

A Formal Analysis of Edit Distance

M. D. Vose

February 1, 2004

1 Overview

The edit distance is a means of measuring the cost of transforming one string into another by way of insert, delete and change operators. The theoretical and computational underpinnings of edit distance are discussed in a formal mathematical framework.

2 Preliminaries

Let Σ be a finite nonempty set of symbols. The set of all finite strings with elements from Σ (including the empty string) is Σ^* . Strings are naturally functions; the string $s = s_0s_1 \dots s_k$ is identified with the set of ordered pairs (i.e., function)

$$s = \{(0, s_0), (1, s_1), \dots, (k, s_k)\}$$

The domain and range of function f , denoted by D_f and R_f respectively, are defined by

$$D_f = \{x : \exists y. (x, y) \in f\}$$

$$R_f = \{y : \exists x. (x, y) \in f\}$$

The notation

$$f : A \longrightarrow B$$

asserts that f is a function whose domain is a subset of A and whose range is a subset of B . The set B is referred to as the *co-domain*, the set A is referred to as the *given-domain*. If $D_f = A$, then f is called a *function from A to B* . If $D_f \neq A$, then f is called a *partial function from A to B* ; for example, every string $s \in \Sigma^*$ is a partial function $s : \mathbb{N} \longrightarrow \Sigma$, where \mathbb{N} is the set of nonnegative integers.

The length $|A|$ of string A refers to the cardinality of (the set of ordered pairs) A . In particular, the empty string \emptyset has length $|\emptyset| = 0$.

The discussion above explains how any given string s is a function, which in turn is a partial function—call it $\Phi(s)$ —from \mathbb{N} to Σ . Given string $s = abc$, for example, $\Phi(s) = \{(0, a), (1, b), (2, c)\}$ is the naturally associated partial function. Consider reversing this process: let

$$f = \{(i_0, s_0), (i_1, s_1), \dots, (i_k, s_k)\}$$

be a partial function from \mathbb{N} to Σ , where $i_0 < i_1 < \dots < i_k$ comprise the domain of f . It is naturally a string, namely $s_0 s_1 \dots s_k$; let $\Psi(f)$ denote this string. Passing from this string to a function in the natural way yields

$$\Phi(\Psi(f)) = \{(0, s_0), (1, s_1), \dots, (k, s_k)\}$$

which is not necessarily the same partial function f as was begun with. The partial function $\Phi(\Psi(f))$ is called the *normal form* of f . A partial function $f : \mathbb{N} \rightarrow \Sigma$ for which $f = \Phi(\Psi(f))$ is said to be *normalized*. The following properties are easily verified

1. $\Psi(\Phi(s)) = s$, for every string $s \in \Sigma^*$.
2. $\Phi(\Psi(f)) = f$, for every normalized partial function $f : \mathbb{N} \rightarrow \Sigma$.

In particular, $\Phi(s)$ is normalized for every $s \in \Sigma^*$, and the functions Φ and Ψ formally describe how strings and normalized partial functions can be thought of interchangeably. *This paper, however, generalizes the concept of string from normalized partial functions to partial functions; any partial function $f : \mathbb{N} \rightarrow \Sigma$ will henceforth be called a string, and $\Psi(f)$ will be called its normal representation.*

As previously observed, from the normal representation of string f it is only the normal form of f which necessarily may be recovered. For instance, the normal representation of string $f = \{(3, b), (1, a), (117, c)\}$ is abc , and given abc it is only the normal form $\Phi(abc) = \{(1, b), (0, a), (2, c)\}$ of f which may be recovered.

An *edit operation* is either a *change*, *insert*, or *delete* operation, and these three types of operations are mutually exclusive. Let \mathcal{E} be the set of all edit operations. Elements of \mathcal{E} map strings (i.e., partial functions from \mathbb{N} to Σ) to strings. Change operators are denoted by $m \xrightarrow{c} b$ where $m \in \mathbb{N}$ and $b \in \Sigma$, and are defined by

$$m \xrightarrow{c} b(f) = (f \setminus \{(m, f(m))\}) \cup \{(m, b)\} \text{ if } m \in D_f$$

Insert operators are denoted by $m \xrightarrow{i} b$ where $m \in \mathbb{N}$ and $b \in \Sigma$, and are defined by

$$m \xrightarrow{i} b(f) = f \cup \{(m, b)\} \text{ if } m \notin D_f$$

Delete operators are denoted by $m \xrightarrow{d}$ where $m \in \mathbb{N}$, and are defined by

$$m \xrightarrow{d}(f) = f \setminus \{(m, f(m))\} \text{ if } m \in D_f$$

Note that edit operations are not defined for every string; for example, $m \xrightarrow{d}(f)$ is undefined whenever $m \notin D_f$. Thus elements of \mathcal{E} are partial functions from the set of all strings to the set of all strings. When defined, the result of applying an edit operation to a string f is different from f except in the case $m \xrightarrow{c} b(f)$ where $(m, b) \in f$. In this case the edit operation (applied to f) is said to be *trivial*.

Let \mathcal{S} be the set of all finite sequences of edit operations. Elements of \mathcal{S} have the form $e_1 e_2 \dots e_n$ where $e_i \in \mathcal{E}$ and $n \in \mathbb{N}$ (if $n = 0$, the sequence is empty). Elements of \mathcal{S} are called *edit sequences* and may be interpreted as mapping strings to strings:

$$e_1 e_2 \dots e_n(f) = e_1(\dots e_{n-1}(e_n(f))\dots)$$

If $\varepsilon \in \mathcal{S}$ is the empty sequence, then $\varepsilon(f) = f$ for all strings f . Note that because elements of \mathcal{E} are partial functions, so too are edit sequences; $e_1 e_2 \dots e_n(f)$ is not necessarily defined.

Lemma 1 *For every string f , there exists an edit sequence s such that*

$$s(f) = \emptyset$$

Proof: If $f = \emptyset$, then let $s = \varepsilon$. Otherwise, let $D_f = \{i_0, \dots, i_n\}$. Note that $i_0 \xrightarrow{d} \dots i_n \xrightarrow{d}(f) = \emptyset$.

□

Lemma 2 *For every string f , there exists an edit sequence s such that*

$$s(\emptyset) = f$$

Proof: If $f = \emptyset$, then let $s = \varepsilon$. Otherwise, let $D_f = \{i_0, \dots, i_n\}$. Note that $i_0 \xrightarrow{i} f(i_0) \dots i_n \xrightarrow{i} f(i_n) (\emptyset) = f$.

□

Theorem 3 *For all strings f and g , there exists an edit sequence s such that $s(f) = g$.*

Proof: Appealing to Lemmas 1 and 2, let s_1 and s_2 be edit sequences such that $s_2(f) = \emptyset$ and $s_1(\emptyset) = g$. Now let $s = s_1 s_2$ (i.e., s is the concatenation of s_1 and s_2). It follows that $s(f) = s_1(s_2(f)) = g$.

□

Each type of edit operation has an associated cost. Let $\gamma_i > 0$ be the cost of an insert operation, $\gamma_d > 0$ be the cost of an delete operation, and $\gamma_c > 0$ be the cost of a nontrivial change operation; trivial edit operations have zero cost. Strictly speaking, a change operator does not have a cost; it's argument is required in order to determine whether it is trivial, in which case it has zero cost. Thus cost is associated with the pair operator and argument, rather than associated with operator alone.

