

訂正可能な削除割合の限界

安永 憲司（金沢大学）

2017.2.23

ICAワークショップ@唐津市


削除訂正

- $010010 \rightarrow 01000, 001101 \rightarrow 00111$ (1個削除)
 $010010 \rightarrow 0110, 001101 \rightarrow 0110$ (2個削除)

- $(r + 1)$ 回繰り返し符号は、 r 個削除訂正可能


- $rep_{r+1}(s) =$

$s_1 \cdots s_1$	$s_2 \cdots s_2$...	$s_t \cdots s_t$
------------------	------------------	-----	------------------



- $s' =$

	s'_2	
--	--------	--



この位置は必ず s_2

- レート $\frac{t}{(r+1)t} = \frac{1}{r+1}$, 削除割合 $\frac{r}{(r+1)t} = O\left(\frac{1}{t}\right)$

削除訂正符号

- $[k] = \{1, 2, \dots, k\}$ 上の p 削除訂正符号 $C \subseteq [k]^n$
 $\Leftrightarrow \forall c_1, c_2 \in C, \text{LCS}(c_1, c_2) < (1 - p)n$
- $\text{LCS}(c_1, c_2)$: 最長共通部分系列の長さ (Longest Common Subsequence)
 - $\text{LCS}(111222, 121212) = 4$
 - $\text{LCS}(111222333, 123123123) = 5$
- c_1, c_2 を最小の削除数で一致させたときの長さが $\text{LCS}(c_1, c_2)$
 - このときの削除数 $= n - \text{LCS}(c_1, c_2)$
 - $n - \text{LCS}(c_1, c_2) > pn$ なら p 削除訂正可能

削除訂正割合の限界

- 正レートの符号で p 削除訂正可能な p の上限は？
 - $|C| \geq \exp(\Omega_k(n))$ で p 削除訂正 \rightarrow 冗長度 $O_k(1)$
- $p^*(k) = \limsup\{ p \in (0,1): \exists [k] \text{ 上の正レート } p \text{ 削除訂正符号} \}$
 - $p^*(k) \leq 1 - 1/k$
 - $p^*(2) \geq 0.17$
 - ランダム符号は $0.788n < \text{LCS} < 0.8263n$
 - $p^*(k) \geq 1 - O(1/\sqrt{k})$ [Guruswami, Wang 2015]
 - $k \rightarrow \infty$ のとき $E[\text{LCS}(c_1, c_2)] \sim \frac{2}{\sqrt{k}} n$ [Kiwi, Loebl, Matousek 2004]

紹介する論文

- Boris Bukh, Venkatesan Guruswami, Johan Hastad:
An improved bound on the fraction of correctable deletions. ECCCC TR15-117.
(SODA 2016 とその改良)

- 主結果：

- $\forall k \geq 2, p^*(k) \geq 1 - \frac{2}{k + \sqrt{k}}$
- $\forall \varepsilon > 0$, レート $r(\varepsilon, k) > 0$ の明示的な k 元符号が存在し、 $\text{LCS}(c_1, c_2) < \frac{2}{k + \sqrt{k}} + \varepsilon$
- 2元するとき、 $\sqrt{2} - 1 - \varepsilon > 0.414 - \varepsilon$ 削除訂正可能

準備

- $w \in [k]^n$ の部分系列 (subsequence) : w から削除して得られる系列
- $w \in [k]^n$ の部分語 (subword) : w 内の連続した系列
- w 内の部分系列 w' のスパン $\text{span}_w(w')$: w' を含む w の最短部分語の長さ
 - $\text{span}_{1121112}(212) = 5$, $\text{span}_{111222333}(122) = 3$
- w_1 と w_2 の共通部分系列 (common subsequence) (w'_1, w'_2) : $w'_1 = w'_2$
- w_1 と w_2 の最長共通部分系列の長さ $\text{LCS}(w_1, w_2)$
- $C \subseteq [k]^n$ に対し、 $\text{LCS}(C) = \max_{c_1 \neq c_2 \in C} \text{LCS}(c_1, c_2)$

