

挿入・削除訂正符号のサイズの上下界式

安永 憲司

東京工業大学

Levenshtein 距離

$d_L(x, y) := \min \{ x \text{ を } y \text{ に変換するのに必要な挿入・削除数} \}$

- 例. $d_L(000, 111) = 6$, $d_L(101, 010) = 2$
- $|x| = |y| = n$ のとき, $0 \leq d_L(x, y) \leq 2n$

符号 C の最小 Levenshtein 距離: $d_L(C) := \min_{c_1 \neq c_2 \in C} d_L(c_1, c_2)$

$d_L(C) \geq d$ のとき, C は合計 $t \leq \left\lfloor \frac{d-1}{2} \right\rfloor$ 個の挿入・削除を訂正できる

- $C \subseteq \Sigma^n$ のとき, $d_L(C)$ は偶数で, $t \leq \frac{d_L(C)}{2} - 1$ 個を訂正できる
- 相対最小距離は $\frac{d_L(C)}{2n} \in [0, 1]$
- 符号の相対最小距離 $\geq \delta \rightarrow \delta$ 割合未満の挿入・削除を訂正

最良の符号サイズ $A_q(n, d)$

$$A_q(n, d) := \max\{ |C| : \exists C \subseteq \Sigma^n \text{ s.t. } |\Sigma| = q, d_L(C) \geq d \}$$

- $A_q(n, d)$ の上界式 \rightarrow 符号として存在しない領域
- $A_q(n, d)$ の下界式 \rightarrow 符号として存在しうる領域

漸近的な評価

- 符号長 $n \rightarrow \infty$ の場合の振る舞いを評価
- $A_q(n, d)$ を達成する符号 $C \subseteq \Sigma^n$ に対し,
符号化率 $R = \log_q |C|$ と相対最小距離 $\delta = \frac{d_L(C)}{2n}$ の
トレードオフを明らかにしたい

基本的な事実（その1）

挿入球 $I_t(\mathbf{x}) := \{ \mathbf{x} \text{ に } t \text{ 挿入してできる } \Sigma \text{ 上の文字列 } \mathbf{y} \}$

- $|I_t(\mathbf{x})| = \sum_{i=0}^t \binom{n+t}{i} (q-1)^i := I_q(n, t) \approx q^{(n+t)H_q\left(\frac{t}{n+t}\right)}$

削除球 $D_t(\mathbf{x}) := \{ \mathbf{x} \text{ から } t \text{ 削除してできる文字列 } \mathbf{y} \}$

- $|D_t(\mathbf{x})| \leq \binom{|\mathbf{x}|}{t}$
- \mathbf{x} のラン数 $= r(\mathbf{x}) \geq 2t$ のとき, $\sum_{i=0}^t \binom{r(\mathbf{x})-t}{i} \leq |D_t(\mathbf{x})| \leq \binom{r(\mathbf{x})+t-1}{t}$
 - 例. $r(0000) = 1, r(0011) = 2, r(0101) = 4$
- $R_q(n, r) := \{ \mathbf{x} \in \Sigma^n : r(\mathbf{x}) = r \}, |R_q(n, r)| = \binom{n-1}{r-1} q (q-1)^{r-1}$

基本的な事実 (その2)

挿入削除球 $L_{t,s}(\mathbf{x}) := \{ \mathbf{x} \text{ に } t \text{ 削除} \cdot s \text{ 挿入してできる文字列 } \mathbf{y} \}$

- $L_{t,0}(\mathbf{x}) = D_t(\mathbf{x}), L_{0,t}(\mathbf{x}) = I_t(\mathbf{x})$

$|L_{t,t}(\mathbf{x})|$ の緊密な評価式は未解決問題

- $|L_{t,t}(\mathbf{x})| \leq |D_t(\mathbf{x})| \cdot I_q(n-t, t)$ が上界では最も良い??