To enable speaking of the cost of an edit sequence (which may contain a potentially trivial change operator) we say that $s \in \mathcal{S}$ takes f to g provided $s(f) = g$. Now the cost $\gamma(s, f)$ of an edit sequence s taking f to g may be inductively defined as follows

$$\begin{aligned} \gamma(\varepsilon, f) &= 0 \\ \gamma(m \xrightarrow{d}, f) &= \gamma_d \\ \gamma(m \xrightarrow{i} b, f) &= \gamma_i \\ \gamma(m \xrightarrow{c} b, f) &= 0 \text{ if } f = m \xrightarrow{c} b(f), \text{ and } \gamma_c \text{ otherwise} \\ \gamma(e_0 \dots e_n, f) &= \gamma(e_0, e_1 \dots e_n(f)) + \gamma(e_1 \dots e_n, f) \end{aligned}$$

When f can be inferred from context, $\gamma(s, f)$ is abbreviated by $\gamma(s)$. Moreover, to assert that s is an edit sequence taking f to g is to establish a context in which $\gamma(s, f)$ may be abbreviated by $\gamma(s)$.

Theorem 4 *Let s be an edit sequence taking f to h , and let s' be an edit sequence taking h to g . The concatenation $s'' = s' s$ takes f to g , and*

$$\gamma(s'') = \gamma(s') + \gamma(s)$$

Proof: Let $s' = e_0 \dots e_k$, and let $s = e_{k+1} \dots e_n$. Then $s'' = e_0 \dots e_n$ and according to the recursive definition for the cost of an edit sequence,

$$\begin{aligned} \gamma(s'', f) &= \sum_{j=0}^n \gamma(e_j, e_{j+1} \dots e_n(f)) \\ &= \sum_{j=0}^k \gamma(e_j, e_{j+1} \dots e_n(f)) + \sum_{j=k+1}^n \gamma(e_j, e_{j+1} \dots e_n(f)) \\ &= \sum_{j=0}^k \gamma(e_j, e_{j+1} \dots e_k(s(f))) + \gamma(s, f) \\ &= \gamma(s', h) + \gamma(s, f) \end{aligned}$$

□

Theorem 5 *Let s be an edit sequence taking f to g . If $\gamma_i = \gamma_d$, then there exists an edit sequence s' taking g to f such that $\gamma(s) = \gamma(s')$. Moreover, s' may be chosen to have the same length as s .*

Proof: Induct on the length of s . Base case: If $s = \varepsilon$ then $f = s(f) = g$. Let $s' = \varepsilon$ and observe that $\gamma(s) = \gamma(s') = 0$.

Inductive step: let $s = e_0 \dots e_k$, and let $h = e_1 \dots e_k(f)$. Since $s(f) = g$, it follows that $e_0(h) = g$. By the inductive hypothesis, there exists $e'_0 \dots e'_{k-1}$ taking h to f such that $\gamma(e'_0 \dots e'_{k-1}) = \gamma(e_1 \dots e_k)$. The proof is completed by showing there exists e'_k taking g to h such that $\gamma(e'_k) = \gamma(e_0)$. This would suffice because then

$$\begin{aligned} s'(g) &= e'_0 \dots e'_{k-1}(e'_k(g)) \\ &= e'_0 \dots e'_{k-1}(h) \\ &= f \end{aligned}$$

and by theorem 4,

$$\begin{aligned} \gamma(s') &= \gamma(e'_0 \dots e'_{k-1}, h) + \gamma(e'_k, g) \\ &= \gamma(e_1 \dots e_k, f) + \gamma(e_0, h) \\ &= \gamma(s) \end{aligned}$$

There are three cases to consider, based on the type of edit operation e_0 is.

Case 1: e_0 is $m \xrightarrow{c} g(m)$ (recall that $e_0(h) = g$; if e_0 is a change operator then g and h agree everywhere except possibly at m). Let e'_k be $m \xrightarrow{c} h(m)$. Note that if e_0 is trivial, then $g = h$ and e'_k is therefore trivial. If e_0 is nontrivial, then $h(m) \neq g(m)$ and e'_k is therefore nontrivial. In either case, $\gamma(e'_k) = \gamma(e_0)$.

Case 2: e_0 is $m \xrightarrow{i} g(m)$ (if e_0 is an insert operator, then g is the disjoint union $h \cup \{(m, g(m))\}$). Let e'_k be $m \xrightarrow{i}$. Since $\gamma_i = \gamma_d$, $\gamma(e'_k) = \gamma(e_0)$.

Case 3: e_0 is $m \xrightarrow{d}$ (if e_0 is a delete operator, then h is the disjoint union $g \cup \{(m, h(m))\}$). Let e'_k be the edit operation $m \xrightarrow{d} h(m)$. Since $\gamma_i = \gamma_d$, $\gamma(e'_k) = \gamma(e_0)$.

□

Edit distance is a function δ which maps two strings to a nonnegative real, and is defined by

$$\delta(f, g) = \min\{\gamma(s) \mid s \text{ is an edit sequence taking } f \text{ to } g\}$$

Given strings f and g , the *distance from f to g* is defined as $\delta(f, g)$. This distance exists, by theorem 3. That this definition is reasonable will now be established.

Lemma 6 *The distance between f and g is nonnegative, and is zero if and only if $f = g$.*

Proof: Distance is nonnegative because edit sequences have nonnegative cost. If $f = g$ then $\delta(f, g) = 0$ since $\varepsilon(f) = g$ and $\gamma(\varepsilon) = 0$. Conversely, let $s = e_0 \dots e_k$ be an edit sequence taking f to g having zero cost. It follows that every edit operator in s must be trivial, since otherwise the sum

$$\gamma(s) = \sum_{j=0}^k \gamma(e_j, e_{j+1} \dots e_k(f))$$

would be positive.

□

Lemma 7 *The triangle inequality holds,*

$$\delta(f, g) \leq \delta(f, h) + \delta(h, g)$$

Proof: Let $\delta(f, h) = \gamma(s)$ and $\delta(h, g) = \gamma(s')$, where $s(f) = h$ and $s'(h) = g$. Note that the concatenation $s'' = s's$ takes f to g , hence $\delta(f, g) \leq \gamma(s'')$. This completes the proof, since by theorem 4, $\gamma(s'') = \gamma(s') + \gamma(s)$.

□

Lemma 8 *If $\gamma_i = \gamma_d$, then distance is symmetric,*

$$\delta(f, g) = \delta(g, f)$$

Proof: It will be shown that $\delta(f, g) \geq \delta(g, h)$. That would complete the proof, since two applications of the inequality yield

$$\delta(f, g) \geq \delta(g, f) \geq \delta(f, g)$$

Let $\delta(f, g) = \gamma(s)$, where s takes f to g . By theorem 5, there exists s' taking g to f such that $\gamma(s') = \gamma(s)$. Hence, $\delta(g, f) \leq \gamma(s') = \gamma(s)$.

□

Theorem 9 *Edit distance is a metric if and only if $\gamma_i = \gamma_d$.*

Proof: To show the edit distance is a metric, it must be established that

1. $\delta(f, g) \geq 0$
2. $\delta(f, g) = 0 \iff f = g$
3. $\delta(f, h) \leq \delta(f, g) + \delta(g, h)$.
4. $\delta(f, g) = \delta(g, f)$.

Appealing to the previous lemmas shows the above properties hold when $\gamma_i = \gamma_d$ (in fact, only symmetry requires $\gamma_i = \gamma_d$).

Conversely, suppose property 4 holds. Let $f = \emptyset$ and $g = \{(0, a)\}$. If s is an edit sequence taking f to g , then s must contain an insert edit operation (otherwise $|s(f)| \leq |f| < |g|$). Since $g = 0 \xrightarrow{i} a(f)$, the distance from f to g is γ_i . If s' is an edit sequence taking g to f , then s' must contain a delete edit operation (otherwise $|s'(g)| \geq |g| > |f|$). Since $f = 0 \xrightarrow{d} (g)$, the distance from g to f is γ_d . Thus $\gamma_i = \delta(f, g) = \delta(g, f) = \gamma_d$. \square

3 Traces

A trace is an ordered triple (p, f, g) where f, g are strings and p is an increasing (i.e., $i < j \implies p(i) < p(j)$) partial function $p : D_f \longrightarrow D_g$. Trace (p, f, g) is referred to as a *trace from f to g* .