準備

- w_1 と w_2 の共通部分系列 (w'_1, w'_2) のスパン：
 $\text{span}(w'_1, w'_2) = \text{span}_{w_1}(w'_1) + \text{span}_{w_2}(w'_2)$
- **事実.** w_1 と w_2 の任意の共通部分系列 (w'_1, w'_2) に対し、
 $\text{span}(w'_1, w'_2) \geq b \text{len}(w'_1) - c$ ならば、

$$\text{LCS}(w_1, w_2) \leq \frac{2n + c}{b}$$

- **証明：**

- $b \text{len}(w'_1) \leq \text{span}(w'_1, w'_2) + c \leq 2n + c$

- $\text{LCS}(w_1, w_2) = \max_{(w'_1, w'_2)} \text{len}(w'_1) \leq \frac{2n+c}{b}$

削除訂正符号のアルファベット削減

- $\forall \varepsilon > 0, K = K(\varepsilon) \gg k$ に対し、
 $C_1 \subseteq [K]^n$ with $\text{LCS}(C_1) \ll \varepsilon n$ を
 $C_2 \subseteq [k]^N$ with $\text{LCS}(C_2) \approx \frac{2}{k+\sqrt{k}} N$ へ変換
 - C_1 の各シンボルを、サイズ K の符号で接続符号化
 - 内符号の構成法により最終的な符号の性質が変わる
- きれいな構成法 (定理 4) : $\text{LCS}(C_2) \approx \frac{2}{k+1} N$
- 汚れた構成法 (定理 3) : $\text{LCS}(C_2) \approx \frac{2}{k+\sqrt{k}} N$

きれいな構成法

定理 4

$C_1 \subseteq [K]^n$ with $\text{LCS}(C_1) = \gamma n, k \geq 2$ に対し

ある $T = T(K, \gamma, k) \leq 32 \left(\frac{2k}{\gamma}\right)^K$ と $\tau: [K] \rightarrow [k]^T$ が存在し、

C_1 の各シンボルを τ で接続符号化して得られる

$C_2 \subseteq [k]^N, N = nT$ は、

2つの異なる $c, \tilde{c} \in C_2$ の共通部分系列 s に対し以下を満たす:

$$\text{span}(s) \geq (k + 1) \text{len}(s) - 4\gamma k N$$

特に、 $\text{LCS}(C_2) \leq \left(\frac{2+4\gamma k}{k+1}\right) N < \left(\frac{2}{k+1} + 4\gamma\right) N$

きれいな構成法の符号化法 (1/2)

- 整数 L を割り切れる整数 A に対し、振幅 A の語 :

$$f_A = (1^A 2^A \dots k^A)^{L/A}$$

- $\text{len}(f_A) = kL$

- $L = 4, k = 2$ のとき

$$f_1 = 12121212, f_2 = 11221122, f_4 = 11112222$$

- $L = 6, k = 3$ のとき

$$f_1 = 123123123123123123$$

$$f_2 = 112233112233112233$$

$$f_3 = 111222333111222333$$

$$f_6 = 111111222222333333$$

きれいな構成法の符号化法 (2/2)

- B/A が大きいとき f_A と f_B が長い共通部分系列を持たないことを利用
- $R \geq k$, $[K]$ 上の語 $w = \ell_1 \ell_2 \dots$ に対し、
$$\tilde{w} = f_{R^{\ell_1-1}} f_{R^{\ell_2-1}} \dots$$
 - $\text{len}(\tilde{w}) = kL \text{len}(w)$
 - $A, B \leq R^{K-1}$
 - \tilde{w} 内のシンボル x が、 w 内のシンボル y から展開されているとき、 y は x の親と呼ぶ