二重数え上げ (double counting) より, 以下が成り立つ

$$\sum_{\mathbf{y} \in \Sigma^{n+t}} |D_t(\mathbf{y})| = \sum_{\mathbf{x} \in \Sigma^n} |I_t(\mathbf{x})| = q^n \cdot I_q(n, t)$$

球充填 (sphere-packing) タイプの上界式

Theorem 1. 最小 Levenshtein 距離 $d = 2(t + 1)$ の $C \subseteq \Sigma^n$ に対し,

$$|C| \leq \left\lfloor \frac{q^{n+t}}{I_q(n, t)} \right\rfloor$$

証明：各 $c \in C$ に対し, $I_t(c) \subseteq \Sigma^{n+t}$ は互いに交わらないため

Corollary 1. 最小 Levenshtein 距離 δn , 符号化率 R の $C \subseteq \Sigma^n$ に対し,

$$R \leq (1 + \delta) \left(1 - H_q \left(\frac{\delta}{1 + \delta} \right) \right) + o(1)$$

主結果その 1 : Elias タイプの上界式

Theorem 2. 最小 Levenshtein 距離 $d < 2n$ の $C \subseteq \Sigma^n$ について,

$$t < \frac{nd}{2n - d}$$

を満たす任意の $t \geq 0$ に対し,

$$|C| \leq \left\lfloor \frac{(n+t)d}{(n+t)d - 2nt} \cdot \frac{q^{n+t}}{I_q(n, t)} \right\rfloor$$

Corollary 2. 最小 Levenshtein 距離 δn , 符号化率 R の $C \subseteq \Sigma^n$ に対し,

$$R \leq \frac{1}{1 - \delta} \left(1 - H_q(\delta) \right) + o(1)$$

Theorem 2 の証明

- 二重数え上げを、符号 C との共通部分に適用すると、

$$\sum_{\mathbf{y} \in \Sigma^{n+t}} |D_t(\mathbf{y}) \cap C| = \sum_{\mathbf{x} \in C} |I_t(\mathbf{x})| = |C| \cdot I_q(n, t)$$

- ランダムに $\mathbf{y} \in \Sigma^{n+t}$ を選ぶと、 $|D_t(\mathbf{y}) \cap C| \geq \frac{|C| \cdot I_q(n, t)}{q^{n+t}}$ を満たす $\mathbf{y} \in \Sigma^{n+t}$ が存在
- $|D_t(\mathbf{y}) \cap C|$ は半径 t のリスト復号のリストサイズ
- [Hayashi, Yasunaga (IEEE IT 2020)] のリストサイズ可能性を適用
 - $t < \frac{nd}{2n-d}$ に対し、 $|D_t(\mathbf{y}) \cap C| \leq \frac{(n+t)d}{(n+t)d-2nt}$

Hamming 距離の符号に対する上界式の適用

符号 C の最小 Hamming 距離 $\leq d$

Hamming 距離での $A_q(n, d)$ の上界式

→ C の最小 Levenshtein 距離 $\leq 2d$

→ Levenshtein 距離での $A_q(n, d)$ の上界式

Theorem 3. 最小 Levenshtein 距離 δn , 符号化率 R の $C \subseteq \Sigma^n$ に対し,

$$R \leq 1 - H_q(\theta - \sqrt{\theta(\theta - \delta)}) + o(1) \quad (\text{Elias 限界})$$

$$R \leq H_q\left(\frac{1}{q}(q-1 - (q-2)\delta - 2\sqrt{\delta(1-\delta)(q-1)})\right) + o(1) \quad (\text{MRRW 限界})$$

ここで, $0 \leq \delta \leq \theta = 1 - \frac{1}{q}$

平均球サイズによる下界式

Tolhuizen (IEEE IT 1997). X 上の距離関数 $\rho: X \times X \rightarrow \mathbb{Z}$ に対し,

- x 中心の半径 d の球サイズ $V_d(x) := |\{y \in X: \rho(x, y) \leq d\}|$
- 平均球サイズ $V_d^{\text{ave}} := \frac{1}{|X|} \sum_{x \in X} V_d(x)$