The cost of trace $t = (p, f, g)$ is defined as

$$\gamma(t) = |D_f \setminus D_p| \gamma_d + |D_g \setminus R_p| \gamma_i + \sum_{(i,j) \in p} \gamma(i \xrightarrow{c} g(j), f)$$

Let \tilde{f} denote the normal representation $\Psi(f)$ of string f . Define an equivalence relation \equiv on the set of strings by $f \equiv g$ if and only if $\tilde{f} = \tilde{g}$. Let $[f]$ denote the equivalence class of f .

Lemma 10 *Strings f and g are equivalent if and only if there exists an increasing onto function $\varphi : D_f \rightarrow D_g$ such that $f = g \circ \varphi$.*

Proof: Let $f = \{(i_0, f_0), \dots, (i_k, f_k)\}$ and $g = \{(j_0, g_0), \dots, (j_l, g_l)\}$, where $i_0 < \dots < i_k$ and $j_0 < \dots < j_l$. Suppose φ exists. Since φ is increasing, it is a one-to-one and onto order-preserving map from D_f to D_g . Hence $k = l$ and $\varphi(i_h) = j_h$. Therefore,

$$f_h = f(i_h) = g \circ \varphi(i_h) = g(j_h) = g_h$$

Conversely, suppose $f \equiv g$. Then $k = l$ and $f_h = g_h$ for $0 \leq h \leq k$. Let $\varphi(i_h) = j_h$, and note that $\varphi : D_f \rightarrow D_g$ is increasing and onto. Moreover,

$$g \circ \varphi(i_h) = g(j_h) = g_h = f_h = f(i_h)$$

□

Note that the function φ in lemma 10 is a one-to-one and onto order-preserving function since an increasing function preserves order and must be one-to-one.

Theorem 11 *Let t be a trace from f to g . There exists $u \in [f]$, $v \in [g]$, and an edit sequence s taking u to v , such that $\gamma(s) = \gamma(t)$.*

Proof: Let $t = (p, f, g)$. Let f be the string $\{(i_0, f_0), \dots, (i_k, f_k)\}$ where $i_0 < \dots < i_k$, and let g be the string $\{(j_0, g_0), \dots, (j_l, g_l)\}$ where $j_0 < \dots < j_l$. Let n be an integer greater than $i_k + j_l$, and let u be

$$u = \bigcup_{h=0}^k \{(i_h + n(1+h), f_h)\}$$

Note that u is equivalent to f . For the purposes of this proof, let the maximum of an empty collection of integers to be 0, and let v be the disjoint union

$$v = \bigcup_{h: i_h \in D_p} \{(i_h + n(1+h), g(p(i_h)))\} \cup \bigcup_{h: j_h \in D_g \setminus R_p} \{(j_h + \max\{i_m + n(1+m) \mid p(i_m) < j_h\}, g_h)\}$$

The union above which is indexed by $h : i_h \in D_p$ (call it the first union) involves mutually disjoint sets because $i_h + n(1+h)$ is an increasing function of h . Similarly, the union above which is indexed by $h : j_h \in D_g \setminus R_p$ (call it the second union) involves mutually disjoint sets. The first union is disjoint from the second union because otherwise there exist h and h' such that

$$i_h + n(1+h) = j_{h'} + \max\{i_m + n(1+m) \mid p(i_m) < j_{h'}\} \quad (1)$$

Reducing modulo n gives

$$i_h = \begin{cases} j_{h'} & \text{if } \emptyset = \{m \mid p(i_m) < j_{h'}\} \\ j_{h'} + i_{\max\{m \mid p(i_m) < j_{h'}\}} & \text{otherwise} \end{cases} \quad (2)$$

In the first case above ($\emptyset = \{m \mid p(i_m) < j_{h'}\}$), equation 1 simplifies to $i_h + n(1+h) = j_{h'}$ which contradicts $i_h = j_{h'}$. Assume therefore that the

second case of equation 2 holds. Substituting for i_h (as given by equation 2) into equation 1 and canceling $j_{h'}$ yields

$$\begin{aligned} & i_{\max\{m \mid p(i_m) < j_{h'}\}} + n(1+h) = \max\{i_m + n(1+m) \mid p(i_m) < j_{h'}\} \\ \Rightarrow & n(1+h) = n \max\{1+m \mid p(i_m) < j_{h'}\} \\ \Rightarrow & h = \max\{m \mid p(i_m) < j_{h'}\} \end{aligned}$$

Substituting for h (as given above) into equation 2 gives

$$\begin{aligned} & i_{\max\{m \mid p(i_m) < j_{h'}\}} = j_{h'} + i_{\max\{m \mid p(i_m) < j_{h'}\}} \\ \Rightarrow & j_{h'} = 0 \\ \Rightarrow & \emptyset = \{m \mid p(i_m) < j_{h'}\} \end{aligned}$$

contradicting the assumption that the second case of equation 2 holds. Note that v has been shown to be a function, since no two elements (of v) have the same first component.

By construction, $v = g \circ \varphi$ where $\varphi : D_v \rightarrow D_g$ is defined as

$$\begin{aligned} \varphi(i_h + n(h+1)) &= p(i_h) \quad \text{for } i_h \in D_p \\ \varphi(j_h + \max\{i_m + n(1+m) \mid p(i_m) < j_h\}) &= j_h \quad \text{for } j_h \in D_g \setminus R_p \end{aligned}$$

Since p is increasing, it is a one-to-one and onto order-preserving function from D_p to R_p . Moreover,

$$D_g = R_p \cup (D_g \setminus R_p) = p(D_p) \cup (D_g \setminus R_p)$$

Hence φ is onto. By lemma 10, $v \equiv g$ provided φ is an increasing function. As has already been observed, $i_h + n(1+h)$ and $p(i_h)$ are increasing functions of h . Thus φ is increasing when restricted to the set

$$A = \{i_h + n(1+h) \mid i_h \in D_p\}$$

Similarly, φ is increasing when restricted to

$$B = \{j_h + \max\{i_m + n(1+m) \mid p(i_m) < j_h\} \mid j_h \in D_g \setminus R_p\}$$

Let $a \in A$ and $b \in B$. To establish that φ is increasing, it remains to show

$$\begin{aligned} a < b &\implies \varphi(a) < \varphi(b) \\ b < a &\implies \varphi(b) < \varphi(a) \end{aligned}$$

Case 1: $a = i_h + n(1+h) < j_{h'} + \max\{i_m + n(1+m) \mid p(i_m) < j_{h'}\} = b$. The desired conclusion is $p(i_h) < j_{h'}$. Note that $a < b$ is contradicted by $\emptyset = \{m \mid p(i_m) < j_{h'}\}$ (since b then simplifies to $j_{h'}$). Therefore let $i_M + n(1+M) = \max\{i_m + n(1+m) \mid p(i_m) < j_{h'}\}$. If $p(i_h) \geq j_{h'}$, then

$h > M = \max\{m \mid p(i_m) < j_{h'}\}$. In particular, $h \geq M + 1$. This yields the contradiction

$$a = i_h + n(1 + h) > i_M + n(1 + M + 1) > j_{h'} + i_M + n(1 + M) = b$$

Case 2: $b = j_{h'} + \max\{i_m + n(1 + m) \mid p(i_m) < j_{h'}\} < i_h + n(1 + h) = a$. The desired conclusion is $j_{h'} < p(i_h)$. If this were not so, then $j_{h'} > p(i_h)$ (equality is impossible; $j_{h'} \in D_g \setminus R_p$). Thus $h \leq M = \max\{m \mid p(i_m) < j_{h'}\}$. This yields the contradiction

$$a = i_h + n(1 + h) \leq i_M + n(1 + M) \leq j_{h'} + i_M + n(1 + M) = b$$

Next an edit sequence s taking u to v will be constructed. Let

$$\begin{aligned} E &= \{b \xrightarrow{d} \mid b \in D_u \setminus A\} \\ C &= \{a \xrightarrow{c} v(a) \mid a \in A\} \\ I &= \{b \xrightarrow{i} v(b) \mid b \in B\} \end{aligned}$$

Let s' be the sequence of elements in E , let s'' be the sequence of elements in C , and let s''' be the sequence of elements in I . Define s as $s'''s''s'$. Note that $s'(u)$ is defined and is simply the restriction of u to A . Therefore, $s''s'(u)$ is defined and is simply the restriction of v to A . Since $A \cup B = D_v$ and A and B are disjoint, $s'''s''s'(u)$ is defined and is v .

