'20]

$T$, $P$ with $t \leq \frac{3}{2}p$; threshold $k$, $P$ appears $O(k)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Main Result (Edits)**  [ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2}p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$



$T$: a a a a a a a a a a a a a a **b** a a **b** a a **b** a a **b** a a

$P$: a a a a a a **b** a a **b** a a

$P$ appears $\approx k^2$ times with $\leq k$ edits    but $P$ and $T$ far from (approx.) periodic

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    25-3

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

$T$, $P$ with $t \leq \frac{3}{2} p$; threshold $k$, $P$ appears $O(k)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Main Result (Edits)** [ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k^2)$ times in $T$    or    $P$ is almost periodic with some period $Q$



$P$ appears $\approx k^2$ times with $\leq k$ edits    but $P$ and $T$ far from (approx.) periodic

**Main Result (Mismatches)** [BKüW'19; ChKoW'20]

$T$, $P$ with $t \leq \frac{3}{2} p$; threshold $k$, $P$ appears $O(k)$ times in $T$ or $P$ is almost periodic with some period $Q$
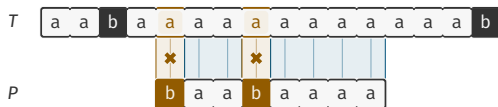
**Main Result (Edits)** [ChKoW'20]

Text $T$, pattern $P$ with $t \leq \frac{3}{2} p$; threshold $k$, one of the following holds

$P$ appears $\leq O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$



$T$ | a a a a a a a a a a a a b a a b a a b a a b a a

$P$ | a a a a a a b a a b a a

$P$ appears $\approx k^2$ times with $\leq k$ edits but $P$ and $T$ far from (approx.) periodic

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  25-5

**Main Result (Edits)** [ChKo**W**'20]

$T$, $P$ with $t \leq \frac{3}{2} p$; threshold $k$, $P$ appears $O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 26-1

NII  Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Main Result (Edits)** [ChKo**W**'20]

$T$, $P$ with $t \le \frac{3}{2} p$; threshold $k$, $P$ appears $O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Key Intermediate Result (Analyze)** [ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ is almost periodic.                          aaacaaaaaaaaacaaaaaaa

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  26-2

**Main Result (Edits)**

[ChKo**W**'20]

*T*, *P* with $t \leq {}^3\!/_2\, p$; threshold *k*, *P* appears $O(k^2)$ times in *T*  or *P* is almost periodic with some period *Q*

**Key Intermediate Result (Analyze)**

[ChKo**W**'20]

Every string *P* satisfies at least one of

◆ *P* has 2*k* disjoint, long breaks           c✿✿✿✿✿a✿✿✿✿✿a

◆ *P* is almost periodic.           aaacaaaaaaaaacaaaaaaa

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Main Result (Edits)** [ChKo**W**'20]

$T$, $P$ with $t \leq {}^3/_2\, p$; threshold $k$, $P$ appears $O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Key Intermediate Result (Analyze)** [ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks      c❀✿❀❀✿❀a❀✿❀❀✿❀a

◆ $P$ has disjoint repetitive regions that cover ${}^3/_8\, P$      ❀❀aaaaaa❀❀ccaccc❀❀

◆ $P$ is almost periodic.      aaacaaaaaaaaacaaaaaaa

**Main Result (Edits)** [ChKo**W**'20]

$T$, $P$ with $t \leq {}^{3}/_{2}\, p$; threshold $k$, $P$ appears $O(k^2)$ times in $T$ or $P$ is almost periodic with some period $Q$

**Key Intermediate Result (Analyze)** [ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks      c✿❀✿❀❀✿❀a✿❀✿❀❀✿❀a

◆ $P$ has disjoint repetitive regions that cover ${}^{3}/_{8}\, P$      ❀❀aaaaaaa❀❀ccaccc❀❀

◆ $P$ is almost periodic.      aaacaaaaaaaaacaaaaaaa

Analyze implies Main Result:

Almost periodicity ⤳ as in Main Result

Breaks, repetitive regions ⤳ good filter (as in [Cole, Hariharan'98])

## Key Intermediate Result (Analyze)

[ChKoW'20]

Every string *P* satisfies at least one of

◆ *P* has 2*k* disjoint, long breaks

◆ *P* has disjoint repetitive regions that cover $^3/_8$ *P*

◆ *P* is almost periodic.

c✿✿✿✿✿✿✿a✿✿✿✿✿✿✿a

✿✿aaaaaaa✿✿ccaccc✿✿

aaacaaaaaaaaacaaaaaaa

*P*

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   27-1

## Key Intermediate Result (Analyze)

[ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks      c❀✿❀❀✳✳a❀✿❀❀✳✳a

◆ $P$ has disjoint repetitive regions that cover $^3/_8$ $P$      ✳✳aaaaaa❀❀ccaccc✳✳

◆ $P$ is almost periodic.      aaacaaaaaaaaacaaaaaaa

$P$ ▢

◆ Process $P$ from left to right, $p/8k$ new characters at a time.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   27-2

**Key Intermediate Result (Analyze)**

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks      c✿✿✿✿✿✿a✿✿✿✿✿✿a

◆ $P$ has disjoint repetitive regions that cover $^3/_8$ $P$      ✿✿aaaaaa✿✿ccaccc✿✿

◆ $P$ is almost periodic.      aaacaaaaaaaaacaaaaaaa

$P$   $\boxed{B_1}$

Breaks $B = \left\{ \ \boxed{B_1} \ \right\}$

Repetitive Regions $R = \left\{ \quad \right\}$

◆ If a fragment is a break, add it to the found breaks.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   27-3

**⟨ Key Intermediate Result (Analyze) ⟩**

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks      c❀✿❀❀✳✳a✳✿❀❀✳✳a

◆ $P$ has disjoint repetitive regions that cover $^3/_8$ $P$      ✳✳aaaaaa❀✳ccaccc✳✳

◆ $P$ is almost periodic.      aaacaaaaaaaaacaaaaaaa

$P$   $B_1$ □

Breaks $B = \left\{ \; B_1 \; \right\}$

Repetitive Regions $R = \left\{ \quad \right\}$

◆ Otherwise, find the shortest prefix (longer than $p/8k$) that is a repetitive region.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   27-4

**Key Intermediate Result (Analyze)**                                    [ChKoW'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks                                    c✿✿✿✿✿✿a✿✿✿✿✿✿a

◆ $P$ has disjoint repetitive regions that cover $\frac{3}{8} P$        ✿✿aaaaaa✿✿ccaccc✿✿

◆ $P$ is almost periodic.                                               aaacaaaaaaaaacaaaaaaa

$P$   $B_1$  $R_1$

Breaks $B = \left\{ \quad B_1 \quad \right\}$

Repetitive Regions $R = \left\{ \quad R_1 \quad \right\}$

◆ Otherwise, find the shortest prefix (longer than $p/8k$) that is a repetitive region.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   27-5

**Key Intermediate Result (Analyze)**

[ChKoW'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks

◆ $P$ has disjoint repetitive regions that cover $\frac{3}{8}$ $P$

◆ $P$ is almost periodic.

c✿✿✿✿✿✿a✿✿✿✿✿✿a

✿✿aaaaaa✿✿ccaccc✿✿

aaacaaaaaaaaacaaaaaaa



$P$   $B_1$ | $R_1$ | $B_2$ | $B_3$ | $R_2$ | $B_4$

Breaks $B$ = { $B_1$ , $B_2$ , $B_3$ , $B_4$ }

Repetitive Regions $R$ = { $R_1$ }

◆ If we found $2k$ breaks, return the breaks.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 27-6

## Key Intermediate Result (Analyze)

[ChKo**W**'20]

Every string $P$ satisfies at least one of

- $P$ has $2k$ disjoint, long breaks
- $P$ has disjoint repetitive regions that cover $\frac{3}{8} P$
- $P$ is almost periodic.

- If the total length of the repetitive regions is $> 3/8 \cdot p$, return the repetitive regions.

## Key Intermediate Result (Analyze)

[ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks

c✿✿✿✿✿✿a✿✿✿✿✿✿a

◆ $P$ has disjoint repetitive regions that cover $^3/_8$ $P$

✿✿aaaaaa✿✿ccaccc✿✿

◆ $P$ is almost periodic.

aaacaaaaaaaaacaaaaaaa

$P$ $\boxed{B_1}$ $\boxed{R_1}$ $\boxed{B_2}$ $\boxed{B_3}$ $\boxed{\phantom{xxxxxxxxxxxxxxxxxxxxxxxx}}$

Breaks $B = \left\{ \boxed{B_1} , \boxed{B_2} , \boxed{B_3} \right\}$

Repetitive Regions $R = \left\{ \boxed{R_1} \right\}$

◆ If we reach the end of $P$, try to find a single repetitive region starting from the end.

