

Edit Distance with Single-Symbol Combinations and Splits

Manolis Christodoulakis¹ and Gerhard Brey²

¹School of Computing & Technology,
University of East London
m.christodoulakis@uel.ac.uk

²Centre for Computing in the Humanities,
Kings College London
gerhard.brey@kcl.ac.uk

Prague Stringology Conference 2008

NCSE – Our Data Repository

- **Nineteenth-Century Serials Edition (NCSE)**

<http://www.ncse.ac.uk>

- Collaborative project
 - Birkbeck College, London
 - British Library
 - King's College London (Centre for Computing in the Humanities)
 - Olive Software
- Digital edition of 6 newspapers / periodicals published between 1806–1890
- Microfilmed versions of the newspapers were scanned and OCR'd by Olive Software

OCR Quality

- OCR quality varies widely across corpus

high intellectual standing. Such was to have been the colleague of Mr. DISRAELI. But his total defeat has drawn forth proofs of weakness ominous for the continuance of a Tory ascendancy. In fact, the Liberals are already calculating that they shall be able to replace two if not three of the members on the next election; and the state of the Registry, which is a perfect archaeological curiosity, will no doubt operate as a stimulus to renewed activity in registration throughout all counties as well as Buckinghamshire.

The financial aspects at the close of the year are indeed remarkable. We find the revenue exhibiting a decrease in comparison with 1856 of £,828,000, at the same time that the Banks of England and

high intellectual standing.' 'Such ^as " ,to jja^been the colleague 'of Mr.' Disraeli. But his tp^al, defeat has' drawn forth proofs of wcAknss ominous for the continuance of a Tory.as,o^jlanay.: In fact, the Liberals are already ealoidaWj! *^cf'nfnc y; shall be able to replace two if notHhreje ofi-hd ^e '^oi-s on the next election; and the state o^ttftfyiglfriy, ».wluoU-i&~a<~porfoot,w doubt oporato as a stimulus to rcnewc\l^tt#fcjlaK : _ registration throughout all counties"" aV^mUEfea? Buckinghamshire. ^^f^mwf) The financial aspects at the clpso oJ^1' i^|y|g^W , indeed remarkable. We find the roYOXJMjigi fffbflGngj1' u decrease in comparison with 1856 °S» ^^ rp|0^3 at the same time that tho Banks of ^ng1^an^

The Heart of the Problem

Approximate Pattern Matching

How to distinguish good approximate matches from random ones?

The Heart of the Problem

Approximate Pattern Matching

How to distinguish good approximate matches from random ones?

Example

Let the pattern be Billington.

The Heart of the Problem

Approximate Pattern Matching

How to distinguish good approximate matches from random ones?

Example

Let the pattern be Billington.

Matches	Edit Distance	True match?
Billington	2	Must be ✓

The Heart of the Problem

Approximate Pattern Matching

How to distinguish good approximate matches from random ones?

Example

Let the pattern be Billington.

Matches	Edit Distance	True match?
Billington	2	Must be ✓
Wellington	2	Not really ✗

The Heart of the Problem

Approximate Pattern Matching

How to distinguish good approximate matches from random ones?

Example

Let the pattern be Billington.

Matches	Edit Distance	True match?
Billmington	2	Must be ✓
Wellington	2	Not really ✗

Remarks

- **m** optically resembles in
- **m** is essentially a merge of **i** and **n**

Problem Requirements

Ideas

An enhanced edit-distance algorithm that allows:

Problem Requirements

Ideas

An enhanced edit-distance algorithm that allows:

- symbols that **look similar** to get smaller penalty

Problem Requirements

Ideas

An enhanced edit-distance algorithm that allows:

- symbols that **look similar** to get smaller penalty
- groups of symbols to be **combined** and matched against a symbol that looks similar to the combination

Problem Requirements

Ideas

An enhanced edit-distance algorithm that allows:

- symbols that **look similar** to get smaller penalty
- groups of symbols to be **combined** and matched against a symbol that looks similar to the combination

We introduce a new edit-distance operation:

Definition

The **single-symbol combination** (hereafter: a combination) operation allows a string, e.g. “in” to be matched against a single symbol, e.g. “m”. The dual of a combination is the **split** operation.

Combination Operation in DP

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Combination Operation in DP

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

Combination Operation

Find the suffix (if any) of $x[1..i]$ that can be combined into symbol $y[j]$ with the **minimum** cost.

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Combination Operation in DP

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

Combination Operation

Find the suffix (if any) of $x[1..i]$ that can be combined into symbol $y[j]$ with the **minimum** cost.

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Combination Operation in DP

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

Combination Operation

Find the suffix (if any) of $x[1..i]$ that can be combined into symbol $y[j]$ with the **minimum** cost.

Problem!

$O(n)$ suffixes per cell!

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Combination Operation in DP

	0	1	2	3	4	5	6	7	8	9	
		B	i	l	l	m	g	t	o	n	
0	0	10	20	30	40	50	60	70	80	90	
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

Combination Operation

Find the suffix (if any) of $x[1..i]$ that can be combined into symbol $y[j]$ with the **minimum** cost.

Problem!

$O(n)$ suffixes per cell!

Idea

Preprocess the valid combinations.

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Combination Lists

Assumptions

For every symbol α , a (possibly empty) list is provided, call it C_α (the *combination list* of α), that contains all the strings that resemble α .

Combination Lists

Assumptions

For every symbol α , a (possibly empty) list is provided, call it C_α (the *combination list* of α), that contains all the strings that resemble α .