The proof is completed by showing that $\gamma(s) = \gamma(t)$. By theorem 4 and the definition of the cost of an edit sequence,

$$\begin{aligned} \gamma(s) &= \gamma(s''', s''s'(u)) + \gamma(s'', s'(u)) + \gamma(s', u) \\ &= |I| \gamma_i + \gamma(s'', s'(u)) + |E| \gamma_d \end{aligned}$$

Note that

$$\begin{aligned} |E| &= |D_u \setminus A| = |D_u| - |A| = |D_u| - |D_p| = |D_f| - |D_p| = |D_f \setminus D_p| \\ |I| &= |B| = |D_g \setminus R_p| \end{aligned}$$

It follows that $\gamma(s) = \gamma(t)$ provided

$$\gamma(s'', s'(u)) = \sum_{(i,j) \in p} \gamma(i \xrightarrow{c} g(j), f) \quad (3)$$

Both sides of equation 3 contain $|C| = |A| = |D_p|$ terms, but some could be zero because change operators may be trivial. The right hand side can be rewritten as

$$\sum_{i_h \in D_p} \gamma(i_h \xrightarrow{c} g(p(i_h)), f)$$

A term (corresponding to $i_h \in D_p$) is trivial exactly when $g(p(i_h)) = f_h$. Since $v(a) = g \circ \varphi(a)$, the left hand side of equation 3 can be rewritten as

$$\sum_{i_h \in D_p} \gamma(i_h + n(1+h) \xrightarrow{c} g(p(i_h)), s'(u))$$

A term (corresponding to $i_h \in D_p$) is trivial exactly when

$$g(p(i_h)) = s'(u)(i_h + n(1+h)) = u(i_h + n(1+h)) = f_h$$

□

Let $t = (p, f, g)$ and $t' = (p', g, h)$ be two traces. Their *composition* $t' \circ t$ is the trace $(p' \circ p, f, h)$ from f to h where the composition $p' \circ p$ of partial functions is defined as

$$p' \circ p = \{(i, k) \mid \exists j. (i, j) \in p \text{ and } (j, k) \in p'\}$$

Note that the composition $t' \circ t$ is defined only when the third component of t is equal to the second component of t' . In this case they are said to be *composable*.

Lemma 12 *Given composable traces t and t' , $\gamma(t' \circ t) \leq \gamma(t) + \gamma(t')$.*

Proof: Let $t = (p, f, g)$ and $t' = (p', g, h)$. Note that

$$\begin{aligned} |D_f \setminus D_p| + |D_g \setminus D_{p'}| &\geq |D_f| - |D_p| + |R_p \setminus D_{p'}| \\ &= |D_f| - (|R_p| - |R_p \setminus D_{p'}|) \\ &= |D_f| - |R_p \cap D_{p'}| \\ &= |D_f \setminus D_{p' \circ p}| \end{aligned}$$

$$\begin{aligned} |D_g \setminus R_p| + |D_h \setminus R_{p'}| &\geq |D_{p'} \setminus R_p| + |D_h| - |R_{p'}| \\ &= |D_h| - (|D_{p'}| - |D_{p'} \setminus R_p|) \\ &= |D_h| - |D_{p'} \cap R_p| \\ &= |D_h| - |D_{p' \circ p}| \\ &= |D_h| - |R_{p' \circ p}| \\ &= |D_h \setminus R_{p' \circ p}| \end{aligned}$$

Therefore

$$\begin{aligned} \gamma(t) + \gamma(t') &= (|D_f \setminus D_p| + |D_g \setminus D_{p'}|) \gamma_d + (|D_g \setminus R_p| + |D_h \setminus R_{p'}|) \gamma_i \\ &+ \sum_{(i,j) \in p} \gamma(i \xrightarrow{c} g(j), f) + \sum_{(j,k) \in p'} \gamma(j \xrightarrow{c} h(k), g) \\ &\geq |D_f \setminus D_{p' \circ p}| \gamma_d + |D_h \setminus R_{p' \circ p}| \gamma_i \\ &+ \sum_{(i,j) \in p} \gamma(i \xrightarrow{c} g(j), f) + \sum_{(j,k) \in p'} \gamma(j \xrightarrow{c} h(k), g) \end{aligned}$$

The proof is completed by showing

$$\sum_{(i,k) \in p' \circ p} \gamma(i \xrightarrow{c} h(k), f) \leq \sum_{(i,j) \in p} \gamma(i \xrightarrow{c} g(j), f) + \sum_{(j,k) \in p'} \gamma(j \xrightarrow{c} h(k), g)$$

In the expression above, an edit operation corresponding to $(i, j) \in p$ is trivial exactly when $f(i) = g(j)$, and an edit operation corresponding to $(j, k) \in p'$ is trivial exactly when $g(j) = h(k)$. When both are trivial, $f(i) = g(j) = h(k)$ and the edit operation corresponding to $(i, k) \in p' \circ p$ is trivial. Consequently, to each nonzero term (i.e., nontrivial edit operation) in the first summation, there is a corresponding nonzero term in the second or third summation.

□

Theorem 13 *Let s be an edit sequence taking f to g . There exists a trace t taking f to g such that $\gamma(t) \leq \gamma(s)$.*

Proof: Induct on the length of s . Base case: $s = \varepsilon$. Note that $D_f = D_g$ since $g = s(f) = f$. Let $id : D_f \rightarrow D_g$ be the identity function. Let $t = (id, f, g)$ and note that $D_f = D_{id}$ and $R_{id} = D_g$. Therefore

$$\gamma(t) = \sum_{(i,j) \in id} \gamma(i \xrightarrow{c} f(j), f)$$

Every term in this sum is zero, since $(i, j) \in id \implies i = j \implies f(i) = f(j)$.

Inductive step: Let $s = e_0 \dots e_k$, and let $h = e_1 \dots e_k(f)$. By the inductive hypothesis, there exists a trace t' taking f to h such that $\gamma(t') \leq \gamma(e_1 \dots e_k)$. Since $s(f) = g$, it follows that e_0 takes h to g . The proof is completed by exhibiting a trace t'' taking h to g such that $\gamma(t'') = \gamma(e_0)$. Then $t = t'' \circ t'$ takes f to g and by lemmas 12 and 4,

$$\gamma(t) \leq \gamma(t'') + \gamma(t') \leq \gamma(e_0) + \gamma(e_1 \dots e_k) = \gamma(s)$$

There are three cases to consider, depending on the type of e_0 .

Case 1: $e_0 = m \xrightarrow{c} g(m)$. Note that $D_h = D_g$ since $m \xrightarrow{c} g(m)(h) = g$. Let $id : D_h \rightarrow D_g$ be the identity function. Let $t'' = (id, h, g)$ and note (as in the base case) that

$$\gamma(t'') = \sum_{i \in D_h} \gamma(i \xrightarrow{c} g(i), h)$$

Every term in this sum is zero, except possibly for $\gamma(m \xrightarrow{c} g(m), h)$, since $i \neq m \implies g(i) = h(i)$.