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 27-8

## Key Intermediate Result (Analyze)

[ChKo**W**'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks

c✿✿✿✿✿✿a✿✿✿✿✿a

◆ $P$ has disjoint repetitive regions that cover $\frac{3}{8}$ $P$

✿✿aaaaaa✿✿ccaccc✿✿

◆ $P$ is almost periodic.

aaacaaaaaaaaacaaaaaaa



$P$      $R_1$

Breaks $B = \left\{ \quad \right\}$

Repetitive Regions $R = \left\{ \boxed{\qquad R_1 \qquad} \right\}$

◆ If we found a repetitive region, return it.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 27-10

⟨ **Key Intermediate Result (Analyze)** ⟩                    [ChKo**W**'20]

Every string *P* satisfies at least one of

◆ *P* has 2*k* disjoint, long breaks                    c✿✿✿✿✿✿✿a✿✿✿✿✿✿a

◆ *P* has disjoint repetitive regions that cover $^3/_8$ *P*          ✿✿aaaaaa✿✿ccaccc✿✿

◆ *P* is almost periodic.                         aaacaaaaaaaaacaaaaaaa



*P*

Breaks *B* = { }

Repetitive Regions *R* = { }

◆ If we again don't obtain a repetitive region, *P* is almost periodic.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 27-11

## Key Intermediate Result (Analyze) ✓

[ChKoW'20]

Every string $P$ satisfies at least one of

◆ $P$ has $2k$ disjoint, long breaks

◆ $P$ has disjoint repetitive regions that cover $\frac{3}{8} P$

◆ $P$ is almost periodic.

```
aaacaaaaaaaaacaaaaaaa
```

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 27-12

How do we turn our insights
into faster algorithms?

How do we turn our insights into faster algorithms?

Need to tackle three cases.

- ◆ $P$ contains $2k$ disjoint breaks;
- ◆ $P$ contains disjoint repetitive regions $R_i$;
- ◆ $P$ is almost periodic

How do we turn our insights into faster algorithms?

Need to tackle ~~three~~ two cases.

◆ $P$ contains $2k$ disjoint breaks;

◆ ~~$P$ contains disjoint repetitive regions $R_i$;~~ ⤳ Follows from the other two cases.

◆ $P$ is almost periodic

How do we turn our insights into faster algorithms?

Need to tackle ~~three~~ one case.

- ~~$P$ contains $2k$ disjoint breaks;~~ ⤳ Adaption of [Cole,Hariharan'98] yields $O(|T| + |T|/|P| \cdot k^3)$.
- ~~$P$ contains disjoint repetitive regions $R_i$;~~ ⤳ Follows from the other two cases.
- $P$ is almost periodic

How do we turn our insights into faster algorithms?

Need to tackle ~~three~~ one case.

◆ ~~*P* contains 2*k* disjoint breaks;~~ ↝ Adaption of [Cole,Hariharan'98] yields $O(|T| + |T|/|P| \cdot k^3)$.

◆ ~~*P* contains disjoint repetitive regions $R_i$;~~ ↝ Follows from the other two cases.

◆ *P* is almost periodic

⟨ **Reduction to Periodic Patterns** ⟩     [ChKo**W**'22]

Algorithm for almost periodic case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$, for $a \geq 3$

$\implies$    Algorithm for general case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$

⟨ **Reduction to Periodic Patterns** ⟩ [ChKo**W**'22]

Algorithm for almost periodic case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$, for $a \geq 3$
$\implies$ Algorithm for general case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$

Need to tackle: $P$ is almost periodic

**Reduction to Periodic Patterns** [ChKo**W**'22]

Algorithm for almost periodic case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$, for $a \geq 3$
$\implies$ Algorithm for general case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$

Need to tackle: $P$ is almost periodic

◆ In [ChKo**W**'20], use elaborate marking scheme for $O(|T| + |T|/|P| \cdot k^4)$ algo
$\rightsquigarrow$ Not faster than [Cole,Hariharan'98]

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 31-2

**Reduction to Periodic Patterns** [ChKo**W**'22]

Algorithm for almost periodic case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$, for $a \geq 3$

$\implies$ Algorithm for general case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$

Need to tackle: $P$ is almost periodic

◆ In [ChKo**W**'20], use elaborate marking scheme for $O(|T| + |T|/|P| \cdot k^4)$ algo
  ⇝ Not faster than [Cole,Hariharan'98]
◆ In [ChKo**W**'22], trade-off between
  ◆ Refinement of algorithm from [ChKo**W**'20]
  ◆ New algorithm based on "Seaweed Technology" of [Tiskin'10,'15]
  ⇝ Yields $O(|T| + |T|/|P| \cdot k^{3.5})$ algorithm

NII Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 31-3

> **Reduction to Periodic Patterns** [ChKo**W**'22]

Algorithm for almost periodic case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$, for $a \geq 3$

$\implies$ Algorithm for general case in time $\tilde{O}(|T| + k^a \cdot |T|/|P|)$

Need to tackle: $P$ is almost periodic

- ◆ In [ChKo**W**'20], use elaborate marking scheme for $O(|T| + |T|/|P| \cdot k^4)$ algo
  ⤳ Not faster than [Cole,Hariharan'98]
- ◆ In [ChKo**W**'22], trade-off between
  - ◆ Refinement of algorithm from [ChKo**W**'20]
  - ◆ ~~New algorithm based on "Seaweed Technology" of [Tiskin'10,'15]~~
    New in [ChKo**W**'25]: Simpler plug-in replacement for "Seaweed Technology" based on SMAWK (efficient (min, +)-multiplication of Monge matrices)
  ⤳ Yields $O(|T| + |T|/|P| \cdot k^{3.5})$ algorithm

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 31-4

◆ All of our algorithms rely on a small set of essential operations:

- ◆ `LCP`($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
- ◆ `LCP`$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
- ◆ `IPM`($P$, $T$): Compute all exact matches of $P$ in $T$.
- ◆ `Length`($S$): Compute the length $s$ of $S$.
- ◆ `Access`($S$, $i$): Retrieve the character $S[i]$.
- ◆ `Extract`($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r)$ from $S$.

⟨ **Main Result (**`PILLAR` **Algorithm)** ⟩ [ChKo**W**'22]

For $t \leq \frac{3}{2} p + k$; ASM can be solved using $O(k^{3.5})$ `PILLAR` operations

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  32-1

◆ All of our algorithms rely on a small set of essential operations:

- ◆ LC**P**($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
- ◆ **L**CP$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
- ◆ **I**PM($P$, $T$): Compute all exact matches of $P$ in $T$.
- ◆ **L**ength($S$): Compute the length $s$ of $S$.
- ◆ **A**ccess($S$, $i$): Retrieve the character $S[i]$.
- ◆ Ext**R**act($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r)$ from $S$.

> **Main Result (**PILLAR **Algorithm)**                                    [ChKoW'22]
>
> For $t \le \frac{3}{2}\,p + k$; ASM can be solved using $O(k^{3.5})$ PILLAR operations

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    32-2

◆ All of our algorithms rely on a small set of essential operations:

- ◆ LC**P**($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
- ◆ **L**CP$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
- ◆ **I**PM($P$, $T$): Compute all exact matches of $P$ in $T$.
- ◆ **L**ength($S$): Compute the length $s$ of $S$.
- ◆ **A**ccess($S$, $i$): Retrieve the character $S[i]$.
- ◆ Ext**R**act($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r]$ from $S$.

**Main Result (**PILLAR **Algorithm)**

[ChKo**W**'22]

For $t \leq \frac{3}{2} p + k$; ASM can be solved using $O(k^{3.5})$ PILLAR operations

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  32-3

◆ All of our algorithms rely on a small set of essential operations:

- ◆ LC**P**($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
- ◆ **L**CP$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
- ◆ **I**PM($P$, $T$): Compute all exact matches of $P$ in $T$.
- ◆ **L**ength($S$): Compute the length $s$ of $S$.
- ◆ **A**ccess($S$, $i$): Retrieve the character $S[i]$.
- ◆ Ext**R**act($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r]$ from $S$.

> **Main Result (**PILLAR **Algorithm)**                    [ChKo**W**'22]
>
> For $t \le \frac{3}{2}\,p + k$; ASM can be solved using $O(k^{3.5})$ PILLAR operations

◆ In [ChKo**W**'20]: efficient PILLAR implementations for

- ◆ Standard setting (each op $O(1)$ after global $O(t)$ preprocessing)
- ◆ Fully-compressed setting ($T$, $P$ grammar-compressed to size $n$; each op $\tilde{O}(1)$ after $\tilde{O}(n)$ prep.)
  (In this setting, already [ChKo**W**'20] obtained faster algorithms)
- ◆ Dynamic Setting

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  32-4

◆ All of our algorithms rely on a small set of essential operations:

   ◆ LC**P**($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
   ◆ **L**CP$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
   ◆ **I**PM($P$, $T$): Compute all exact matches of $P$ in $T$.
   ◆ **L**ength($S$): Compute the length $s$ of $S$.
   ◆ **A**ccess($S$, $i$): Retrieve the character $S[i]$.
   ◆ Ext**R**act($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r]$ from $S$.