のとき, 最小距離 d の符号 C として, $|C| \geq \frac{|X|}{V_{d-1}^{\text{ave}}}$ が存在

Levenshtein (ISIT 2002). 任意の $1 \leq t \leq n$ に対し,
以下を満たす最小距離 $d = 2(t + 1) = 2\delta n$ の符号 C が存在.

$$|C| \geq \frac{q^{n+t}}{I_q(n-t, t)^2} \quad \text{つまり} \quad \text{符号化率 } R \geq 1 + \delta - 2H_q(\delta)$$

グラフ理論による下界式

集合 $X \subseteq \Sigma^n$ に対し, グラフ $G = (V, E)$ を

- 各 $x \in X$ が頂点, $d_L(x, y) \leq d - 1$ なら辺 (x, y) が存在と定めると,

- G の独立集合の最大サイズ $= \alpha(G) \Leftrightarrow A_q(n, d) = \alpha(G)$

- 独立数 (independence number) $\alpha(G)$ に関するグラフ理論の結果が利用できる

- Turán の定理より GV 限界が導かれる ([Tohlutzen \(1997\)](#))

- [Jiang, Vardy \(IEEE IT 2004\)](#) による GV 限界の改良もこれ

Caro-Wei 限界による下界式

Caro-Wei 限界

$$\alpha(G) \geq \sum_{x \in V} \frac{1}{1 + \deg(x)}$$

Theorem 5. 任意の $d = 2(t + 1) < n$, 整数 $1 \leq r \leq n$ に対し,

$$A_q(n, d) \geq \left\lfloor \frac{\left(q^n - \sum_{i=r}^n \binom{n-1}{i-1} q (q-1)^{i-1} \right)^2}{I_q(n-t, t) \left(q^{n-t} \cdot I_q(n, t) - \sum_{i=r}^n \binom{n-1}{i-1} q (q-1)^{i-1} \cdot \left(\sum_{j=1}^t \binom{i-t}{j} \right) \right)} \right\rfloor$$

証明：集合 $X \subseteq \Sigma^n$ をラン数 r 以上の文字列集合とし、Caro-Wei 限界に $\deg(x) \leq |L_{t,t}(x)| \leq |D_t(x)| \cdot I_q(n-t, t)$ を適用

Sala, Gabrys, Dolecek (ISIT 2014) の下界式 (式自体は省略)

グラフの三角形数 T を用いた以下を利用 (Δ は最大次数)

$$\alpha(G) \geq \frac{|V|}{10\Delta} \left(\log \Delta - \frac{1}{2} \log \left(\frac{T}{|V|} \right) \right)$$

Jiang, Vardy (2004) は $\frac{T}{|V|}$ の上界を与えて GV 限界を $\log n$ 倍改善

Sala, Gabrys, Dolecek (ISIT 2014) も漸近的に $\log n$ 倍改善

However, the present work is focused on non-asymptotic bounds. To the best of the authors' knowledge, Theorem 1 is so far the strongest lower bound on deletion-correcting codes, with an improvement on the order of $\log n$ over all existing bounds, in both the asymptotic and non-asymptotic cases.

との記載はあるが . . .

平均挿入・削除球サイズ上界の改良による下界

Levenshtein (ISIT 2002) は

$$|L_{t,t}(\mathbf{x})| \leq |D_t(\mathbf{x})| \cdot I_q(n-t, t)$$

から、以下を利用

$$V_d^{\text{ave}} := \frac{1}{|\Sigma^n|} \sum_{\mathbf{x} \in \Sigma^n} |L_{t,t}(\mathbf{x})| \leq \frac{I_q(n-t, t)}{q^n} \sum_{\mathbf{x} \in \Sigma^n} |D_t(\mathbf{x})| = \frac{I_q(n-t, t)^2}{q^t}$$

→ 平均挿入・削除球サイズ $|L_{t,t}(\mathbf{x})|$ の上界の改良を目指す

$|L_{t,t}(x)|$ の上界の改良

アイデア：数え上げの重複を考える

$D_t(00011110000111)$ において,

$x = 00\mathbf{0}111\mathbf{1}000\mathbf{0}111$ の黒い部分を削除したものを y とおく.