きれいな構成法の分析

補題 5

$f_A^\infty = (1^A 2^A \dots k^A)^*$ とし、 $kA \leq B$ のとき、
 f_A^∞ と f_B^∞ の共通部分系列 $s = (w'_1, w'_2)$ に対し、

$$\text{span}(s) \geq \left(k + 1 - \frac{kA}{B}\right) \text{len}(s) - 2(A + B)$$

■ 用語の定義

- **チャンク** : ℓ^A, ℓ^B の形をした部分語
- f_A^∞ 内のチャンクが w'_1 によってスパンされる
⇔ w'_1 のスパンに対応する部分語が、
そのチャンクのシンボルを1つ以上含む
- 例. $f_A^\infty = 111122223333111122223333 \dots$
 $w'_1 = \quad 11 \quad \quad \quad 31 \quad 2$
- スパンされるチャンクは、1111, 2222, 3333, 1111, 2222²

補題5の証明 (1/2)

- 方針： $s = (w'_1, w'_2)$ でスパンされるチャンク数を見積もる
- $s = k_1^{p_1} k_2^{p_2} \dots k_t^{p_t}$ の形をしている ($k_\ell \neq k_{\ell+1}$)
- $k_\ell^{p_\ell}$ は f_A^∞ 内で $k \left\lfloor \frac{p_\ell - A}{A} \right\rfloor + 1$ 個以上のチャンクをスパンする
 - k_ℓ^A というチャンクは含まれる
 - $p_\ell > A$ なら、その他の k_ℓ^A を $\left\lfloor \frac{p_\ell - A}{A} \right\rfloor$ 個含む

補題 5 の証明 (2/2)

- f_A^∞ と f_B^∞ の両方で $k_\ell^{p_\ell}$ によってスパンされるチャンクに含まれるシンボル数は

$$\phi(p_\ell) = A \left(k \left\lfloor \frac{p_\ell - A}{A} \right\rfloor + 1 \right) + B \left(k \left\lfloor \frac{p_\ell - B}{B} \right\rfloor + 1 \right)$$

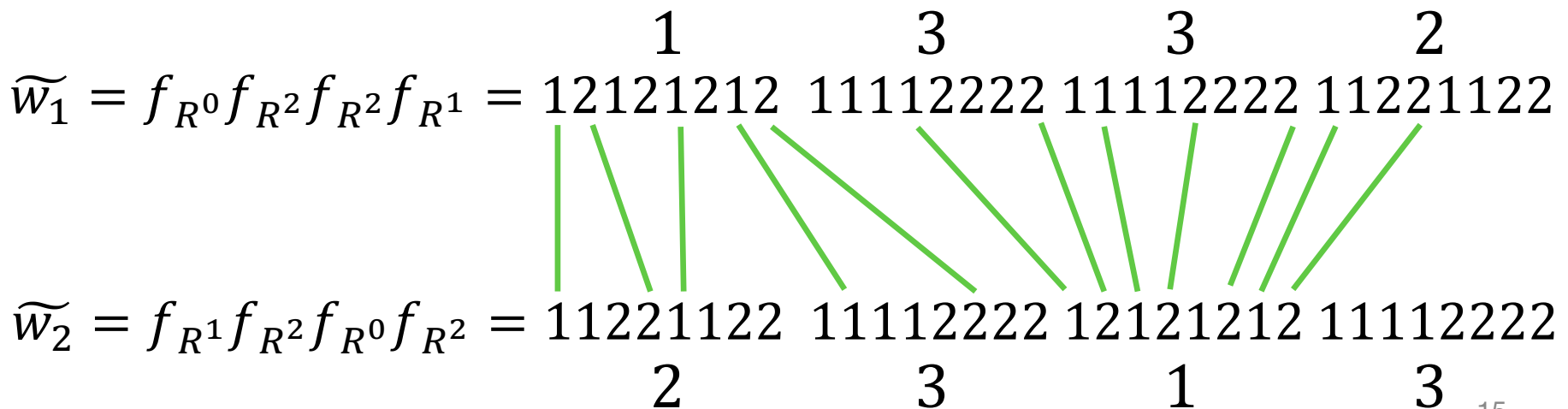
- $$\phi(p_\ell) \geq \begin{cases} k(p_\ell - A) + B & p_\ell \leq B \\ k(p_\ell - A) + k(p_\ell - B) + B & p_\ell > B \end{cases}$$
$$\geq \left(k + 1 - \frac{kA}{B} \right) p_\ell$$