> **Main Result (** `PILLAR` **Algorithm)**          [ChKo**W**'22]
>
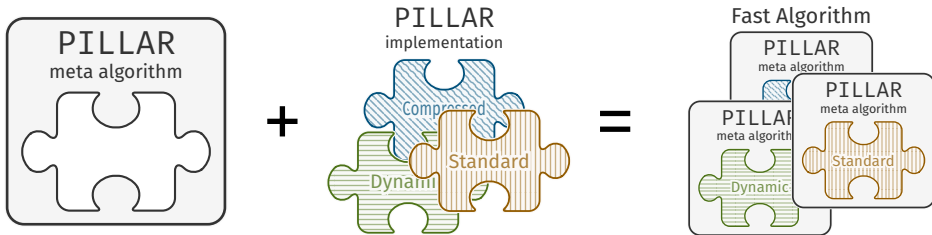> For $t \leq \frac{3}{2}p + k$; ASM can be solved using $O(k^{3.5})$ `PILLAR` operations

◆ In [ChKo**W**'20]: efficient `PILLAR` implementations for
   ◆ Standard setting (each op $O(1)$ after global $O(t)$ preprocessing)
   ◆ Fully-compressed setting ($T$, $P$ grammar-compressed to size $n$; each op $\tilde{O}(1)$ after $\tilde{O}(n)$ prep.)
     (In this setting, already [ChKo**W**'20] obtained faster algorithms)
   ◆ Dynamic Setting
◆ Since then: quantum setting (more details soon!), packed setting, read-only setting, …
  `PILLAR` algorithms for other related problems, etc. [ChKo**W**'20] has > 60 citations by now

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms          32-5

◆ All of our algorithms rely on a small set of essential operations:

   ◆ LC**P**($S$, $T$): Compute the length of the longest common prefix of $S$ and $T$.
   ◆ **L**CP$^R$($S$, $T$): Compute the length of the longest common suffix of $S$ and $T$.
   ◆ **I**PM($P$, $T$): Compute all exact matches of $P$ in $T$.
   ◆ **L**ength($S$): Compute the length $s$ of $S$.
   ◆ **A**ccess($S$, $i$): Retrieve the character $S[i]$.
   ◆ Ext**R**act($S$, $\ell$, $r$): Extract the fragment (or substring) $S[\ell..r]$ from $S$.

> **Main Result (`PILLAR` Algorithm)**                                    [ChKo**W**'22]
>
> For $t \le \frac{3}{2}p + k$; ASM can be solved using $O(k^{3.5})$ `PILLAR` operations

◆ In [ChKo**W**'20]: efficient `PILLAR` implementations for
   ◆ Standard setting (each op $O(1)$ after global $O(t)$ preprocessing)
   ◆ Fully-compressed setting ($T$, $P$ grammar-compressed to size $n$; each op $\tilde{O}(1)$ after $\tilde{O}(n)$ prep.)
     (In this setting, already [ChKo**W**'20] obtained faster algorithms)
   ◆ Dynamic Setting
◆ Since then: quantum setting (more details soon!), packed setting, read-only setting, ...
   `PILLAR` algorithms for other related problems, etc. [ChKo**W**'20] has > 60 citations by now
◆ Bottom-line, useful design principle: meta algorithm (`PILLAR`) + simple-to-solve subproblems

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms       32-6

Spotlight Extension: Quantum Algorithms for ASM

**Approximate String Matching**　　　　early 1980's

Given: text $T$, pattern $P$, integer $k$;　Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

- ◆ $P$/$T$ given as oracle ⤳ queries possible in superposition
- ◆ Query complexity $Q(n)$: number of oracle queries
- ◆ Time complexity $T(n)$: number of elementary gates

## Approximate String Matching

early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ *P*/*T* given as oracle ⇝ queries possible in superposition
- ◆ Query complexity $Q(n)$: number of oracle queries
- ◆ Time complexity $T(n)$: number of elementary gates
- ◆ Goal: algorithms with $Q(n) \ll n$ and $T(n) \ll n$

**Approximate String Matching**

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

◆ *P*/*T* given as oracle ⤳ queries possible in superposition
◆ Query complexity $Q(n)$: number of oracle queries
◆ Time complexity $T(n)$: number of elementary gates
◆ Goal: algorithms with $Q(n) \ll n$ and $T(n) \ll n$

Toy example: "Standard Trick" for ASM decision version (∃? occ)
Split *T* into overlapping fragments of len $\ell + p + k$
⤳ $O(t/\ell)$ instances, each "responsible" for its first $\ell$ positions
⤳ Classically: $O(t/\ell)$ time overhead (evaluate each inst. 1 by 1)

## Approximate String Matching

early 1980's

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ *P*/*T* given as oracle ⇝ queries possible in superposition
- ◆ Query complexity $Q(n)$: number of oracle queries
- ◆ Time complexity $T(n)$: number of elementary gates
- ◆ Goal: algorithms with $Q(n) \ll n$ and $T(n) \ll n$

Toy example: "Standard Trick" for ASM decision version (∃? occ)

Split *T* into overlapping fragments of len $\ell + p + k$

⇝ $O(t/\ell)$ instances, each "responsible" for its first $\ell$ positions

⇝ Classically: $O(t/\ell)$ time overhead (evaluate each inst. 1 by 1)

Quantum setting: Use amplitude amp. / Grover's algo.

$\tilde{O}(\sqrt{t/\ell})$ time and $O(\sqrt{t/\ell})$ evaluations of single inst.

## Amplitude Amplification  [Gro96, BHMT02]

Given function $f : [\,0\mathinner{..}n\,) \to \{0, 1\}$

Can (wp ≥ 2/3) obtain $x \in f^{-1}(1)$ in $\tilde{O}(\sqrt{n})$ time + $O(\sqrt{n})$ queries to $f$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    34-4

**Approximate String Matching** early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.
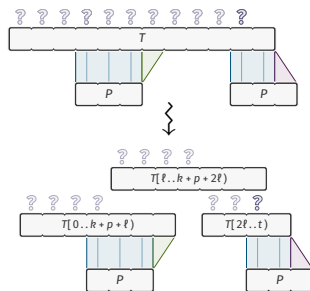
◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Idea: Use existing quantum algorithms to implement `PILLAR` operations.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 35-1

**Approximate String Matching**                              early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

- ◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Idea: Use existing quantum algorithms to implement PILLAR operations.

$$\text{LC}\textbf{P}(S,T); \; \textbf{I}\text{PM}(P,T); \; \textbf{L}\text{ength}(S); \; \textbf{A}\text{ccess}(S,i); \; \text{Ext}\textbf{R}\text{act}(S,\ell,r)$$

PILLAR **Algorithm**   [ChKoW'22]

For $t \leq \frac{3}{2}\,p + k$;
ASM can be solved
using $O(k^{3.5})$ PILLAR operations

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz   35-2
Revitalizing Research on Approximate String Matching Algorithms

**Approximate String Matching**                    early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Idea: Use existing quantum algorithms to implement PILLAR operations.