オレンジのペアの片方を一つずつ削除すると,
その結果 z は $2^3 = 8$ 通りあり, すべて y の部分文字列

→ 各 $|I_t(z)|$ において, $|I_{t-3}(y)|$ は重複して数え上げられている

→ $7|I_{t-3}(y)|$ は数えなくてよい

主結果その2：平均挿入・削除球サイズ上界の改良による下界

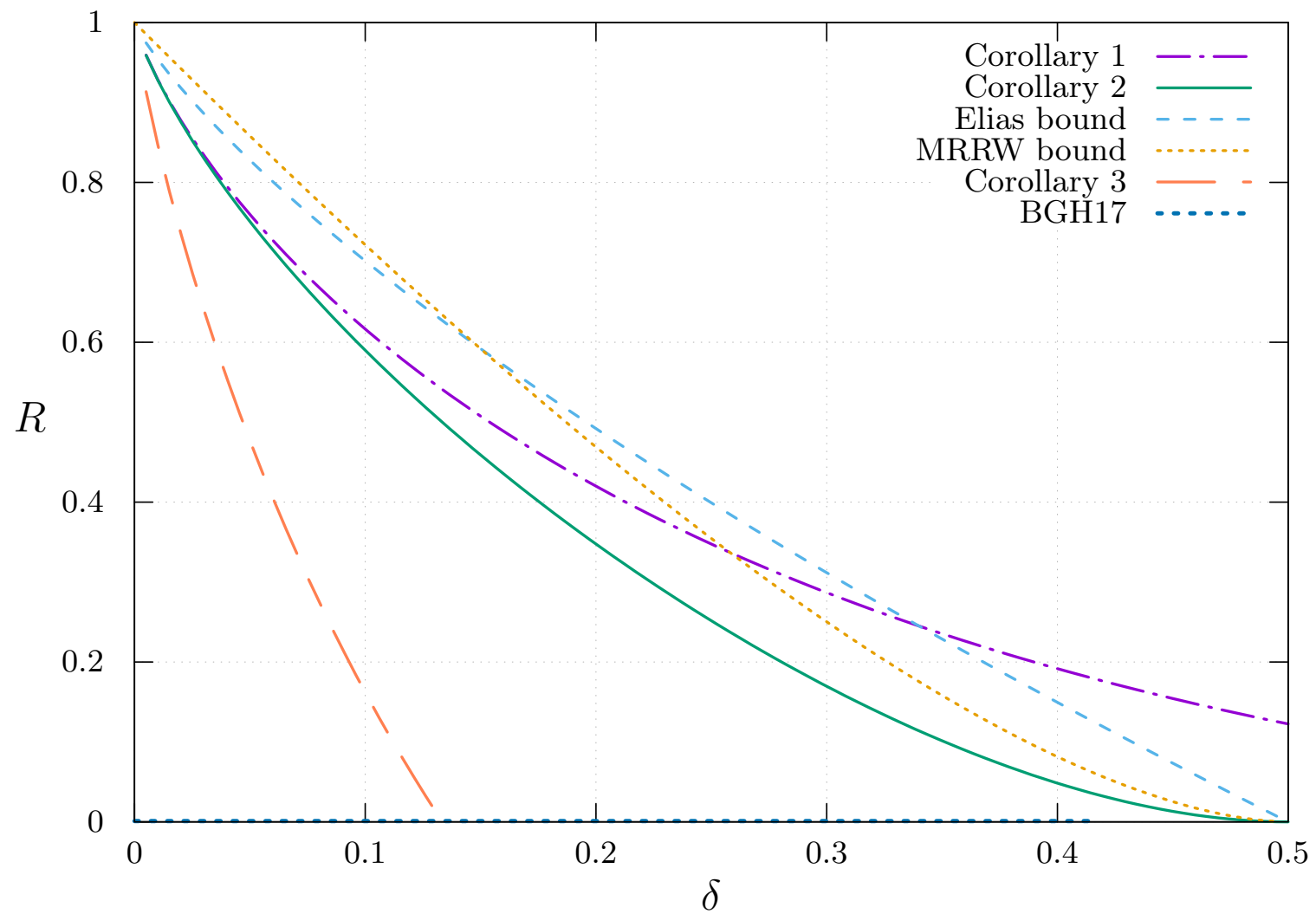
ラン数 $r(x) \geq 2$ のときオレンジペアが1つ以上あることを利用

