# A Simple Algorithm for Transposition-Invariant Amplified (δ, γ)-Matching

**Inbok LEE**[†a)], *Member*

**SUMMARY**    Approximate pattern matching plays an important role in various applications. In this paper we focus on $(\delta, \gamma)$-matching, where a character can differ at most $\delta$ and the sum of these errors is smaller than $\gamma$. We show how to find these matches when the pattern is transformed by $y = \alpha x + \beta$, without knowing $\alpha$ and $\beta$ in advance.
***key words:*** *combinatorics, pattern matching, fast Fourier transform*

## 1.    Introduction

Approximate pattern matching plays an important role in various applications, such as bioinformatics, computer-aided music analysis and computer vision where the pattern does not appear exactly but within small differences.

Let $T$ and $P$ be strings over a positive integer alphabet $\Sigma$. $T[i]$ denotes the $i$-th character of $T$. $T[i, j]$ denotes the substring $T[i]T[i+1]\cdots T[j]$. We focus on $(\delta, \gamma)$-matching which is defined as follows.

**Definition 1:** Given a text $T = T[1, n]$, a pattern $P = P[1, m]$, and two integer parameters $\delta$ and $\gamma$, $(\delta, \gamma)$-matching refers to the problem of finding all the substrings $T[i, i+m-1]$ satisfying two conditions.

- $\forall 1 \le j \le m, |T[i+j-1] - P[j]| \le \delta$ ($\delta$-matching).
- $\sum_{j=1}^{m} |T[i+j-1] - P[j]| \le \gamma$ ($\gamma$-matching).

Usually the size of alphabet $|\Sigma|$ (the number of elements in $\Sigma$) is large and the value of $\delta$ is small.

In addition, we consider *transposition-invariant amplified matching* where each character of $P$ is multiplied by an arbitrary integer $\alpha$ (amplified) and added by another integer $\beta$ (transposition-invariant).

**Definition 2:** Transposition-invariant amplified $(\delta, \gamma)$ matching refers to the problem of finding all the substrings $T[i, i+m-1]$ satisfying two conditions with two integers $\alpha$ and $\beta$ which are not known in advance.

- $\forall 1 \le j \le m, |T[i+j-1] - (\alpha P[j] + \beta)| \le \delta$ ($\delta$-matching).
- $\sum_{j=1}^{m} |T[i+j-1] - (\alpha P[j] + \beta)| \le \gamma$ ($\gamma$-matching).

Figure 1 shows an example with $\delta = 1$ and $\gamma = 2$. The original pattern $P$ is $(1, 3, 2, 1, 2)$ in (a). In (b), $T = (2, 5, 4, 3, 4)$. We can find an occurrence of $(1, 2)$-matching
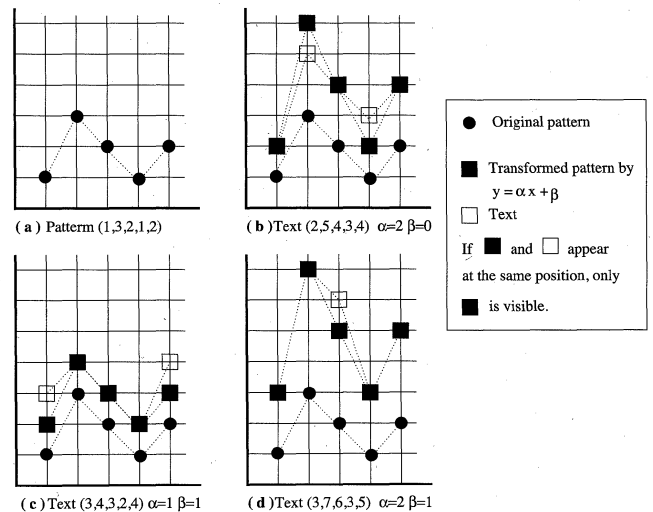
**Fig. 1**    (a) The original pattern, (b) amplified, (c) transposition-invariant, and (d) transposition-invariant amplified occurrence of the pattern.

if $P$ is transformed into $(2, 6, 4, 2, 4)$ by $y = 2x$. Similarly, in (c), $T = (3, 4, 3, 2, 4)$ and by $y = x + 1$, we get an occurrence of $(1, 2)$-matching of $(2, 4, 3, 2, 3)$. Finally, $T = (3, 7, 6, 3, 5)$ and we can find $(1, 2)$-matching of $(3, 7, 5, 3, 5)$ by $y = 2x + 1$ in (d). Note that we know just $P, T, \delta$, and $\gamma$. We need to determine $\alpha$ and $\beta$ if such $(\delta, \gamma)$-matching by some transform $y = \alpha x + \beta$ exists.