LC**P**$(S,T)$; **I**PM$(P,T)$; **L**ength$(S)$; **A**ccess$(S,i)$; Ext**R**act$(S,\ell,r)$

PILLAR **Algorithm**                [ChKoW'22]

For $t \leq \frac{3}{2} p + k$;
ASM can be solved
using $O(k^{3.5})$ PILLAR operations

**Quantum** PILLAR **Impl**        implied by [HV'03]

For length-$n$ strings,
can implement all PILLAR ops
in $\tilde{O}(\sqrt{n})$ time and queries

Yields $\tilde{O}(t/p \cdot k^{3.5} \cdot \sqrt{p})$ quantum time algo for ASM; ($\tilde{O}(\sqrt{t/p} \cdot k^{3.5} \cdot \sqrt{p})$ for $\exists$?) (algos are correct whp.)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   35-3

**Approximate String Matching**                          early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

◆ *P*/*T* given as oracle ⤳ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Idea: Use existing quantum algorithms to implement PILLAR operations.

$$\text{LC}\textbf{P}(S,T); \textbf{I}\text{PM}(P,T); \textbf{L}\text{ength}(S); \textbf{A}\text{ccess}(S,i); \text{Ext}\textbf{R}\text{act}(S,\ell,r)$$

PILLAR **Algorithm** [ChKoW'22]

For $t \leq \frac{3}{2}p + k$;
ASM can be solved
using $O(k^{3.5})$ PILLAR operations

**Quantum** PILLAR **Impl** implied by [HV'03]

For length-*n* strings,
can implement all PILLAR ops
in $\tilde{O}(\sqrt{n})$ time and queries

Yields $\tilde{O}(t/p \cdot k^{3.5} \cdot \sqrt{p})$ quantum time algo for ASM; ($\tilde{O}(\sqrt{t/p} \cdot k^{3.5} \cdot \sqrt{p})$ for ∃?) (algos are correct whp.)
Are these algorithms optimal? ⤳ No.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   35-4

> **Approximate String Matching** early 1980's

Given: text *T*, pattern *P*, integer *k*;   Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ *P*/*T* given as oracle ⤳ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all *k*-edit occ's of *P* in *T*

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   36-1

**Approximate String Matching**

early 1980's

Given: text $T$, pattern $P$, integer $k$;   Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

- ◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all $k$-edit occ's of $P$ in $T$

$\rightsquigarrow$ $k$-edit-aware lossy compression of $P$ and $T$ (think: filter + discarding filtered-out parts)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   36-2

> **Approximate String Matching** early 1980's

Given: text $T$, pattern $P$, integer $k$;    Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

- ◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all $k$-edit occ's of $P$ in $T$
 $\rightsquigarrow$ $k$-edit-aware lossy compression of $P$ and $T$ (think: filter + discarding filtered-out parts)

In [KoN**W**'24, KoN**W**'25] for $T$ and $P$ with $t \leq \frac{3}{2}\,p + k$:

- ◆ Show: any $k$-edit-aware lossy compr. of $P$ and $T$ needs $\Omega(k \log(t/k))$ bits

Inter-University Research Institute Corporation
Research Organization of Information and Systems
**National Institute of Informatics**

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  36-3

⟨ **Approximate String Matching** ⟩      early 1980's

Given: text $T$, pattern $P$, integer $k$;    Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all $k$-edit occ's of $P$ in $T$
 $\rightsquigarrow$ $k$-edit-aware lossy compression of $P$ and $T$ (think: filter + discarding filtered-out parts)

In [KoN**W**'24, KoN**W**'25] for $T$ and $P$ with $t \leq \frac{3}{2}\, p + k$:

◆ Show: any $k$-edit-aware lossy compr. of $P$ and $T$ needs $\Omega(k \log(t/k))$ bits

◆ Obtain $k$-edit-aware lossy compr. of $P$ and $T$ as $\tilde{O}(k)$ substring equations;
  computable via quantum algorithm using $\hat{O}(\sqrt{kp})$ queries/time

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   36-4

# If You Want to Be Fast, Take a Detour

⟨ **Approximate String Matching** ⟩       early 1980's

Given: text $T$, pattern $P$, integer $k$;    Find: (starting pos. of) substrings of $T$ at edit distance $\leq k$ to $P$.

- ◆ $P$/$T$ given as oracle ⤳ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all $k$-edit occ's of $P$ in $T$
   ⤳ $k$-edit-aware lossy compression of $P$ and $T$ (think: filter + discarding filtered-out parts)

In [KoN**W**'24, KoN**W**'25] for $T$ and $P$ with $t \leq \frac{3}{2}\, p + k$:

- ◆ Show: any $k$-edit-aware lossy compr. of $P$ and $T$ needs $\Omega(k \log(t/k))$ bits
- ◆ Obtain $k$-edit-aware lossy compr. of $P$ and $T$ as $\tilde{O}(k)$ substring equations;
  computable via quantum algorithm using $\hat{O}(\sqrt{kp})$ queries/time
- ◆ Give classical $\tilde{O}(b^2)$-time algo to compute sol. to $b$ substr. eqn's as $O(b)$-size grammar

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   36-5

⟨ **Approximate String Matching** ⟩        early 1980's

Given: text *T*, pattern *P*, integer *k*;    Find: (starting pos. of) substrings of *T* at edit distance ≤ *k* to *P*.

- ◆ *P*/*T* given as oracle ⤳ queries possible in superposition
- ◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$

Challenge: query as few char's as possible, but enough to recover all *k*-edit occ's of *P* in *T*
  ⤳ *k*-edit-aware lossy compression of *P* and *T* (think: filter + discarding filtered-out parts)

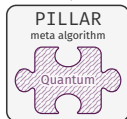In [KoN**W**'24, KoN**W**'25] for *T* and *P* with $t \leq \frac{3}{2} p + k$:

- ◆ Show: any *k*-edit-aware lossy compr. of *P* and *T* needs $\Omega(k \log(t/k))$ bits
- ◆ Obtain *k*-edit-aware lossy compr. of *P* and *T* as $\tilde{O}(k)$ substring equations; computable via quantum algorithm using $\hat{O}(\sqrt{kp})$ queries/time
- ◆ Give classical $\tilde{O}(b^2)$-time algo to compute sol. to *b* substr. eqn's as $O(b)$-size grammar
  ⤳ Can use PILLAR algorithm for fully-compressed setting on lossy compression!

![NII logo] Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   36-6

◆ $P/T$ given as oracle ⤳ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$



**Easy Result**

ASM on $T$, $P$ w/ $t \leq \frac{3}{2} p + k$ admits qu-algo
using $\tilde{O}(k^{3.5}\sqrt{n})$ queries/time

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    37-1

◆ $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
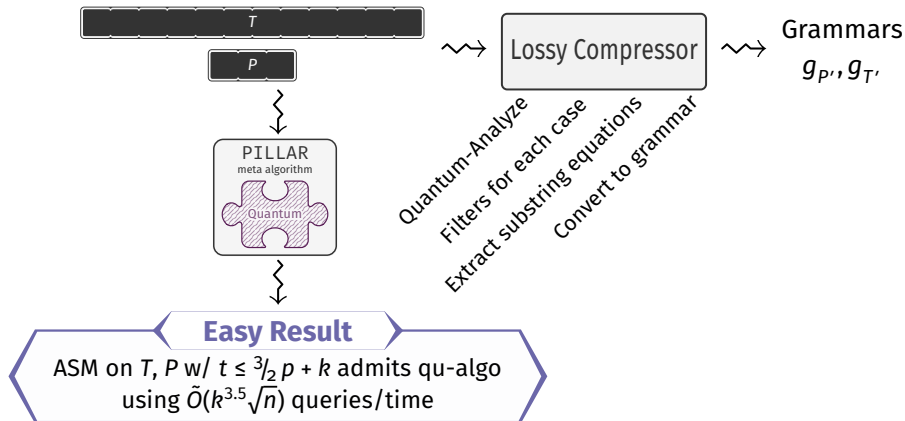◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$



**Easy Result**

ASM on $T$, $P$ w/ $t \leq \frac{3}{2} p + k$ admits qu-algo
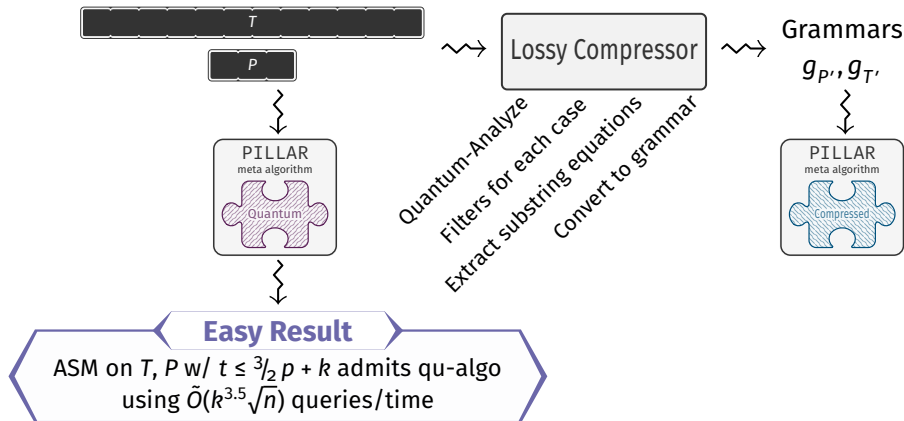using $\tilde{O}(k^{3.5}\sqrt{n})$ queries/time