Case 2: $e_0 = m \xrightarrow{i} g(m)$. Note that $D_g = \{m\} \cup D_h$ since $m \xrightarrow{i} g(m)(h) = g$. Let $id : D_h \rightarrow D_g$ be the identity function. Let $t'' = (id, h, g)$ and note that $D_g \setminus R_{id} = \{m\}$. Therefore

$$\gamma(t'') = \gamma_i + \sum_{i \in D_h} \gamma(i \xrightarrow{c} g(i), h)$$

Every term in the summation over D_h is zero, since $i \in D_h \implies g(i) = h(i)$.

Case 3: $e_0 = m \xrightarrow{d}$. Note that $D_h = \{m\} \cup D_g$ since $m \xrightarrow{d} (h) = g$. Let $id^{-1} : D_g \rightarrow D_h$ be the identity function (observe that id is a partial function from D_h to D_g). Let $t'' = (id, h, g)$ and note that $D_h \setminus D_{id} = \{m\}$. Therefore

$$\gamma(t'') = \gamma_d + \sum_{i \in D_g} \gamma(i \xrightarrow{c} g(i), h)$$

Every term in the summation over D_h is zero, since $i \in D_g \implies g(i) = h(i)$.
□

An edit sequence s is said to take $[f]$ to $[g]$ if there exists $u \in [f]$ and $v \in [g]$ such that $s(u) = v$. A trace t is said to take $[f]$ to $[g]$ if its second component is equivalent to f and its third component is equivalent to g .

To refer to the cost of an edit sequence s taking $[f]$ to $[g]$ is to refer to

$$\gamma(s, [f], [g]) = \min\{\gamma(s, u) \mid u \in [f], s(u) \in [g]\}$$

An edit sequence s is called a minimal cost edit sequence taking $[f]$ to $[g]$ if it takes $[f]$ to $[g]$ and among all such edit sequences its cost (as given by the expression above) is minimal.

Theorem 14 *If s is a minimal cost edit sequence taking $[f]$ to $[g]$, and t is a minimal cost trace taking $[f]$ to $[g]$, then $\gamma(s, [f], [g]) = \gamma(t)$.*

Proof: By assumption, there exist $u \in [f]$ and $v \in [g]$ such that $s(u) = v$. Moreover, $\gamma(s, u)$ is minimal in the sense that it cannot decrease by changing s or u subject to the constraints that $u \in [f]$ and $s(u) \in [g]$. By theorem 13, there exists a trace t' from u to v such that $\gamma(t') \leq \gamma(s, u)$. By theorem 11, there exists an edit sequence s' taking $[u] = [f]$ to $[v] = [g]$ with cost no greater than $\gamma(t')$. That cost must in fact be $\gamma(s, u)$, since otherwise the minimality of $\gamma(s, u)$ would be contradicted. Therefore $\gamma(t') = \gamma(s, u)$. The proof is completed by showing $\gamma(t') \leq \gamma(t)$ ($\gamma(t) = \gamma(t') = \gamma(s, u)$ would then follow by minimality of $\gamma(t)$).

Let $t = (t_1, t_2, t_3)$. By theorem 11, there exists an edit sequence s'' taking $[t_2] = [f]$ to $[t_3] = [g]$ with cost no greater than $\gamma(t)$. Its cost contradicts

the minimality of $\gamma(s, u)$ if $\gamma(t) < \gamma(t')$.

□

Extend the concept of edit distance to equivalence classes by

$$\delta([f], [g]) = \min\{\delta(u, v) \mid u \in [f], v \in [g]\}$$

Since $\delta(u, v)$ is the minimum with respect to s of $\gamma(s, u)$ subject to $s(u) = v$, it follows that $\delta([f], [g])$ is the cost of a minimal cost edit sequence taking $[f]$ to $[g]$. By theorem 14, that coincides with the cost of a minimal cost trace taking $[f]$ to $[g]$.

Lemma 15 *Let $t = (t_1, t_2, t_3)$ be a trace, and let t' be a trace taking $[t_2]$ to $[t_3]$. There exists a partial function $t_1^* : D_{t_2} \rightarrow D_{t_3}$ such that $t^* = (t_1^*, t_2, t_3)$ is a trace and $\gamma(t^*) = \gamma(t')$.*

Proof: Let $t' = (t'_1, t'_2, t'_3)$. Since $t'_2 \equiv t_2$ and $t'_3 \equiv t_3$, by lemma 10 there exist one-to-one and onto order-preserving functions $\varphi_2 : D_{t_2} \rightarrow D_{t'_2}$ and $\varphi_3 : D_{t'_3} \rightarrow D_{t_3}$ such that $t_2 = t'_2 \circ \varphi_2$ and $t'_3 = t_3 \circ \varphi_3$. Let $t_1^* = \varphi_3 \circ t'_1 \circ \varphi_2$, and note that t_1^* is an increasing partial function from D_{t_2} to D_{t_3} . Hence $t^* = (t_1^*, t_2, t_3)$ is a trace. Because φ_2 and φ_3 are isomorphisms,

$$\begin{aligned} |D_{t_2} \setminus D_{t_1^*}| &= |D_{t_2}| - |D_{\varphi_3 \circ t'_1 \circ \varphi_2}| = |D_{t'_2}| - |D_{t'_1}| = |D_{t'_2} \setminus D_{t'_1}| \\ |D_{t_3} \setminus R_{t_1^*}| &= |D_{t_3}| - |R_{\varphi_3 \circ t'_1 \circ \varphi_2}| = |D_{t'_3}| - |R_{t'_1}| = |D_{t'_3} \setminus R_{t'_1}| \end{aligned}$$

It follows that $\gamma(t^*) = \gamma(t')$ provided the following equality holds

$$\sum_{(u,v) \in \varphi_3 \circ t'_1 \circ \varphi_2} \gamma(u \xrightarrow{c} t_3(v), t_2) = \sum_{(i,j) \in t'_1} \gamma(i \xrightarrow{c} t'_3(j), t'_2)$$

Replacing t_3 with $t'_3 \circ \varphi_3^{-1}$, replacing t_2 with $t'_2 \circ \varphi_2$ and making the change of variables $i = \varphi_2(u)$, $v = \varphi_3(j)$ puts the left hand side of the equality into the form

$$\sum_{(i,j) \in t'_1} \gamma(\varphi_2^{-1}(i) \xrightarrow{c} t'_3 \circ \varphi_3^{-1}(\varphi_3(j)), t'_2 \circ \varphi_2)$$

A term corresponding to (i, j) (in the left hand side of the equality) is therefore zero exactly when

$$t'_3(j) = t'_2 \circ \varphi_2(\varphi_2^{-1}(i)) = t'_2(i)$$

A term corresponding to (i, j) in the right hand side of the equality is zero exactly when

$$t'_3(j) = t'_2(i)$$

□

Theorem 16 *Edit distance is a metric on equivalence classes of \equiv if and only if $\gamma_i = \gamma_d$. Moreover, the following are equal*

1. $\delta([f], [g])$
2. *The cost of a minimal cost edit sequence taking $[f]$ to $[g]$.*
3. *The cost of a minimal cost trace taking $[f]$ to $[g]$.*
4. $\min\{\gamma((p, f, g)) \mid p : D_f \rightarrow D_g \text{ is an increasing partial function}\}$

Proof: The equality of the first three quantities listed was already noted in the discussion preceding lemma 15. Note that (p, f, g) (where $p : D_f \rightarrow D_g$ is an increasing partial function) is a trace taking $[f]$ to $[g]$. Hence the fourth quantity listed is at least $\delta([f], [g])$. Let t' be a minimal cost trace taking $[f]$ to $[g]$. By lemma 15, there exists an increasing partial function $p : D_f \rightarrow D_g$ such that $\gamma((p, f, g)) = \gamma(t')$. Therefore all four quantities listed above are equal.