Corollary 3. 任意の $1 \leq t \leq n$ に対し，以下を満たす
最小距離 $d = 2(t + 1)$ の符号 C が存在.

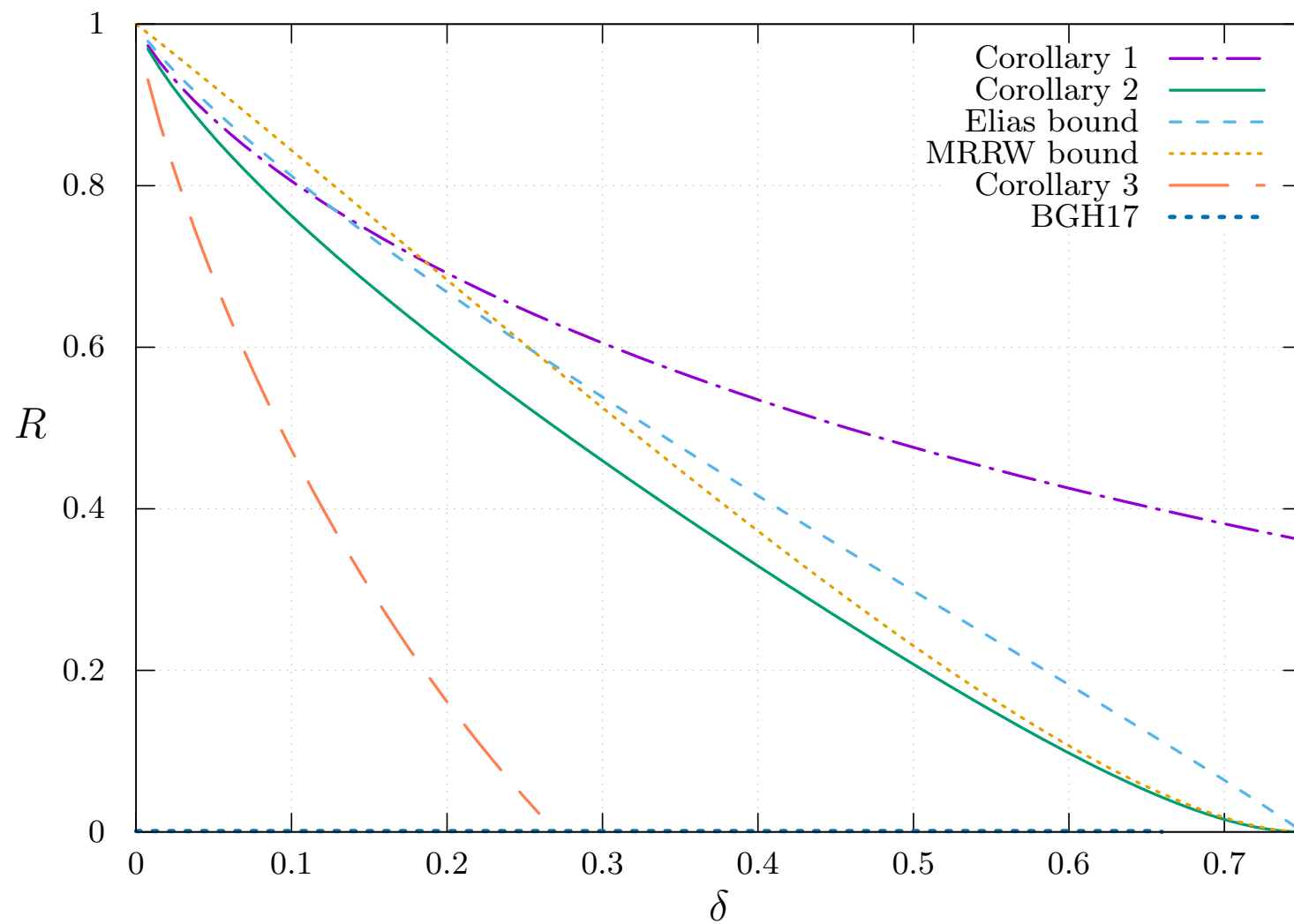
$$A_q(n, d) \geq \left\lfloor \frac{q^{n+t}}{I_q(n-t, t)^2 - q^{-(n-t)}(q^n - q)I_q(n-t+1, t-1)} \right\rfloor$$