- $P/T$ given as oracle $\rightsquigarrow$ queries possible in superposition
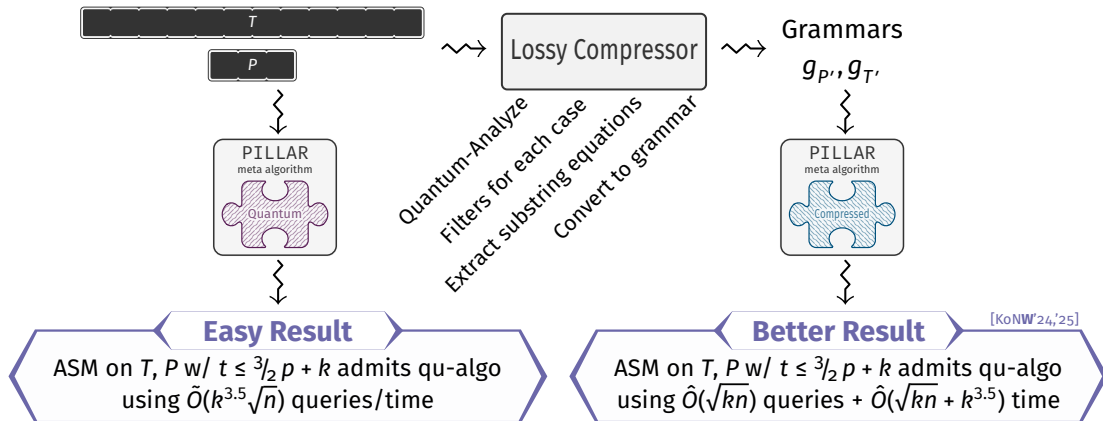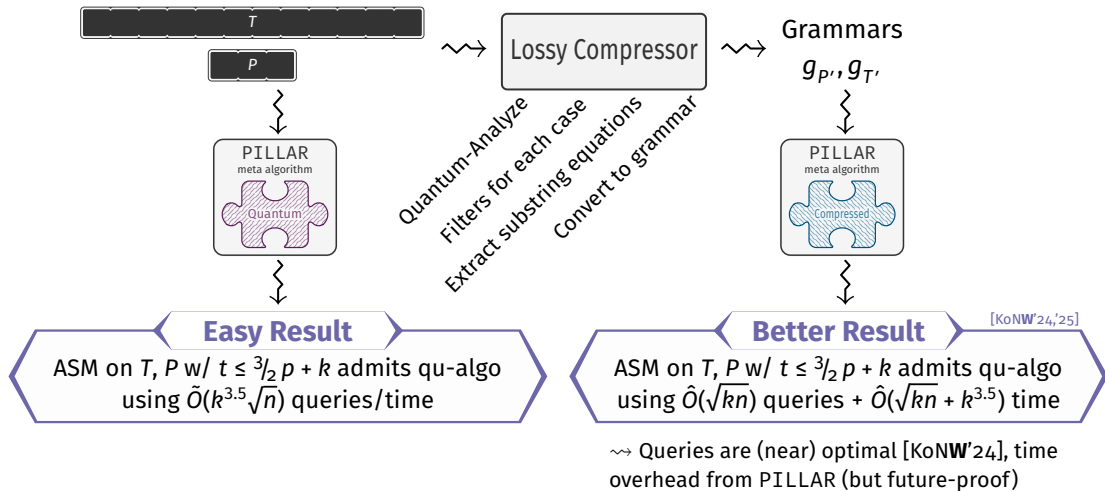- Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$



T

P

PILLAR
meta algorithm

Quantum

$\rightsquigarrow$ Lossy Compressor $\rightsquigarrow$ Grammars $g_{P'}, g_{T'}$

Quantum-Analyze
Filters for each case
Extract substring equations
Convert to grammar

PILLAR
meta algorithm

Compressed

**Easy Result**

ASM on $T$, $P$ w/ $t \leq \frac{3}{2} p + k$ admits qu-algo
using $\tilde{O}(k^{3.5}\sqrt{n})$ queries/time

Philip Wellnitz 37-3
Revitalizing Research on Approximate String Matching Algorithms

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

◆ $P/T$ given as oracle ⤳ queries possible in superposition
◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$



```
              T
            P                ⤳    Lossy Compressor    ⤳    Grammars
                                                            g_P', g_T'
            ⎰
       ┌─────────┐                                      ┌─────────┐
       │ PILLAR  │                                      │ PILLAR  │
       │meta algorithm│                                 │meta algorithm│
       │ Quantum │                                      │Compressed│
       └─────────┘                                      └─────────┘
            ⎰                                                ⎰
```

Quantum-Analyze
Filters for each case
Extract substring equations
Convert to grammar

**Easy Result**
ASM on $T$, $P$ w/ $t \leq {}^3\!/_2\, p + k$ admits qu-algo
using $\tilde{O}(k^{3.5}\sqrt{n})$ queries/time

**Better Result**                    [KoNW'24,'25]
ASM on $T$, $P$ w/ $t \leq {}^3\!/_2\, p + k$ admits qu-algo
using $\hat{O}(\sqrt{kn})$ queries + $\hat{O}(\sqrt{kn} + k^{3.5})$ time

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    37-4

◆ $P/T$ given as oracle ⤳ queries possible in superposition

◆ Goal: algorithms with query complexity $Q(n) \ll n$ and time complexity $T(n) \ll n$



**Easy Result**

ASM on $T$, $P$ w/ $t \leq \frac{3}{2} p + k$ admits qu-algo using $\tilde{O}(k^{3.5}\sqrt{n})$ queries/time

**Better Result** [KoNW'24,'25]

ASM on $T$, $P$ w/ $t \leq \frac{3}{2} p + k$ admits qu-algo using $\hat{O}(\sqrt{kn})$ queries + $\hat{O}(\sqrt{kn} + k^{3.5})$ time

⤳ Queries are (near) optimal [KoNW'24], time overhead from PILLAR (but future-proof)

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   37-5

Spotlight Extension: String Matching with Weighted Edits

How similar are two strings *X* and *Y*?

How similar are two strings *X* and *Y*?

How similar are two strings *X* and *Y*?

**Edit Distance** ED($X, Y$)

Min number of character insertions, deletions, and substitutions that transform $X$ to $Y$

**Edit Distance** ED($X, Y$)

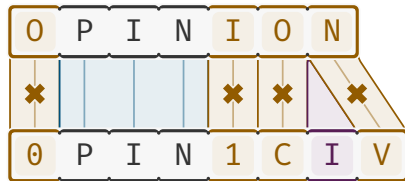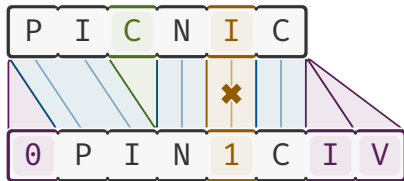Min number of character insertions, deletions, and substitutions that transform $X$ to $Y$

ED(0PIN1CIV, OPINION) = 5

Edit Distance $\qquad$ ED($X$, $Y$)

Min number of character insertions, deletions, and substitutions that transform $X$ to $Y$

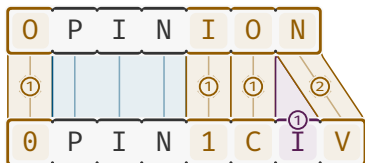ED(0PIN1CIV, OPINION) = 5

ED(0PIN1CIV, PICNIC) = 5

**Weighted Edit Distance** $\text{ED}^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:

◆ inserting $y$ costs $w(\varepsilon, y)$; ◆ deleting $x$ costs $w(x, \varepsilon)$;
◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(0, O) := 1 \quad w(1, I) := 1 \quad w(C, O) := 1 \quad\quad w(*, *) := 2 \quad w(*, \varepsilon) := 1 \quad w(\varepsilon, *) := 10$

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  41-1

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics
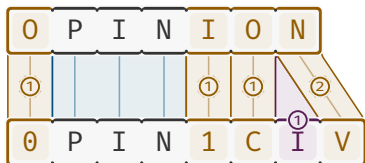
## Weighted Edit Distance $\text{ED}^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:
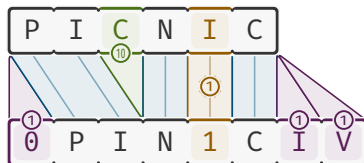
◆ inserting $y$ costs $w(\varepsilon, y)$;      ◆ deleting $x$ costs $w(x, \varepsilon)$;

◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(0, O) := 1$   $w(1, I) := 1$   $w(C, O) := 1$     $w(*, *) := 2$   $w(*, \varepsilon) := 1$   $w(\varepsilon, *) := 10$



$\text{ED}^w(\texttt{0PIN1CIV}, \texttt{OPINION}) = 6$

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   41-2

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics
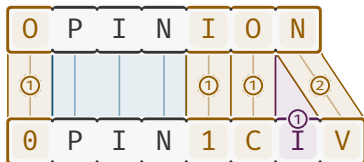
**Weighted Edit Distance** $\text{ED}^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:

◆ inserting $y$ costs $w(\varepsilon, y)$;     ◆ deleting $x$ costs $w(x, \varepsilon)$;

◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(\mathtt{0}, \mathtt{O}) := 1$   $w(\mathtt{1}, \mathtt{I}) := 1$   $w(\mathtt{C}, \mathtt{O}) := 1$     $w(*, *) := 2$   $w(*, \varepsilon) := 1$   $w(\varepsilon, *) := 10$