This problem arises from computer-aided music analysis. A simple motif (short melody) in music can evolve into different variations by changing the frequency or duration of each note in the motif. In these variations, each character $x$ (either frequency or duration) is transformed by a linear equation $y = \alpha x + \beta$.

$(\delta, \gamma)$-matching can be solved using the Fast Fourier Transform (FFT, [1]) in $O(\delta n \log m + occ \cdot m)$ time where $occ$ is the number of $\delta$-matches of $P$ in $T$. First, we find $\delta$-matching of $P$. The key to our solution is that all the inner-products between $P[1, m]$ and $T[i, i+m-1]$ ($1 \le i \le n-m+1$)

$$P[1, m] \cdot T[i, i+m-1] = \sum_{j=1}^{m} P[j]T[i+j-1]$$

can be calculated in $O(n \log m)$ time [2, Chap 32]. If there is an exact match between $P[1, m]$ and $T[i, i+m-1]$,

$$\sum_{j=1}^{m} (P[j] - T[i+j-1])^2$$

$$= \sum_{j=1}^{m} P[j]^2 - 2P[1,m] \cdot T[i, i+m-1]$$

$$+ \sum_{j=1}^{m} T[i+j-1]^2$$

should be zero. The total time complexity is $O(n \log m)$ because the first and last terms can be computed in $O(n+m)$ time.

To solve the $\delta$-matching problem, we need a function $g(x,y)$ where $\sum_{j=1}^{m} g(P[j], T[i+j-1]) = 0$ if and only if there is a $\delta$-matching between $P[1,m]$ and $T[i, i+m-1]$. We briefly explain how to define such $g(x,y)$ when $\delta = 1$. Let $g(x,y) = (x-y)^2 + 0.5 \times (-1)^{x+y} - 0.5$. It is easy to show that $g(x,y) = 0$ if $|x-y| \le 1$ and $g(x,y) > 0$ otherwise. To compute the second term, we define two strings $\epsilon_T$ and $\epsilon_P$: $\epsilon_T[i] = 1$ if $T[i]$ is even and $\epsilon_T[i] = -1$ if $T[i]$ is odd. $\epsilon_P$ can be defined similarly. Then, using $P[1,m] \cdot T[i, i+m-1]$ and $\epsilon_P[1,m] \cdot \epsilon_T[i, i+m-1]$, we can calculate $\sum_{j=1}^{m} g(P[j], T[i+j-1])$ in $O(n \log m)$ time for all $1 \le i \le n - m + 1$. The idea can be extended to the general case where $\delta > 1$. In that case the time complexity is $O(\delta n \log m)$. We skip the details of the general case here. Interested readers are directed to [1, pages 70–74].

For the transposition-invariant matching, several algorithms have been proposed recently in [3]–[5] based on sparse dynamic programming, which is not easy to understand and implement.

## 2. Algorithms

Our aim is to develop a simple and efficient algorithm for finding transposition-invariant amplified $(\delta, \gamma)$-matches. To do so, we will show how to modify the original FFT matching algorithm in [1].

We first find occurrences of transposition-invariant amplified $\delta$-matches. If $T[i, i+m-1]$ is a $\delta$-match of $P$ by some linear transformation $y = \alpha x + \beta$, we can check whether they are also $\gamma$-matches using the technique in [5]. What we need to know is just $\alpha$. If there exists such a transform $y = \alpha x + \beta$, $\beta$ should minimize the value of $\sum_{j=1}^{m} |T[i+j-1] - (\alpha P[j] + \beta)|$. For $1 \le j \le m$, we compute $T[i+j-1] - \alpha P[j]$. From these $m$ values, we choose the median $T[i+j'-1] - \alpha P[j'] (1 \le j \le m)$ in $O(m)$ time. Then we get $\beta = -T[i+j'-1] + \alpha P[j']$ which minimizes the sum of differences. Using this $\beta$, we check whether it is also a $\gamma$-match of $P$ in $O(m)$ time.

Now our problem is to find $\delta$-matches of $P$ after the transformation $y = \alpha x + \beta$. We first explain how to solve amplified $\delta$-mating (by $y = \alpha x$) and move to transposition-invariant amplified matching (by $y = \alpha x + \beta$).