Assume edit distance is a metric on equivalence classes of \equiv . Then it is symmetric. The same argument as given in the proof of theorem 9 shows

$$\delta([\emptyset], [\{(0, a)\}]) = \delta([\{(0, a)\}], [\emptyset]) \implies \gamma_i = \gamma_d$$

Assume $\gamma_i = \gamma_d$. Since $\delta(f, g)$ is nonnegative, so too is $\delta([f], [g])$. Since $\delta(f, g)$ is symmetric, so too is $\delta([f], [g])$. If $\delta([f], [g]) = 0$, then there exist $u \in [f]$ and $v \in [g]$ such that $\delta(u, v) = 0$. Hence $u = v$ and $[f] = [g]$. Moreover, $\delta([f], [f]) \leq \delta(f, f) = 0$. It remains to establish the triangle inequality.

By what has already been established, let $t = (p, f, h)$ be a trace such that $\delta([f], [h]) = \gamma(t)$. Similarly, let $t' = (p', h, g)$ be a trace such that $\delta([h], [g]) = \gamma(t')$. Note that $t' \circ t$ is a trace taking $[f]$ to $[g]$. By lemma 12,

$$\delta([f], [g]) \leq \gamma(t' \circ t) \leq \gamma(t) + \gamma(t') = \delta([f], [h]) + \delta([h], [g])$$

□

Trace t is said to be *minimal* if $\gamma(t) = \delta([t_2], [t_3])$. Note that it makes sense to speak of a minimal trace from f to g ; that refers to a trace (p, f, g) for which $\gamma((p, f, g)) = \delta([f], [g])$. According to theorem 16, such a trace exists.

Traces $t = (t_1, t_2, t_3)$ and $t' = (t'_1, t'_2, t'_3)$ are said to be *compatible* provided that $D_{t_2} \cap D_{t'_2} = \emptyset = D_{t_3} \cap D_{t'_3}$, and for all $(i, j) \in t_1$ and all $(i', j') \in t'_1$

$$i < i' \implies j < j' \quad \text{and} \quad i' < i \implies j' < j$$

If t and t' are compatible, then $t \oplus t'$ is defined as

$$t \oplus t' = (t_1 \cup t'_1, t_2 \cup t'_2, t_3 \cup t'_3)$$

If t and t' are not compatible, then $t \oplus t'$ is undefined.

Lemma 17 *If $t \oplus t'$ is defined, then it is a trace and $\gamma(t \oplus t') = \gamma(t) + \gamma(t')$.*

Proof: Let $t = (t_1, t_2, t_3)$ and $t' = (t'_1, t'_2, t'_3)$ be compatible traces. Then $\emptyset = D_{t_2} \cap D_{t'_2} = D_{t_3} \cap D_{t'_3}$. Since t_2 and t'_2 have disjoint domains, $t_2 \cup t'_2$ is a function and

$$D_{t_2 \cup t'_2} = D_{t_2} \cup D_{t'_2}$$

Likewise, $t_3 \cup t'_3$ is a function and

$$D_{t_3 \cup t'_3} = D_{t_3} \cup D_{t'_3}$$

Likewise, $t_1 \cup t'_1$ is a function which is partitioned by t_1 and t'_1 . Hence $t_1 \cup t'_1$ is a partial function from $D_{t_2 \cup t'_2}$ to $D_{t_3 \cup t'_3}$. Moreover, $t_1 \cup t'_1$ is increasing since for all $(i, j) \in t_1$ and all $(i', j') \in t'_1$

$$i < i' \implies j < j' \quad \text{and} \quad i' < i \implies j' < j$$

If (i, j) and (i', j') are both in either t_1 or t'_1 , the above implications hold because t_1 and t'_1 are increasing. Therefore, $t \oplus t'$ is a trace. Note that, because of the disjoint unions involved,

$$\begin{aligned} |D_{t_2 \cup t'_2} \setminus D_{t_1 \cup t'_1}| &= |D_{t_2}| + |D_{t'_2}| - |D_{t_1}| - |D_{t'_1}| \\ &= |D_{t_2}| - |D_{t_1}| + |D_{t'_2}| - |D_{t'_1}| \\ &= |D_{t_2} \setminus D_{t_1}| + |D_{t'_2} \setminus D_{t'_1}| \end{aligned}$$

$$\begin{aligned} |D_{t_3 \cup t'_3} \setminus R_{t_1 \cup t'_1}| &= |D_{t_3}| + |D_{t'_3}| - |R_{t_1}| - |R_{t'_1}| \\ &= |D_{t_3}| - |R_{t_1}| + |D_{t'_3}| - |R_{t'_1}| \\ &= |D_{t_3} \setminus R_{t_1}| + |D_{t'_3} \setminus R_{t'_1}| \end{aligned}$$

$$\begin{aligned} \sum_{(i,j) \in t_1 \cup t'_1} \gamma(i \xrightarrow{c} (t_3 \cup t'_3)(j), t_2 \cup t'_2) &= \sum_{(i,j) \in t_1} \gamma(i \xrightarrow{c} t_3(j), t_2) + \\ &\quad \sum_{(i',j') \in t'_1} \gamma(i' \xrightarrow{c} t'_3(j'), t'_2) \end{aligned}$$

Therefore, $\gamma(t \oplus t') = \gamma(t) + \gamma(t')$.

□

Trace $t = (t_1, t_2, t_3)$ is said to *precede* trace $t' = (t'_1, t'_2, t'_3)$, denoted $t \prec t'$, provided

$$\begin{aligned} \max\{i \mid i \in D_{t_2}\} &< \min\{i \mid i \in D_{t'_2}\} \\ \max\{i \mid i \in D_{t_3}\} &< \min\{i \mid i \in D_{t'_3}\} \end{aligned}$$

where $\max \emptyset = -\infty$ and $\min \emptyset = +\infty$.

Theorem 18 *If $t = (t_1, t_2, t_3) \prec t' = (t'_1, t'_2, t'_3)$, then $t \oplus t'$ is a trace. If $t \oplus t'$ is minimal, then so are t and t' .*

Proof: By lemma 17, $t \oplus t'$ is a trace provided t and t' are compatible. Since t precedes t' , it follows that $D_{t_2} \cap D_{t'_2} = \emptyset = D_{t_3} \cap D_{t'_3}$. Moreover, if $(i, j) \in t_1$ and $(i', j') \in t'_1$, then $i < i'$ and $j < j'$. Hence t and t' are compatible. Note that the compatibility of traces t and t' is not influenced by either t_1 or t'_1 , because the compatibility follows from $t \prec t'$ which is defined independent of t_1 and t'_1 (whether t precedes t' depends only on $D_{t_2}, D_{t_3}, D_{t'_2}, D_{t'_3}$). Therefore (by lemma 17)

$$\gamma(t \oplus t') = \gamma(t) + \gamma(t')$$

and this equality remains valid when t_1 and t'_1 are treated as parameters and are allowed to change. Suppose $t \oplus t'$ is minimal. Then the left hand side of the equality is $\delta([t_2 \cup t'_2], [t_3 \cup t'_3])$. By theorem 16, it cannot decrease by changing $t_1 \cup t'_1$. However, if either t or t' were not minimal, then (by theorem 16) the right hand side of the equality could decrease by changing t_1 or t'_1 .