$\text{ED}^w(\mathtt{0PIN1CIV}, \mathtt{OPINION}) = 6$     $\text{ED}^w(\mathtt{0PIN1CIV}, \mathtt{PICNIC}) \leq 14$

## Weighted Edit Distance $ED^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:
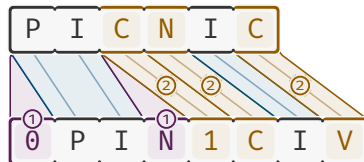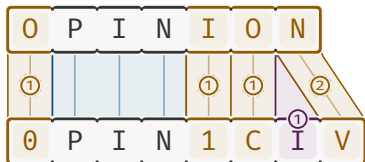
◆ inserting $y$ costs $w(\varepsilon, y)$;      ◆ deleting $x$ costs $w(x, \varepsilon)$;

◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(0, O) := 1 \quad w(1, I) := 1 \quad w(C, O) := 1 \qquad w(*, *) := 2 \quad w(*, \varepsilon) := 1 \quad w(\varepsilon, *) := 10$



$ED^w(\texttt{0PIN1CIV}, \texttt{OPINION}) = 6$



$ED^w(\texttt{0PIN1CIV}, \texttt{PICNIC}) = 8$

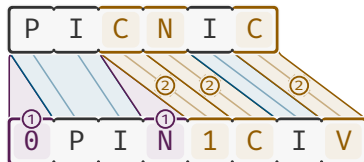**Weighted Edit Distance** $\text{ED}^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:

◆ inserting $y$ costs $w(\varepsilon, y)$;    ◆ deleting $x$ costs $w(x, \varepsilon)$;
◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(0, O) := 1$  $w(1, I) := 1$  $w(C, O) := 1$    $w(*, *) := 2$  $w(*, \varepsilon) := 1$  $w(\varepsilon, *) := 10$

$\text{ED}^w(\texttt{0PIN1CIV}, \texttt{OPINION}) = 6$    $\text{ED}^w(\texttt{0PIN1CIV}, \texttt{PICNIC}) = 8$

Justified Assumption:    $w$ is normalized,    $w(x, y) \geq 1$ for all $x \neq y$.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
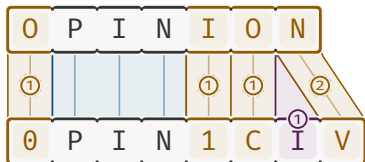Revitalizing Research on Approximate String Matching Algorithms    41-5
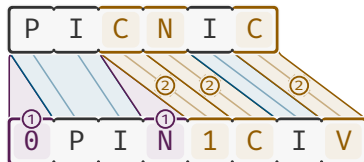
## Weighted Edit Distance

$ED^w(X, Y)$

Min cost of transforming $X$ to $Y$ using character edits, where:

◆ inserting $y$ costs $w(\varepsilon, y)$;     ◆ deleting $x$ costs $w(x, \varepsilon)$;

◆ substituting $x$ for $y$ costs $w(x, y)$.

$w(0, O) := 1$   $w(1, I) := 1$   $w(C, O) := 1$     $w(*, *) := 2$   $w(*, \varepsilon) := 1$   $w(\varepsilon, *) := 10$



$ED^w(\texttt{0PIN1CIV}, \texttt{OPINION}) = 6$     $ED^w(\texttt{0PIN1CIV}, \texttt{PICNIC}) = 8$

Justified Assumption:   $w$ is normalized,   $w(x, y) \geq 1$ for all $x \neq y$.

Otherwise: could scale weights and condition $ED^w_{\leq k}(X, Y) \leq k$ (e.g. in ASM) becomes meaningless.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   41-6

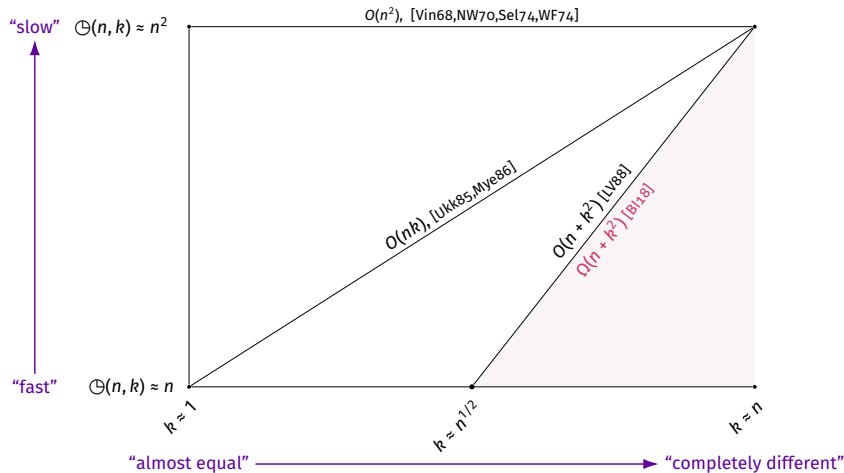Recall crucial subroutine for ASM:

Verify that occurrence starts at given (interval of) position of T

Recall crucial subroutine for ASM:

Verify that occurrence starts at given (interval of) position of T

⤳ Need to compute Bounded Weighted Edit Distance

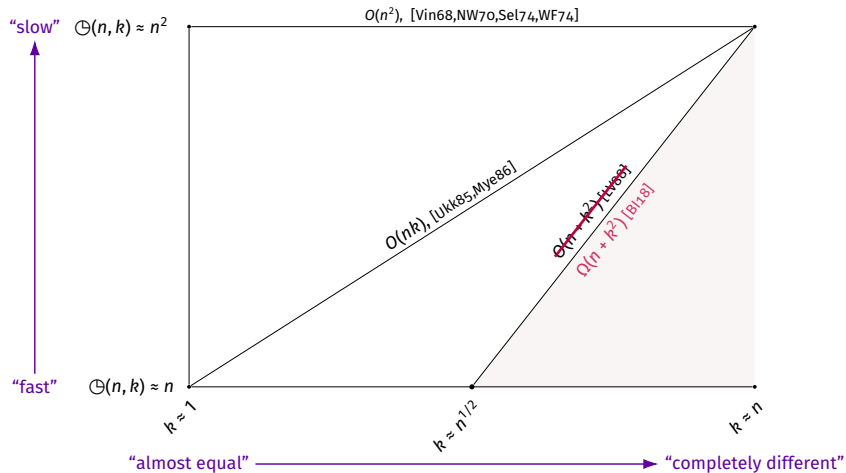Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  42-2

"slow" $\mathbb{O}(n,k) \approx n^2$

$O(n^2)$, [Vin68,NW70,Sel74,WF74]

$O(nk)$, [Ukk85,Mye86]

$O(n + k^2)$ [LV88]

$\Omega(n + k^2)$ [BI18]

"fast" $\mathbb{O}(n,k) \approx n$

$k \approx 1$

$k \approx n^{1/2}$

$k \approx n$

"almost equal" ⟶ "completely different"

Existing algorithms for Edit Distance ED(X, Y), where |X|, |Y| ≤ n

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    43

What about Weighted Edit Distance?

Existing algorithms for Weighted Edit Distance $ED^w(X, Y)$, where $|X|, |Y| \leq n$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
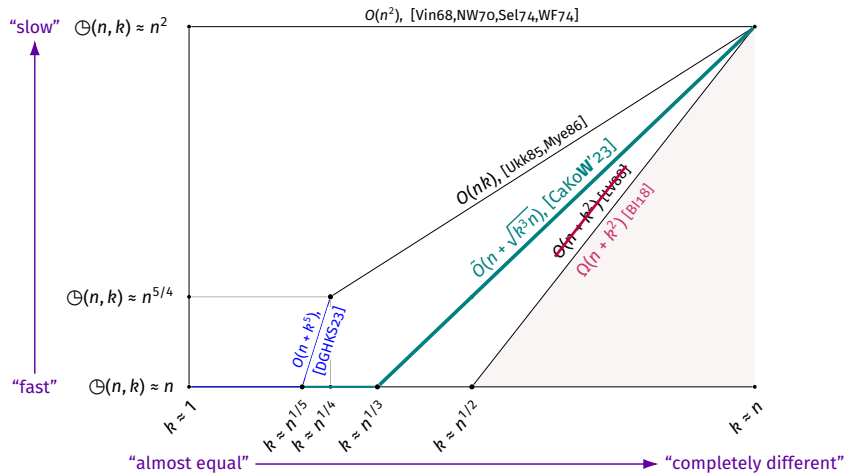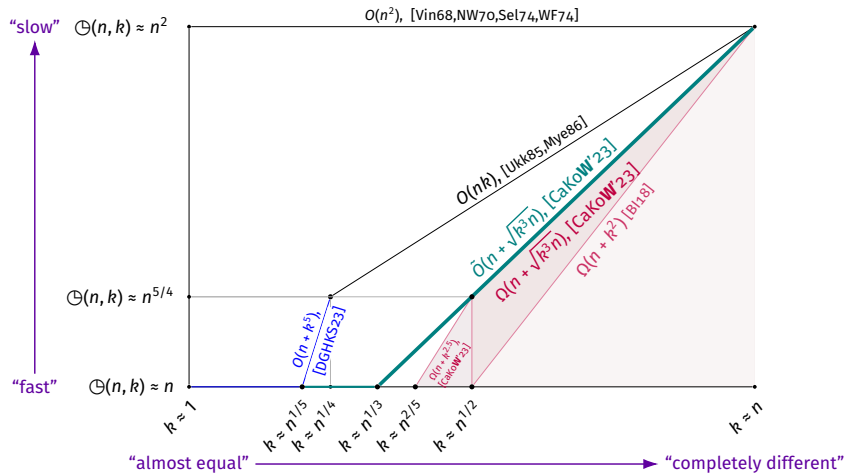Revitalizing Research on Approximate String Matching Algorithms    45-1