符号化率は [Levenshtein\(2002\)](#) と同様に， $R \geq 1 + \delta - 2H_q(\delta)$

符号化率と相対最小距離のトレードオフ： $q = 2$



符号化率と相対最小距離のトレードオフ： $q = 4$



数值計算結果：上界

q	n	d	UB of [12]	Theorem 1	Theorem 2	q	n	d	UB of [12]	Theorem 1	Theorem 2
2	10	4	190	170	311	4	10	4	62 908	123 361	159 529
2	10	6	148	51	93	4	10	6	17 792	26 588	34 357
2	10	8	148	21	38	4	10	8	9 600	7 928	7 547
2	10	10	156	11	17	4	10	10	5 504	2 925	1 659
2	10	12	292	6	9	4	10	12	11 504	1 257	372
2	10	14	528	4	5	4	10	14	51560	608	87
2	10	16	772	3	4	4	10	16	173840	322	24
2	10	18	936	2	3	4	10	18	418736	184	10
2	20	4	97 453	95 325	181 643	4	20	4	30 003 945 118	68 719 476 736	90 174 299 388
2	20	6	33 903	16 513	31 402	4	20	6	2 902 217 544	8 197 663 580	10 754 599 022
2	20	8	26 456	4 096	7 772	4	20	8	752 550 391	1 402 773 785	1 839 884 106
2	20	10	26 456	1 295	2 452	4	20	10	360 221 648	306 647 351	316 287 316
2	20	14	26 456	213	287	4	20	14	146 887 008	24 329 793	11 105 216
2	20	18	91 688	56	54	4	20	18	108 563 408	3 094 985	409 222
2	20	22	340 556	20	16	4	20	22	2 517 203 000	549 256	16 755
2	20	26	709 300	9	7	4	20	26	31 608 638 744	125 240	851
2	20	30	961 048	5	4	4	20	30	192 278 071 952	34 771	71
2	20	34	1 038 520	3	3	4	20	34	587 772 208 784	11 321	15
2	20	38	1 048 198	2	2	4	20	38	977 086 753 268	4 206	7
2	40	4	47 498 012 376	52 357 696 560	102 167 009 660	4	40	10	$\approx 6 113 \times 10^{15}$	$\approx 27 238 \times 10^{15}$	$\approx 36 015 \times 10^{15}$
2	40	8	2 063 338 945	661 957 632	1 290 214 063	4	40	12	$\approx 1 502 \times 10^{15}$	$\approx 4 001 \times 10^{15}$	$\approx 5 290 \times 10^{15}$
2	40	12	1 130 893 408	25 385 916	49 417 671	4	40	16	$\approx 350 829 \times 10^{12}$	$\approx 135 888 \times 10^{12}$	$\approx 89 050 \times 10^{12}$
2	40	16	1 122 371 648	1 867 567	2 644 775	4	40	20	$\approx 133 526 \times 10^{12}$	$\approx 7 269 \times 10^{12}$	$\approx 2 172 \times 10^{12}$
2	40	20	1 122 371 648	215 900	203 859	4	40	30	$\approx 14 173 \times 10^{12}$	$\approx 18 554 \times 10^9$	$\approx 296 437 \times 10^6$
2	40	30	13 097 807 352	3 735	1 195	4	40	40	$\approx 34 641 \times 10^{15}$	$\approx 173 431 \times 10^6$	$\approx 69 \times 10^6$
2	40	40	287 193 094 240	231	43	4	40	50	$\approx 10 426 \times 10^{18}$	$\approx 3 882 \times 10^6$	33 642
2	40	50	914 362 931 844	33	8	4	40	60	$\approx 306 026 \times 10^{18}$	164 423 496	108
2	40	60	1 094 302 526 208	8	4	4	40	70	$\approx 1 074 \times 10^{21}$	11 354 434	10
2	40	70	1 099 503 766 738	3	3	4	40	78	$\approx 1 207 \times 10^{21}$	1 777 074	5
2	40	78	1 099 511 626 218	2	2						