□

Lemma 19 *Let $t = (t_1, t_2, t_3)$ be a trace. Let $t_2 = \{(i_0, f_0), \dots, (i_k, f_k)\}$ where $i_0 < \dots < i_k$, and let $t_3 = \{(j_0, g_0), \dots, (j_l, g_l)\}$ where $j_0 < \dots < j_l$. If t_2 and t_3 are nonempty, then t can be expressed as $t = t'' \oplus t'$ where one of the following cases hold.*

1. $t' = (\{(i_k, j_l)\}, \{(i_k, f_k)\}, \{(j_l, g_l)\})$
2. $t' = (\emptyset, \{(i_k, f_k)\}, \emptyset)$
3. $t' = (\emptyset, \emptyset, \{(j_l, g_l)\})$

Moreover, if t is minimal, then so is t'' .

Proof: Let $t' = (t'_1, t'_2, t'_3)$. The three cases correspond to a case decomposition based on t_1 . The first case is $(i_k, j_l) \in t_1$, which can be described by saying both $i_k \in D_{t_1}$ and $j_l \in R_{t_1}$. The second case is $i_k \notin D_{t_1}$ and $j_l \in R_{t_1}$. The third case is $j_l \notin R_{t_1}$. In each case $t'' = (t''_1, t''_2, t''_3)$ must (by the definition of \oplus) be defined by

$$\begin{aligned} t''_1 &= t_1 \setminus t'_1 \\ t''_2 &= t_2 \setminus t'_2 \\ t''_3 &= t_3 \setminus t'_3 \end{aligned}$$

In every case, t' is clearly a trace. Moreover, $t'' \oplus t'$ is a trace (via theorem 18), assuming that t'' is a trace, since $t'' \prec t'$.

In case 1, t''_1 is a partial function from $D_{t_2} \setminus \{i_k\} = D_{t_2 \setminus t'_2}$ to $D_{t_3} \setminus \{j_l\} = D_{t_3 \setminus t'_3}$. Thus t'' is a trace.

In case 2, $t''_1 = t_1$ is a partial function from $D_{t_2} \setminus \{i_k\} = D_{t_2 \setminus t'_2}$ to D_{t_3} , because $i_k \notin D_{t_1}$. Thus t'' is a trace.

In case 3, $t''_1 = t_1$ is a partial function from D_{t_2} to $D_{t_3} \setminus \{j_l\} = D_{t_3 \setminus t'_3}$, because $j_l \notin R_{t_1}$. Thus t'' is a trace.

If t is minimal, then by theorem 18 so is t'' .

□

4 The Normal Distance Matrix

Define the *distance* between normal representations \bar{f} and \bar{g} as

$$d(\bar{f}, \bar{g}) = \delta([f], [g])$$

Note that distance is well-defined since $\bar{f} = \bar{h} \iff [f] = [h]$. By theorem 16, distance is a metric on normal representations if and only if $\gamma_i = \gamma_d$.

Let $s = e_0 \dots e_k$ be an edit sequence taking u to v . Normal representation \bar{u} is regarded as being transformed to \bar{v} through the following sequence \bar{s} of normal representations

$$\bar{s} = \Psi(u) \Psi(e_k(u)) \Psi(e_{k-1}e_k(u)) \dots \Psi(e_0 \dots e_k(u))$$

Each step in the sequence (from one element to the next) corresponds to one of three types of operations on normal representations. Let $\bar{w} = w_0 \dots w_n$. If e_i is a delete operation, then

$$w_0 \dots w_n \mapsto \Psi(e_i(w)) = w'_0 \dots w'_{n-1}$$

where there exists $0 \leq l \leq n$ such that

$$w'_j = \begin{cases} w_j & \text{if } j < l \\ w_{j+1} & \text{if } j > l \end{cases}$$

In other words, the l th element of \bar{w} has been removed. If $e_i = m \xrightarrow[c]{\quad} b$ is a change operation, then

$$w_0 \dots w_n \mapsto \Psi(e_i(w)) = w'_0 \dots w'_n$$

where there exists $0 \leq l \leq n$ such that

$$w'_j = \begin{cases} w_j & \text{if } j \neq l \\ b & \text{if } j = l \end{cases}$$

In other words, the l th element of \bar{w} has been changed; this is called a *trivial change* when $b = w_l$. If $e_i = m \xrightarrow{i} b$ is an insert operation, then

$$w_0 \dots w_n \mapsto \Psi(e_i(w)) = w'_0 \dots w'_{n+1}$$

where there exists $0 \leq l \leq n + 1$ such that

$$w'_j = \begin{cases} w_j & \text{if } j < l \\ b & \text{if } j = l \\ w_{j-1} & \text{if } j > l \end{cases}$$

In other words, b has been inserted into \bar{w} at position l .

To streamline exposition, refer to the three types of operations (on normal representations described above) as delete, change, and insert operations. Let them have respective costs γ_d , γ_c , and γ_i , except that the cost of a trivial change is zero. To distinguish these operators (which act on normal representations) from previously discussed operators, they are called *normal operators*. The sum of the costs of the normal operators corresponding to the steps (from one element to the next) in the sequence \bar{s} is therefore $\gamma(\bar{s})$.

Given any sequence r of normal representations

$$r = \bar{h}_0 \bar{h}_1 \dots \bar{h}_n$$

such that \bar{h}_{i+1} is the result of some normal operator o_i applied to \bar{h}_i , define its cost $\gamma(r)$ as the sum (over $0 \leq i < n$) of the costs of the operators o_i . Such a sequence is referred to as a *normal sequence*, and is described as being *from \bar{h}_0 to \bar{h}_n* . By the discussion above, if $s = e_0 \dots e_k$ is an edit sequence taking u to v , then the sequence

$$\bar{s} = \Psi(u) \Psi(e_k(u)) \dots \Psi(e_0 \dots e_k(u))$$

is a normal sequence taking \bar{u} to \bar{v} . Moreover, $\gamma(\bar{s}) = \gamma(s)$.

A minimum normal sequence from \bar{f} to \bar{g} is a minimal cost normal sequence from \bar{f} to \bar{g} . Define $m(\bar{f}, \bar{g})$ as the cost of a minimum normal sequence from \bar{f} to \bar{g} . Hence $m(\bar{f}, \bar{g}) \leq \gamma(\bar{s}) = \gamma(s)$ where s is a minimal cost edit sequence taking $[f]$ to $[g]$. It follows (via theorem 16) that

$$m(\bar{f}, \bar{g}) \leq d(\bar{f}, \bar{g})$$

There may be question as to whether $m(\bar{f}, \bar{g}) = d(\bar{f}, \bar{g})$, because it has not yet been established that *every* normal sequence r can be expressed as \bar{s} for some edit sequence s . The following lemma shows that to be the case, and therefore $m(\bar{f}, \bar{g})$ and $d(\bar{f}, \bar{g})$ coincide.

Lemma 20 *Given nonempty normal sequence r , there exists u and v and an edit sequence s taking u to v such that $r = \bar{s}$.*

Proof: To facilitate induction on the length of r , a stronger result will be proved; in addition, u may be chosen such that the distance between consecutive elements of D_u is arbitrarily large.

Base case: $r = \bar{h}$ where $h = \{(i_0, f_0), \dots, (i_l, f_l)\}$. Let $n \in \mathbb{Z}^+$ be arbitrary, and let $u = \{(i_0 n, f_0), \dots, (i_l n, f_l)\}$. Let $s = \varepsilon$ so that $\bar{s} = \bar{u} = \bar{h} = r$. Moreover, the distance between consecutive elements of D_u is at least n .