- $k_\ell^{p_\ell}$ と $k_{\ell'}^{p_{\ell'}}$ でスパンされるチャンクは異なるため、
(s でスパンされるチャンクに含まれるシンボル数)
$$\geq \sum_\ell \phi(p_\ell) \geq \left(k + 1 - \frac{kA}{B} \right) \text{len}(s)$$

- 最初と最後のチャンクは全て入っていないかもしれないので $2(A + B)$ を引く (証明終)

共通部分系列のマッチ

- (w'_1, w'_2) : \widetilde{w}_1 と \widetilde{w}_2 の共通部分系列
- (w'_1, w'_2) 内の i 番目シンボルが **well-matched**
 $\Leftrightarrow w'_1[i]$ と $w'_2[i]$ の親が $[K]$ 内で同じ記号
- 共通部分系列が **badly-matched**
 \Leftrightarrow どのシンボルも well-matched でない
- 例. $L = 4, k = 2, R = 2, w_1 = 1332, w_2 = 2313,$



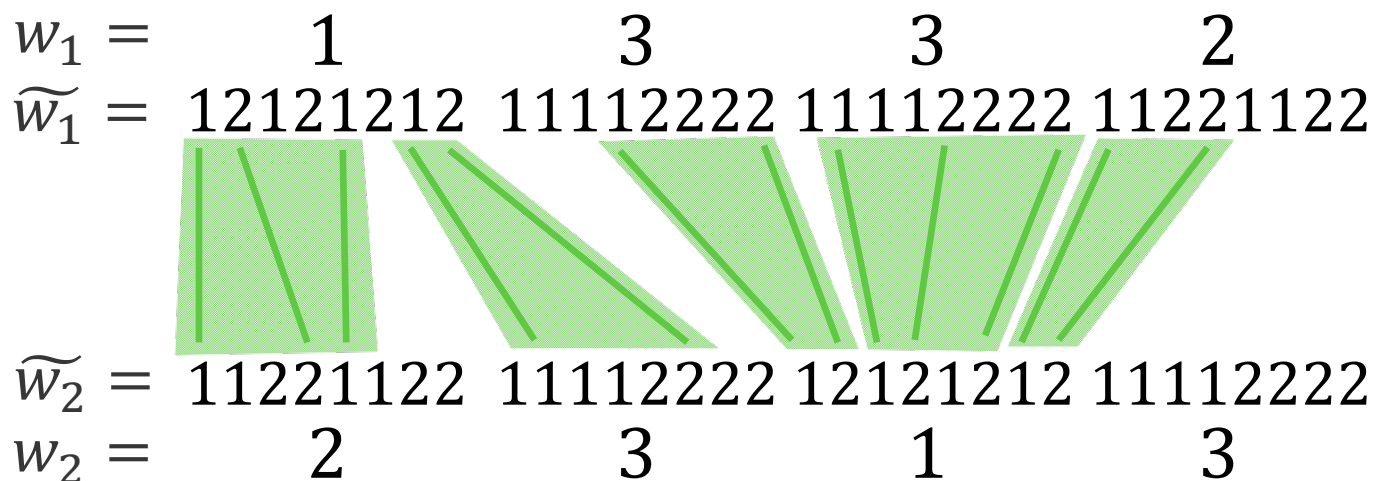
補題 6

$w_1, w_2 \in [K]^*$ であり、 \widetilde{w}_1 と \widetilde{w}_2 の共通部分系列 $s = (w'_1, w'_2)$ が badly-matched のとき

$$\text{span}(s) \geq \left(k + 1 - \frac{k}{R} - \frac{8R^{K-1}}{L} \right) \text{len}(s) - 16R^{K-1}$$