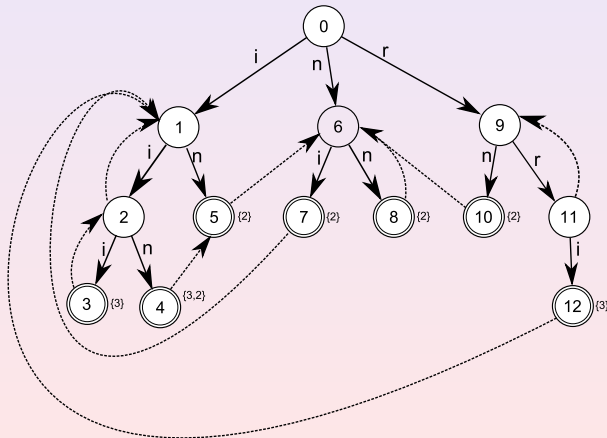
Preprocessing Algorithm

For each symbol α

- 1 Build **Aho-Corasick** automaton, T_α , from the list C_α
- 2 Replace the **output** values of final nodes with the lengths of the strings matched at those nodes

Combination Lists

T_m built from the combination list $C_m = \{iii, iin, in, ni, nn, rn, rri\}$



Incorporating the Combination Lists

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Incorporating the Combination Lists

Remark 1

For cell $D(6, 5)$ we need to spell out "Billin" on T_m

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Incorporating the Combination Lists

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

Remark 1

For cell $D(6, 5)$ we need to spell out “Billin” on T_m

Remark 2

For cell $D(7, 5)$ we need to spell out “Billing” on T_m (only one letter longer than the previous)

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Incorporating the Combination Lists

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Remark 1

For cell $D(6, 5)$ we need to spell out "Billin" on T_m

Remark 2

For cell $D(7, 5)$ we need to spell out "Billing" on T_m (only one letter longer than the previous)

Idea

Store a pointer to the node of T_m where "Billin" ended. When processing $D(7, 5)$ proceed from that node following the letter 'g' only.

Incorporating the Combination Lists

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

	1	2	3	4	5	6	7	8	9
	B	i	l	l	m	g	t	o	n
	root TB	root Ti	root Tl	root Tl	root Tm	root Tg	root Tt	root To	root Tn

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Incorporating the Combination Lists

	0	1	2	3	4	5	6	7	8	9	
		B	i	l	l	m	g	t	o	n	
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

1	2	3	4	5	6	7	8	9
B	i	l	l	m	g	t	o	n
"B" TB	"B" Ti	"B" Tl	"B" Tl	"B" Tm	"B" Tg	"B" Tt	"B" To	"B" Tn

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Incorporating the Combination Lists

		0	1	2	3	4	5	6	7	8	9
			B	i	l	l	m	g	t	o	n
0		0	10	20	30	40	50	60	70	80	90
1	B	10	0	10	20	30	40	50	60	70	80
2	i	20	10	0	10	20	30	40	50	60	70
3	l	30	20	10	0	10	20	30	40	50	60
4	l	40	30	20	10	0	10	20	30	40	50
5	i	50	40	30	20	10	10	20	30	40	50
6	n	60	50	40	30	20					
7	g										
8	t										
9	o										
10	n										

1	2	3	4	5	6	7	8	9
B	i	l	l	m	g	t	o	n
"Bi"	"Bi"	"Bi"	"Bi"	"Bi"	"Bi"	"Bi"	"Bi"	"Bi"
TB	Ti	Tl	Tl	Tm	Tg	Tt	To	Tn

$$d_{sub} = 10, d_{indel} = 10, d_{comb} = 1$$

Running Time and Space

Preprocessing Time

$$O\left(\sum_{\forall \alpha} l_{\alpha}\right) \quad \text{where } l_{\alpha} = \sum_{x \in C_{\alpha}} |x|$$

Running Time and Space

Preprocessing Time

$$O\left(\sum_{\forall\alpha} l_{\alpha}\right) \quad \text{where } l_{\alpha} = \sum_{x \in C_{\alpha}} |x|$$

Edit Distance Algorithm Time

$$O(mnk) \quad \text{where } k \text{ is an upper bound on } |C_{\alpha}| \text{ over all } \alpha$$

Running Time and Space

Preprocessing Time

$$O\left(\sum_{\forall\alpha} l_{\alpha}\right) \quad \text{where } l_{\alpha} = \sum_{x \in C_{\alpha}} |x|$$

Edit Distance Algorithm Time

$$O(mnk) \quad \text{where } k \text{ is an upper bound on } |C_{\alpha}| \text{ over all } \alpha$$

Space Requirements

$$O\left(mn + \sum_{\forall\alpha} l_{\alpha}\right)$$

Open Problems

Current Algorithm

Single-symbol combinations allow a string (e.g. “in”) to be matched against a *single* symbol (e.g. “m”).

Open Problems

Current Algorithm

Single-symbol combinations allow a string (e.g. “in”) to be matched against a *single* symbol (e.g. “m”).

Definition (Edit Distance with Re-Combinations)

How about allowing a string (e.g. “in”) to be matched against another string (e.g. “ni”)?

Open Problems

Current Algorithm

Single-symbol combinations allow a string (e.g. “in”) to be matched against a *single* symbol (e.g. “m”).

Definition (Edit Distance with Re-Combinations)

How about allowing a string (e.g. “in”) to be matched against another string (e.g. “ni”)?

Definition (Edit Distance with Transitive Combinations)

Let “in” $\in C_m$ and “ri” $\in C_n$. Write an algorithm that can infer that “iri” resembles “m” despite the fact that “iri” $\notin C_m$.