Inductive step: $r = \bar{h}_0 \dots \bar{h}_k$ where $\bar{h}_0 = a_0 \dots a_q$. Let o be the normal operator taking \bar{h}_0 to \bar{h}_1 , and let p be the location at which a change, insertion, or deletion takes place in \bar{h}_0 . Let $n \in \mathbb{Z}^+$ be arbitrary, and let s be an edit sequence taking h_1 to h_k such that $\bar{s} = \bar{h}_1 \dots \bar{h}_k$. Let $h_1 = \{(i_0, f_0), \dots, (i_l, f_l)\}$ where $i_0 < \dots < i_l$ and the distance between consecutive elements of D_{h_1} is greater than $2n$. The proof is completed by showing there exists an edit operation e taking u to h_1 where u may be chosen such that $\bar{u} = \bar{h}_0$ and the distance between consecutive elements of D_u is at least n . The required edit sequence is then se . There are three cases to consider, depending on the type of o .

Case 1: o is a change operator. Then $\bar{h}_0 = \bar{h}_1$ except perhaps at position p . Let $u = i_p \xrightarrow{c} a_p (h_1)$ and let $e = i_p \xrightarrow{c} f_p$. Then $\bar{u} = \bar{h}_0$ and $e(u) = h_1$ as required. Moreover, the distance between consecutive elements of $D_u = D_{h_1}$ is at least n .

Case 2: o is an insert operator. Then the element inserted by o is f_p and $\bar{h}_0 = f_0 \dots f_{p-1} f_{p+1} \dots f_q$. Let $u = i_p \xrightarrow{d} (h_1)$ and let $e = i_p \xrightarrow{i} f_p$. Then $\bar{u} = \bar{h}_0$ and $e(u) = h_1$ as required. Moreover, the distance between consecutive elements of $D_u = D_{h_1} \setminus \{i_p\}$ is at least n .

Case 3: o is a delete operator. Then $\bar{h}_1 = a_0 \dots a_{p-1} a_{p+1} \dots a_q$. Let i be $\lfloor (i_{p-1} + i_p)/2 \rfloor$. Let $u = i \xrightarrow{i} a_p (h_1)$ and let $e = i \xrightarrow{d}$. Then $\bar{u} = \bar{h}_0$ and $e(u) = h_1$ as required. Moreover, the distance between consecutive elements of $D_u = D_{h_1} \cup \{i_p\}$ is at least n .

□

Theorem 21 *Distance $d(\bar{f}, \bar{g})$ defined as the cost of a minimal trace taking f to g is a metric on normal representations if and only if $\gamma_i = \gamma_d$. Moreover, $d(\bar{f}, \bar{g})$ is equal to the cost of a minimal normal sequence from \bar{f} to \bar{g} .*

Proof: Theorem 16 established the claims regarding distance being a metric. Lemma 20 and the discussion preceding it complete the proof.

□

Given string $f = \{(i_0, f_0), \dots, (i_k, f_k)\}$ where $i_0 < \dots < i_k$, define $\sigma_j(f)$ for $0 < j \leq |f|$ to be the normal representation of $\{(i_0, f_0), \dots, (i_{j-1}, f_{j-1})\}$,

$$\sigma_j(f) = f_0 \dots f_{j-1}$$

Let $\sigma_0(f)$ be the empty sequence ε . Note that if $f \equiv h$ then $\sigma_j(f) = \sigma_j(h)$, so f may as well be normalized. Moreover, $\sigma_{|f|}(f)$ is the normal representation of f .

Given strings f and g , their *normal distance matrix* is the $1 + |f| \times 1 + |g|$ matrix D with i, j entry (for $0 \leq i \leq |f|$ and $0 \leq j \leq |g|$)

$$D_{i,j} = d(\sigma_i(f), \sigma_j(g))$$

In particular, $D_{|f|,|g|}$ is the distance between the normal representations of f and g .

The notation $[expression]$ is used in the following theorem to simplify exposition. It is defined as

$$[expression] = \begin{cases} 1 & \text{if } expression \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

Theorem 22 *Let f and g be nonempty normalized strings. For $0 < u \leq |f|$ and $0 < v \leq |g|$, their normal distance matrix D satisfies*

$$\begin{aligned} D_{0,0} &= 0 \\ D_{0,v} &= v \gamma_i \\ D_{u,0} &= u \gamma_d \\ D_{u,v} &= \min\{\gamma_c[f(u-1) \neq g(v-1)] + D_{u-1,v-1}, \gamma_d + D_{u-1,v}, \gamma_i + D_{u,v-1}\} \end{aligned}$$

Proof: Let $\bar{f} = f_0 \dots f_k$ and $\bar{g} = g_0 \dots g_l$. Using normal operators on normal representations (which is justified by theorem 21), it is clear that

$$\begin{aligned} D_{0,0} &= d(\varepsilon, \varepsilon) = 0 \\ D_{0,v} &= d(\varepsilon, g_0 \dots g_{v-1}) = v \gamma_i \\ D_{u,0} &= d(f_0 \dots f_{u-1}, \varepsilon) = u \gamma_d \end{aligned}$$

Note that $D_{u,v} = d(f_0 \dots f_{u-1}, g_0 \dots g_{v-1}) = \gamma(t)$ where $t = (t_1, t_2, t_3)$ is a minimal trace taking $\Phi(f_0 \dots f_{u-1})$ to $\Phi(g_0 \dots g_{v-1})$. Appealing to lemmas 17 and 19, $D_{u,v} = \gamma(t'') + \gamma(t')$ where $t = t'' \oplus t'$ and one of the following cases holds.

Case 1: If $t' = (\{(i_{u-1}, j_{v-1})\}, \{(i_{u-1}, f_{u-1})\}, \{(j_{v-1}, g_{v-1})\})$, then

$$\gamma(t') = \gamma(i_{u-1} \xrightarrow{c} g_{v-1}, \{(i_{u-1}, f_{u-1})\}) = \gamma_c[f_{u-1} \neq g_{v-1}]$$

Moreover, t'' is a minimal trace taking $\Phi(f_0 \dots f_{u-2})$ to $\Phi(g_0 \dots g_{v-2})$. Thus $\gamma(t'') = D_{u-1, v-1}$.

Case 2: If $t' = (\emptyset, \{(i_{u-1}, f_{u-1})\}, \emptyset)$, then $\gamma(t') = \gamma_d$ and t'' is a minimal trace taking $\Phi(f_0 \dots f_{u-2})$ to $\Phi(g_0 \dots g_{v-1})$. Thus $\gamma(t'') = D_{u-1, v}$.

Case 3: If $t' = (\emptyset, \emptyset, \{(j_{v-1}, g_{v-1})\})$, then $\gamma(t') = \gamma_i$ and t'' is a minimal trace taking $\Phi(f_0 \dots f_{u-1})$ to $\Phi(g_0 \dots g_{v-2})$. Thus $\gamma(t'') = D_{u, v-1}$.

It follows that $D_{u, v}$ is equal to some element of the set

$$\{\gamma_c[f(u-1) \neq g(v-1)] + D_{u-1, v-1}, \gamma_d + D_{u-1, v}, \gamma_i + D_{u, v-1}\}$$

If each element in this set is the cost of some sequence of normal operators taking $f_0 \dots f_{u-1}$ to $g_0 \dots g_{v-1}$, then the proof is complete by the minimality of $D_{u, v}$.

The first element in the set is the cost of changing f_{u-1} to g_{v-1} by a normal change operator followed by the cost of \tilde{s} where s is a minimal cost edit sequence from $[\Phi(f_0 \dots f_{u-2})]$ to $[\Phi(g_0 \dots g_{v-2})]$.

The second element in the set is the cost of deleting f_{u-1} by a normal delete operator followed by the cost of \tilde{s} where s is a minimal cost edit sequence from $[\Phi(f_0 \dots f_{u-2})]$ to $[\Phi(g_0 \dots g_{v-1})]$.

The third element in the set is the cost of \tilde{s} where s is a minimal cost edit sequence from $[\Phi(f_0 \dots f_{u-1})]$ to $[\Phi(g_0 \dots g_{v-2})]$ followed by the cost of inserting g_{v-1} by a normal insert operator.

□