数值計算結果：下界

q	n	d	LB of [12]	LB of [14]	Theorem 5	Corollary 3	q	n	d	LB of [12]	LB of [14]	Theorem 5	Corollary 3
2	10	4	16	2	18	17	4	10	4	4 364	842	4 489	4 382
2	10	6	1	0	1	1	4	10	6	88	14	78	88
2	10	8	0	0	0	0	4	10	8	4	0	3	4
2	10	10	0	0	0	0	4	10	10	0	0	0	0
2	10	12	0	0	0	0	4	10	12	0	0	0	0
2	10	14	—	—	—	—	4	10	14	—	—	—	—
2	10	16	—	—	—	—	4	10	16	—	—	—	—
2	10	18	—	—	—	—	4	10	18	—	—	—	—
2	10	20	—	—	—	—	4	10	20	—	—	—	—
2	20	4	4755	783	4968	4777	4	20	4	1 181 952 838	269 953 863	1 200 316 339	1 183 224 781
2	20	6	94	10	94	94	4	20	6	5 608 964	1 260 800	5 257 096	5 610 710
2	20	8	4	0	4	4	4	20	8	66 412	11 205	52 137	66 419
2	20	10	0	0	0	0	4	20	10	1 558	167	937	1 558
2	20	12	0	0	0	0	4	20	12	64	3	27	64
2	20	14	0	0	0	0	4	20	14	4	0	1	4
2	20	16	0	0	0	0	4	20	16	0	0	0	0
2	20	18	0	0	0	0	4	20	18	0	0	0	0
2	20	20	0	0	0	0	4	20	20	0	0	0	0
2	20	22	0	0	0	0	4	20	22	0	0	0	0
2	40	4	1 308 163 745	244 663 405	1 339 190 459	1 309 722 010	4	40	10	5 251 871 945 006	968 893 684 250	4 033 370 043 313	5 251 878 194 182
2	40	6	6 524 894	881 891	6 532 808	6 526 482	4	40	12	4 408 536 581	5 621 632 730	28 003 006 604	4 408 537 815
2	40	8	76 814	6 032	71 601	76 818	4	40	14	562 976 279	46 013 071	281 585 593	562 976 323
2	40	10	1 687	68	1 396	1 687	4	40	16	10 353 270	506 907	3 886 646	10 353 271
2	40	12	60	1	42	60	4	40	18	263 527	7 256	70 970	263 527
2	40	14	3	0	1	3	4	40	20	9 010	131	1 673	9 010
2	40	16	0	0	0	0	4	40	22	404	2	50	404
2	40	18	0	0	0	0							
2	40	20	0	0	0	0							
2	40	22	0	0	0	0							