■ 証明：

- 以下を満たすように $s = (s_1, s_2, \dots, s_r)$ と分割
 - 各 s_i の w'_j における親が同じであり、 r は最小



補題6の証明 (続き)

- 各 s_i に対する w_1 と w_2 における親は異なる
- w_1 と w_2 における $r - 4$ 個以上のシンボルを展開したものは、 w'_1 と w'_2 のスパンに含まれる
 - w_1 と w_2 の左右両端4つを除けばよい
- $kL(r - 4) \leq \text{span}(s)$ であり、 $r \leq \frac{\text{span}(s)}{kL} + 4$

■ 補題5より

$$\begin{aligned} \text{span}(s) &\geq \left(k + 1 - \frac{k}{R}\right) \text{len}(s) - 2(A + B)r \\ &\geq \left(k + 1 - \frac{k}{R}\right) \text{len}(s) - 4rR^{K-1} \\ &\geq \left(k + 1 - \frac{k}{R}\right) \text{len}(s) - 4R^{K-1} \left(\frac{\text{span}(s)}{kL} + 4\right) \end{aligned}$$

(証明終)

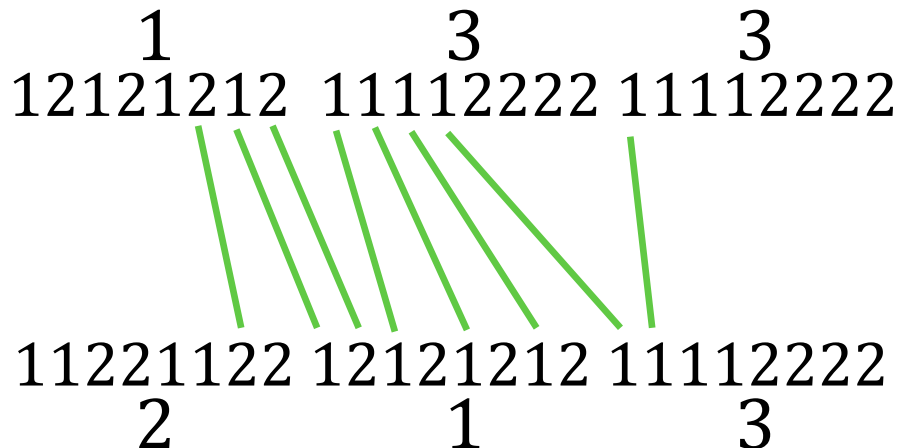
補題 7

$w_1, w_2 \in [K]^*$ であり、 \widetilde{w}_1 と \widetilde{w}_2 の共通部分系列 $s = (w'_1, w'_2)$ に対し

$$\text{span}(s) \geq \left(k + 1 - \frac{k}{R} - \frac{8R^{K-1}}{L} \right) \text{len}(s) - 2Lk(k+1)\text{LCS}(w_1, w_2) - 16R^{K-1}$$

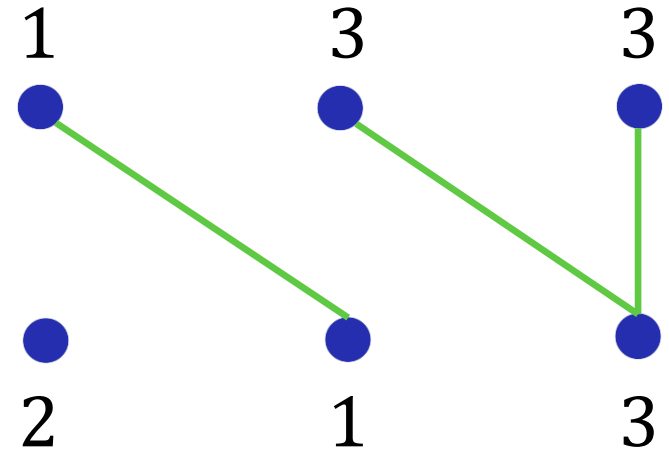
■ 証明 :