Existing algorithms for Weighted Edit Distance $ED^w(X, Y)$, where $|X|, |Y| \leq n$

Existing algorithms for Weighted Edit Distance $ED^w(X, Y)$, where $|X|, |Y| \leq n$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  45-3

Existing algorithms for Weighted Edit Distance $ED^w(X, Y)$, where $|X|, |Y| \leq n$

**Main Theorem (Upper Bound)** [CaKo**W**'23]

Strings $X$, $Y$ each of length at most $n$

Oracle access to (normalized) weight function $w$

Can compute $k = \text{ED}^w(X, Y)$ in time $O(n + \sqrt{nk^3} \log^3 n)$

[CaKo**W**'23]

## Main Theorem (Upper Bound)

Strings $X$, $Y$ each of length at most $n$

Oracle access to (normalized) weight function $w$

Can compute $k = \mathrm{ED}^w(X, Y)$ in time $O(n + \sqrt{n k^3} \log^3 n)$

[CaKo**W**'23]

## Main Theorem (Lower Bound)

Assuming the APSP Hypothesis and for $\sqrt{n} \le k \le n$,

Main Theorem (Upper Bound) is tight (up to $n^{o(1)}$-factors)

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 46-2

**Main Theorem (Upper Bound)** [CaKo**W**'23]

Strings $X$, $Y$ each of length at most $n$

Oracle access to (normalized) weight function $w$

Can compute $k = \text{ED}^w(X, Y)$ in time $O(n + \sqrt{nk^3} \log^3 n)$

**Main Theorem (Lower Bound)** [CaKo**W**'23]

Assuming the APSP Hypothesis and for $\sqrt{n} \le k \le n$,

Main Theorem (Upper Bound) is tight (up to $n^{o(1)}$-factors)

Can extend UB to verify $\le k$ consecutive positions

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    46-3

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

**Main Theorem (Upper Bound)**                                    [CaKo**W**'23]

Strings $X$, $Y$ each of length at most $n$

Oracle access to (normalized) weight function $w$

Can compute $k = \text{ED}^w(X, Y)$ in time $O(n + \sqrt{n}k^3 \log^3 n)$

**Main Theorem (Lower Bound)**                                    [CaKo**W**'23]

Assuming the APSP Hypothesis and for $\sqrt{n} \leq k \leq n$,

Main Theorem (Upper Bound) is tight (up to $n^{o(1)}$-factors)

Can extend UB to verify $\leq k$ consecutive positions

$\rightsquigarrow$ Good enough for $\tilde{O}(t + t/p \cdot k^4)$ algo for weighted ASM (long, boring) [ChKo**W**'25]

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   46-4

**Main Theorem (Upper Bound)** [CaKo**W**'23]

Strings $X$, $Y$ each of length at most $n$

Oracle access to (normalized) weight function $w$

Can compute $k = \text{ED}^w(X, Y)$ in time $O(n + \sqrt{nk^3} \log^3 n)$

**Main Theorem (Lower Bound)** [CaKo**W**'23]

Assuming the APSP Hypothesis and for $\sqrt{n} \le k \le n$,

Main Theorem (Upper Bound) is tight (up to $n^{o(1)}$-factors)

Can extend UB to verify $\le k$ consecutive positions

⤳ Good enough for $\tilde{O}(t + t/p \cdot k^4)$ algo for weighted ASM (long, boring) [ChKo**W**'25]

⤳ What about $\tilde{O}(kt)$-time for weighted ASM?

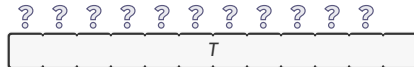With Standard Trick, UB yields $\tilde{O}(t/k \cdot \sqrt{pk^3}) = \tilde{O}(t\sqrt{pk})$

LB rules out $\tilde{O}(kt)$ via better WED algo ⤳ need different approach

NII Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 46-5

⟨**Weighted Approximate String Matching**⟩

Given: $T$, $P$, $k$, oracle to $w$;    Find: (starting pos. of) substr. of $T$ at weighted ED ≤ $k$ to $P$.
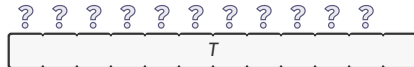
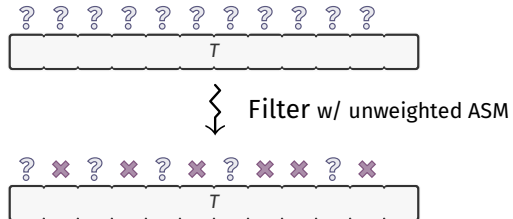Focus: Obtain starting positions of occurrences

⟨**Weighted Approximate String Matching**⟩

Given: $T$, $P$, $k$, oracle to $w$;    Find: (starting pos. of) substr. of $T$ at weighted ED ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)



Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  47-2

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

⟨**Weighted Approximate String Matching**⟩

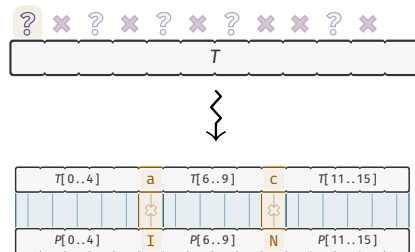Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
⤳ Use unweighted $O(kn)$-time ASM as Filter



Filter w/ unweighted ASM

National Institute of Informatics
Inter-University Research Institute Corporation
Research Organization of Information and Systems

⟨ **Weighted Approximate String Matching** ⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $\text{ED}^w(X, Y) \geq \text{ED}(X, Y)$ (by norm.)
⤳ Use unweighted $O(kn)$-time ASM as Filter

If $\text{ED}(P, T[a..b]) \leq k$, also get cheap alignment
$P \rightsquigarrow T[a..b]$ that matches exactly most of $P$

⟨ **Weighted Approximate String Matching** ⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

Focus: Obtain starting positions of occurrences
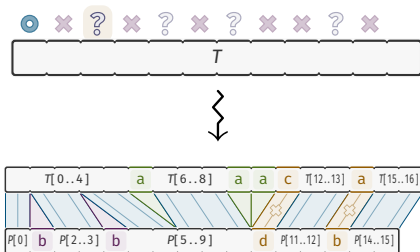
Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
⤳ Use unweighted $O(kn)$-time ASM as Filter

If $ED(P, T[a..b]) \leq k$, also get cheap alignment
$P \rightsquigarrow T[a..b]$ that matches exactly most of $P$

All such alignments very similar to each other
But how to exploit this fact?