### 2.1 Amplified $\delta$-Matching

Now we want to find occurrences of $\delta$-matches of amplified $P$ where every character of $P$ is multiplied by an integer $\alpha$. If $T[i, i+m-1]$ is such an occurrence, it should satisfy

- $\forall 1 \le j \le m, |T[i+j-1] - \alpha \cdot P[j]| \le \delta$, and
- $\sum_{1 \le j \le m} |T[i+j-1] - \alpha \cdot P[j]| \le \gamma$.

The problem is that we do not know $\alpha$ in advance. Once we know the exact value of $\alpha$, the following equation can be computed in $O(n \log m)$ time and its value should be zero when there is an exact match.

$$\sum_{j=1}^{m} (\alpha \cdot P[j] - T[i+j-1])^2$$

$$= \alpha^2 \sum_{j=1}^{m} P[j]^2 - 2\alpha P[1,m] \cdot T[i, i+m-1]$$

$$+ \sum_{j=1}^{m} T[i+j-1]^2.$$

To find $\delta$-matches, we do the same as we did in the original $\delta$-matching problem in $O(\delta n \log m)$ time. The difference is that now we use $\alpha P[j]$ instead of $P[j]$.

We do not know the exact value of $\alpha$ in advance. Therefore, for each substring $T[i, i+m-1]$ $(1 \le i \le n - m + 1)$, we create an integer array $\alpha[1, n-m+1]$ and store the candidate of $\alpha$ between $P[1,m]$ and $T[i, i+m-1]$ at $\alpha[i]$. To compute $\alpha[i]$, we first select a base element $P[k]$. For simplicity, assume that $P[k]$ is the greatest in $P[1,m]$. Then

$$\alpha[i] = \left\lceil \frac{T[i+k-1]}{P[k]} \right\rceil (1 \le i \le n - m + 1).$$

The base element $P[k]$ should meet one condition. From the above equation, $\alpha[i] - 0.5 \le T[i+k-1]/P[k] < \alpha[i] + 0.5$. Also, if there is a $\delta$-match between $P[1,m]$ and $T[i, i+m-1]$, $|T[i+k-1] - \alpha[i] \cdot P[k]| \le \delta$. We obtain

$$\alpha[i] \cdot P[k] - \delta \le T[i+k-1] \le \alpha[i] \cdot P[k] + \delta$$

and by eliminating $T[i+k-1]$, we obtain

$$\alpha[i] - 0.5 \le \frac{\alpha[i] \cdot P[k] - \delta}{P[k]} < \alpha[i] + 0.5 \text{ and}$$

$$\alpha[i] - 0.5 \le \frac{\alpha[i] \cdot P[k] + \delta}{P[k]} < \alpha[i] + 0.5.$$

After some calculation using $\delta \ge 0$, we get $P[k] > 2\delta$: at least one character in $P$ should be greater than $2\delta$. In real applications this condition can be met easily.

**Theorem 1:** The amplified $(\delta, \gamma)$-matching can be solved in $O(\delta n \log m + occ \cdot m)$ time, where $occ$ is the number of candidates.

*Proof.* Computing the array $\alpha[1, n-m+1]$ takes $O(n)$ time. The FFT runs in $O(\delta n \log m)$ time. After finding $occ$ $\delta$-matches (candidates), each of them requires $O(m)$ time verification. □

### 2.2 Transposition-Invariant Amplified $\delta$-Matching

We show a simple algorithm without using sparse-dynamic programming. We create two new strings $T' = T'[1, n-1]$

and $P' = P'[1, m - 1]$ such that $T'[i] = T[i + 1] - T[i]$ and $P'[i] = P[i + 1] - P[i]$. Then the following simple lemma holds.

**Lemma 1:** If there is a $(\delta, \gamma)$-matching between $P[1, m]$ and $T[i, i + m - 1]$, then there is a $(2\delta, 2\gamma)$-matching between $P'[1, m - 1]$ and $T'[i, i + m - 2]$.