- s はスパン内で長さ最大と仮定
- $\text{len}(s)$ を増やせば $\text{span}(s)$ も増える



補題 7 の証明 (続き)

- マッチしたシンボル同士を連結した二部グラフ G を構成
 - 次数は 2 以下にできる
- グラフの最大マッチングは、 $|E(G)|/2$ 以上, $\text{LCS}(w_1, w_2)$ 以下
 $\rightarrow |E(G)| \leq 2\text{LCS}(w_1, w_2)$



- s 内のマッチしたシンボルをすべて削除 $\rightarrow s'$
 $\text{len}(s') \geq \text{len}(s) - Lk|E(G)| \geq \text{len}(s) - 2Lk \text{LCS}(w_1, w_2)$
- s' は badly-matched であり、補題 6 より

$$\text{span}(s) \geq \text{span}(s')$$

$$\geq \left(k + 1 - \frac{k}{R} - \frac{8R^{K-1}}{L} \right) \text{len}(s) - 2Lk(k+1)\text{LCS}(w_1, w_2) - 16R^{K-1}$$

(証明終)

きれいな構成法 (再掲)

定理 4

$C_1 \subseteq [K]^n$ with $\text{LCS}(C_1) = \gamma n, k \geq 2$ に対し

ある $T = T(K, \gamma, k) \leq 32 \left(\frac{2k}{\gamma}\right)^K$ と $\tau: [K] \rightarrow [k]^T$ が存在し、

C_1 の各シンボルを τ で接続符号化して得られる

$C_2 \subseteq [k]^N, N = nT$ は、

2つの異なる $c, \tilde{c} \in C_2$ の共通部分系列 s に対し以下を満たす:

$$\text{span}(s) \geq (k + 1) \text{len}(s) - 4\gamma k N$$

特に、 $\text{LCS}(C_2) \leq \left(\frac{2+4\gamma k}{k+1}\right) N < \left(\frac{2}{k+1} + 4\gamma\right) N$

定理 4 の証明

■ $C_1 \subseteq [K]^n$ with $\text{LCS}(C_1) = \gamma n$, $\varepsilon > 0$, $k \geq 2$ に対し
 $R = \left\lfloor \frac{2k}{\gamma} \right\rfloor$, $L = 16R^{K-1} \left\lfloor \frac{1}{\gamma} \right\rfloor$ とする

■ $\tau: [K] \rightarrow [k]^T$, $T = kL$ を $\tau(\ell) = f_{R^{\ell-1}}$ と定義し
 $C_2 \subseteq [k]^N$, $N = nkL$ とする

● $T = 16 \left\lfloor \frac{2k}{\gamma} \right\rfloor^{K-1} \left\lfloor \frac{1}{\gamma} \right\rfloor$ $k \leq 32 \left(\frac{2k}{\gamma} \right)^K$ である