⟨**Weighted Approximate String Matching**⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
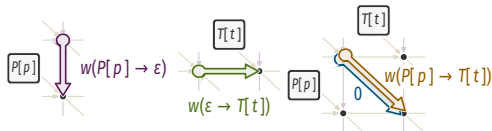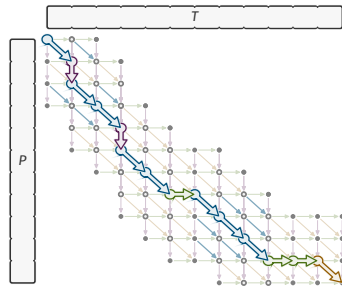⇝ Use unweighted $O(kn)$-time ASM as Filter

If $ED(P, T[a..b]) \leq k$, also get cheap alignment
$P \rightsquigarrow T[a..b]$ that matches exactly most of $P$

All such alignments very similar to each other
But how to exploit this fact?
⇝ Cast as shortest path prob in alignment graphs

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   47-6

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

⟨ **Weighted Approximate String Matching** ⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
⇝ Use unweighted $O(kn)$-time ASM as Filter

If $ED(P, T[\,a..b\,]) \leq k$, also get cheap alignment
$P \rightsquigarrow T[\,a..b\,]$ that matches exactly most of $P$

All such alignments very similar to each other
But how to exploit this fact?
⇝ Cast as shortest path prob in alignment graphs
⇝ Copy matched parts from (trivial) align. $P \rightsquigarrow P$

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

⟨ **Weighted Approximate String Matching** ⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED ≤ $k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
⤳ Use unweighted $O(kn)$-time ASM as Filter

If $ED(P, T[a..b]) \leq k$, also get cheap alignment
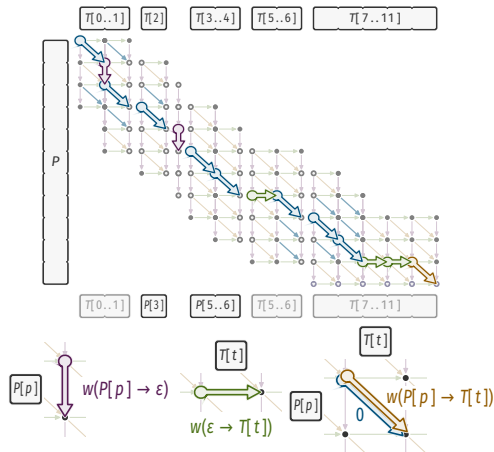$P \rightsquigarrow T[a..b]$ that matches exactly most of $P$

All such alignments very similar to each other
But how to exploit this fact?
⤳ Cast as shortest path prob in alignment graphs
⤳ Copy matched parts from (trivial) align. $P \rightsquigarrow P$

Comp. bound-to-bound dist for each block [Klein'05];
then combine using fast (min, +)-product [SMAWK]

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   47-8

⟨**Weighted Approximate String Matching**⟩

Given: $T$, $P$, $k$, oracle to $w$;   Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

Focus: Obtain starting positions of occurrences

Observation: $ED^w(X, Y) \geq ED(X, Y)$ (by norm.)
⤳ Use unweighted $O(kn)$-time ASM as Filter

If $ED(P, T[a..b]) \leq k$, also get cheap alignment
$P \rightsquigarrow T[a..b]$ that matches exactly most of $P$
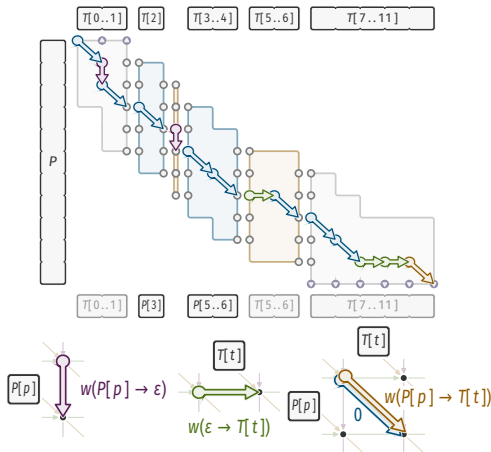
All such alignments very similar to each other
But how to exploit this fact?
⤳ Cast as shortest path prob in alignment graphs
⤳ Copy matched parts from (trivial) align. $P \rightsquigarrow P$

Comp. bound-to-bound dist for each block [Klein'05];
then combine using fast (min, +)-product [SMAWK]
Precompute all $P$ vs $P$ blocks [Klein'05];
⤳ Verify $\leq k$ pos in $\tilde{O}(k^2)$ after $\tilde{O}(kp)$ prep. ⤳ $\tilde{O}(kt)$ total

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   47-9

**⟨Weighted Approximate String Matching⟩**

Given: $T$, $P$, $k$, oracle to $w$;    Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

**⟨Main Theorem⟩** [ChKo**W**'25]

Weighted ASM is in time $\tilde{O}(t + t/p \cdot k^4)$ and in $O(k^4)$ PILLAR time for $t < \frac{3}{2} p + k$.

**⟨Main Theorem⟩** [ChKo**W**'25]

Weighted ASM is in time $\tilde{O}(kt)$.

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms   48-1

⟨ **Weighted Approximate String Matching** ⟩

Given: $T$, $P$, $k$, oracle to $w$;  Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

⟨ **Main Theorem** ⟩ [ChKo**W**'25]

Weighted ASM is in time $\tilde{O}(t + t/p \cdot k^4)$ and in $O(k^4)$ PILLAR time for $t < \tfrac{3}{2} p + k$.

⟨ **Main Theorem** ⟩ [ChKo**W**'25]

Weighted ASM is in time $\tilde{O}(kt)$.

Generalization of PILLAR ASM algorithm to weighted setting "easy".

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms  48-2

⟨**Weighted Approximate String Matching**⟩

Given: $T$, $P$, $k$, oracle to $w$;    Find: (starting pos. of) substr. of $T$ at weighted ED $\leq k$ to $P$.

**Main Theorem** [ChKoW'25]

Weighted ASM is in time $\tilde{O}(t + t/p \cdot k^4)$ and in $O(k^4)$ PILLAR time for $t < \frac{3}{2} p + k$.

**Main Theorem** [ChKoW'25]

Weighted ASM is in time $\tilde{O}(kt)$.

Generalization of PILLAR ASM algorithm to weighted setting "easy".

But, also interesting non-PILLAR algorithms out there.

**NII** Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms    48-3

- ◆ Big Question: Close gap between UB $O(n + k^{3.5})$ / LB $\Omega(n + k^2)$ for ASM.
  - ◆ An $\tilde{O}(k^2)$-time PILLAR algo would imply optimal algorithms in many other settings (e.g. quantum, compressed, etc.)
  - ◆ First step: $\tilde{O}(k^3)$-time PILLAR algo
  - ◆ Stronger (conditional) LB would require new techniques
  - ◆ Similar gap in weighted setting

◆ Big Question: Close gap between UB $O(n + k^{3.5})$ / LB $\Omega(n + k^2)$ for ASM.
- ◆ An $\tilde{O}(k^2)$-time `PILLAR` algo would imply optimal algorithms in many other settings (e.g. quantum, compressed, etc.)
- ◆ First step: $\tilde{O}(k^3)$-time `PILLAR` algo
- ◆ Stronger (conditional) LB would require new techniques
- ◆ Similar gap in weighted setting

◆ Big Question: Reporting substrings instead of starting positions
- ◆ Alignment graph-based methods seem to fare better
- ◆ First step: $k^{4-\varepsilon}$-time PILLAR algo for ASM
- ◆ First step: $\tilde{O}(kp)$-time for WASM

◆ Big Question: Close gap between UB $O(n + k^{3.5})$ / LB $\Omega(n + k^2)$ for ASM.
   ◆ An $\tilde{O}(k^2)$-time `PILLAR` algo would imply optimal algorithms in many other settings (e.g. quantum, compressed, etc.)
   ◆ First step: $\tilde{O}(k^3)$-time `PILLAR` algo
   ◆ Stronger (conditional) LB would require new techniques
   ◆ Similar gap in weighted setting

◆ Big Question: Reporting substrings instead of starting positions
   ◆ Alignment graph-based methods seem to fare better
   ◆ First step: $k^{4-\varepsilon}$-time PILLAR algo for ASM
   ◆ First step: $\tilde{O}(kp)$-time for WASM

◆ Big Question: Simpler algorithms with similar guarantees
   ◆ Especially almost periodic case highly involved
   ◆ Structural insights currently come with impractical constants

NII Inter-University Research Institute Corporation Research Organization of Information and Systems National Institute of Informatics

- ◆ Big Question: Close gap between UB $O(n + k^{3.5})$ / LB $\Omega(n + k^2)$ for ASM.
  - ◆ An $\tilde{O}(k^2)$-time PILLAR algo would imply optimal algorithms in many other settings (e.g. quantum, compressed, etc.)
  - ◆ First step: $\tilde{O}(k^3)$-time PILLAR algo
  - ◆ Stronger (conditional) LB would require new techniques
  - ◆ Similar gap in weighted setting

- ◆ Big Question: Reporting substrings instead of starting positions
  - ◆ Alignment graph-based methods seem to fare better
  - ◆ First step: $k^{4-\varepsilon}$-time PILLAR algo for ASM
  - ◆ First step: $\tilde{O}(kp)$-time for WASM

- ◆ Big Question: Simpler algorithms with similar guarantees
  - ◆ Especially almost periodic case highly involved
  - ◆ Structural insights currently come with impractical constants

- ◆ Big Question: Generalizing idea behind PILLAR to other problem families, settings

Inter-University Research Institute Corporation
Research Organization of Information and Systems
National Institute of Informatics

Philip Wellnitz
Revitalizing Research on Approximate String Matching Algorithms 49-4

# Navigation