*Proof.* We first prove $2\delta$-matching part. If there is an occurrence of $\delta$-matching at position $i$ of $T$, for $1 \le j \le m - 1$, it is evident that

$$-\delta \le T[i + j - 1] - P[j] \le \delta \text{ and}$$
$$-\delta \le T[i + j] - P[j + 1] \le \delta.$$

It follows that

$$-2\delta \le (T[i + j] - T[i + j - 1]) - (P[j + 1] - P[j])$$
$$= T'[i + j - 1] - P'[j] \le 2\delta.$$

Now we prove $2\gamma$-matching part. If there is $\gamma$-matching between $P[1, m]$ and $T[i, i+m-1]$, $\sum_{j=1}^{m} |T[i+j-1] - P[j]| \le \gamma$. By using the simple fact $|A| + |B| \ge |A + B|$, we obtain

$$\sum_{j=1}^{m-1} |T'[i + j - 1] - P'[j]|$$

$$= \sum_{j=1}^{m-1} |(T[i + j] - T[i + j - 1]) - (P[i + 1] - P[i])|$$

$$= \sum_{j=1}^{m-1} |(T[i + j] - P[i + 1]) + (P[i] - T[i + j - 1])|$$

$$\le \sum_{j=1}^{m-1} (|(T[i + j] - P[i + 1])| + |P[i] - T[i + j - 1])|)$$

$$\le 2 \cdot \sum_{i=1}^{m} |T[i + j - 1] - P[i]| \le 2\gamma.$$

$\square$

For simplicity, we used the basic $(\delta, \gamma)$-matching. Once we compute $\alpha$ array from $T'$ and $P'$, amplified $(\delta, \gamma)$-matching can be solved straightforward. Using this fact, we find occurrences of $(2\delta, 2\gamma)$-matching of $P'$ from $T'$. The results are candidates for $(\delta, \gamma)$-matching of $P$ from $T$. Then we check whether they are real $(\delta, \gamma)$-matches or not.

**Theorem 2:** The transposition-invariant $(\delta, \gamma)$-matching can be solved in $O(\delta n \log m + occ \cdot m)$ time, where $occ$ is the number of candidates.

*Proof.* Computing $T'$ and $P'$ takes $O(n)$ time $(n > m)$. The FFT runs in $O(\delta n \log m)$ time. Again, verifying each of $occ$ candidates takes $O(m)$ time. $\square$

Now we consider the size of $occ$. If $T$ and $P$ are drawn randomly from $\Sigma$, it is easy to show that the probability that $T'[i]$ and $P'[j]$ can have a $2\delta$-matching is $(4\delta+1)/|\Sigma|$. Hence, the probability is $((4\delta + 1)/|\Sigma|)^{m-1}$. The expected number of candidates is $n((4\delta + 1)/|\Sigma|)^{m-1}$, which is small when $\delta$ is small and $|\Sigma|$ is large.

## 3. Conclusion

We showed a simple $O(\delta n \log m + occ \cdot m)$ time algorithm for transposition-invariant amplified $(\delta, \gamma)$-matching. Its space complexity is $O(n)$ [2]. It is an improvement over [5] which requires $O(mn)$ time and space. Furthermore, we also consider the amplified $(\delta, \gamma)$-matching. The results in [3], [4] cannot be compared directly because their problem is transposition-invariant approximate pattern matching under edit distance. The algorithm in [3] requires $O(n + d^3)$ time and space where $d$ is the maximal edit distance between $T$ and $P$. The one in [4] runs in $O((mn + \log |\Sigma|)|\Sigma|)$ time and space.

The merit of our algorithm lies in its simplicity. Using the FFT libraries available, it is also easy to implement. Further research includes finding $(\delta, \gamma)$-matches of $P$ after more complex transformations.

**Acknowledgment**

**References**

[1] P. Clifford, R. Clifford, and C.S. Iliopoulos, "Faster algorithms for $(\delta, \gamma)$-matching and related problems," Proc. 16th Combinatorial Pattern Matching (CPM '05), pp.68–78, 2005.

[2] T.H. Cormen, C.E. Leiserson, and R.L. Rivest, Introduction to Algorithms, MIT Press, 1990.

[3] H. Hyrrö, "Restricted transposition invariant approximate string matching," Proc. 12th String Processing and Information Retrieval (SPIRE '05), pp.257–267, 2005.

[4] K. Lemström, G. Navarro, and Y. Pinzón, "Practical algorithms for transposition-invariant string-matching," J. Discrete Algorithms, vol.3, no.2-4, pp.267–292, 2005.

[5] V. Mäkinen, G. Navarro, and E. Ukkonen, "Transposition invariant string matching," J. Algorithms, vol.56, no.2, pp.124–153, 2005.