参考文献

- [2] B. Bukh, V. Guruswami, and J. Håstad. An improved bound on the fraction of correctable deletions. *IEEE Trans. Inf. Theory*, 63(1):93–103, 2017.
- [4] D. Cullina and N. Kiyavash. An improvement to Levenshtein’s upper bound on the cardinality of deletion correcting codes. *IEEE Trans. Inf. Theory*, 60(7):3862–3870, 2014.
- [5] V. Guruswami, X. He, and R. Li. The zero-rate threshold for adversarial bit-deletions is less than $1/2$. In *62nd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2021, Denver, CO, USA, February 7-10, 2022*, pages 727–738. IEEE, 2021.
- [7] T. Hayashi and K. Yasunaga. On the list decodability of insertions and deletions. *IEEE Trans. Inf. Theory*, 66(9):5335–5343, 2020.
- [10] A. A. Kulkarni and N. Kiyavash. Nonasymptotic upper bounds for deletion correcting codes. *IEEE Trans. Inf. Theory*, 59(8):5115–5130, 2013.
- [11] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966.
- [12] V. I. Levenshtein. Bounds for deletion/insertion correcting codes. In *Proceedings IEEE International Symposium on Information Theory*, page 370, 2002.
- [14] F. Sala, R. Gabrys, and L. Dolecek. Gilbert-varshamov-like lower bounds for deletion-correcting codes. In *2014 IEEE Information Theory Workshop, ITW 2014, Hobart, Tasmania, Australia, November 2-5, 2014*, pages 147–151. IEEE, 2014.

数值計算結果：上界

q	n	d	UB of [12]	Theorem 1	Theorem 2	q	n	d	UB of [12]	Theorem 1	Theorem 2
2	10	4	190	170	311	4	10	4	62 908	123 361	159 529
2	10	6	148	51	93	4	10	6	17 792	26 588	34 357
2	10	8	148	21	38	4	10	8	9 600	7 928	7 547
2	10	10	156	11	17	4	10	10	5 504	2 925	1 659
2	10	12	292	6	9	4	10	12	11 504	1 257	372
2	10	14	—	—	5	4	10	14	—	—	87
2	10	16	—	—	4	4	10	16	—	—	24
2	10	18	—	—	3	4	10	18	—	—	10
2	10	20	—	—	—	4	10	20	—	—	—
2	20	4	97 453	95 325	181 643	4	20	4	30 003 945 118	68 719 476 736	90 174 299 388
2	20	6	33 903	16 513	31 402	4	20	6	2 902 217 544	8 197 663 580	10 754 599 022
2	20	8	26 456	4 096	7 772	4	20	8	752 550 391	1 402 773 785	1 839 884 106
2	20	10	26 456	1 295	2 452	4	20	10	360 221 648	306 647 351	316 287 316
2	20	12	26 456	490	768	4	20	12	226 003 920	80 420 755	60 876 276
2	20	14	26 456	213	287	4	20	14	146 887 008	24 329 793	11 105 216
2	20	16	41 520	104	118	4	20	16	79 778 144	8 267 148	2 003 849
2	20	18	91 688	56	54	4	20	18	108 563 408	3 094 985	409 222
2	20	20	190 416	32	28	4	20	20	536 774 720	1 258 226	79 926
2	20	22	340 556	20	16	4	20	22	2 517 203 000	549 256	16 755
2	40	4	47 498 012 376	52 357 696 560	102 167 009 660	4	40	10	6 113 592 833 576 549 294	27 238 444 999 469 732 769	36 015 920 136 083 097 898
2	40	6	6 561 107 408	4 865 095 698	9 488 059 400	4	40	12	1 502 985 120 942 552 080	4 001 768 904 009 233 099	5 290 974 980 372 832 578
2	40	8	2 063 338 945	661 957 632	1 290 214 063	4	40	14	694 833 352 099 147 362	690 127 171 352 978 162	628 749 366 837 598 021
2	40	10	1 279 636 864	117 292 187	228 473 245	4	40	16	350 829 440 062 284 900	135 888 933 076 115 960	89 050 511 830 036 160
2	40	12	1 130 893 408	25 385 916	49 417 671	4	40	18	205 583 936 785 834 080	29 940 446 039 885 620	13 234 636 413 304 032
2	40	14	1 122 371 648	6 445 783	12 539 381	4	40	20	133 526 342 747 906 144	7 269 429 537 145 809	2 172 089 508 608 137
2	40	16	1 122 371 648	1 867 567	2 644 775						
2	40	18	1 122 371 648	605 094	754 358						
2	40	20	1 122 371 648	215 900	203 859						
2	40	22	1 122 371 648	83 817	65 257						