■ 補題 7 より、 C_2 の任意の共通部分系列 s は

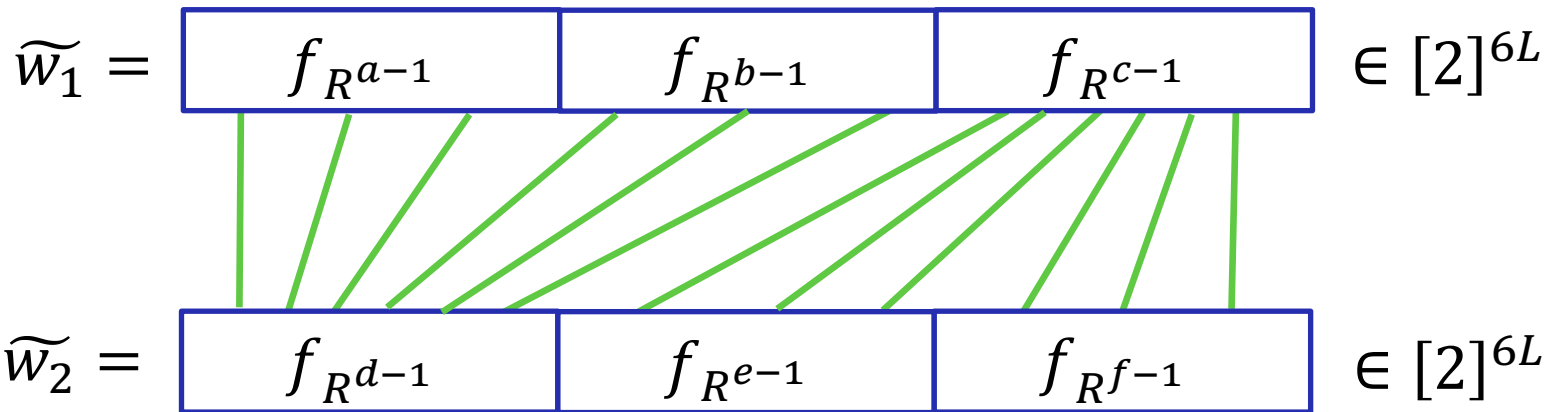
$$\begin{aligned} \text{span}(s) &\geq \left(k + 1 - \frac{k}{R} - \frac{8R^{K-1}}{L} \right) \text{len}(s) - 2Lk(k+1)\text{LCS}(w_1, w_2) \\ &\quad - 16R^{K-1} \\ &\geq \left(k + 1 - \frac{\gamma}{2} - \frac{\gamma}{2} \right) \text{len}(s) - 2(k+1)\gamma N - \gamma N \\ &\geq (k+1)\text{len}(s) - \gamma N - 2(k+1)\gamma N - \gamma N \\ &= (k+1)\text{len}(s) - 2(k+2)\gamma N \\ &\geq (k+1)\text{len}(s) - 4k\gamma N \end{aligned}$$

(証明終)

きれいな構成法の限界

- $k = 2$ で誤り割合 $p < 1/3$ しか達成できない理由

- $w_1 = abc, w_2 = def \in [K]^3$ が
 $d > a, b$ および $c > e, f$ を満たすとき



- $f_{R^{d-1}}$ は $f_{R^{a-1}}, f_{R^{b-1}}$ にマッチする ($f_{R^{c-1}}$ も同様)
 - $f_{R^{a-1}}, f_{R^{b-1}}$ の方が、変動が早いため
 - 長さ $4L$ の共通部分系列が存在
 - a, b, c, d, e, f に共通シンボルが存在していなくても
- $k > 2$ でも同様の議論より、 $p < \frac{k-1}{k+1}$ が限界

改善のアイデア (1/2)

- $f_A^\infty = (1^A 2^A \dots k^A)^*$, $f_B^\infty = (1^B 2^B \dots k^B)^*$, $A \ll B$
- スパンの短い共通部分系列がネック
 - 補題5ではチャンク内シンボル数の下界を導出

$$\phi(p_\ell) = A \left(k \left\lfloor \frac{p_\ell - A}{A} \right\rfloor + 1 \right) + B \left(k \left\lfloor \frac{p_\ell - B}{B} \right\rfloor + 1 \right)$$
 - $p_\ell = B$ のときに値が小さそう
 - $s = 1^B$ とすると、 f_B^∞ 内のスパンは f_A^∞ 内の半分

$$f_A^\infty = \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \dots$$

$$s = \boxed{1^A} \boxed{} \boxed{1^A} \boxed{} \boxed{1^A} \boxed{} \boxed{1^A}$$

$$f_B^\infty = \boxed{1^B} \boxed{2^B} \boxed{1^B} \boxed{2^B} \dots$$

改善のアイデア (2/2)

- アイディア： f_B^∞ 内の 1^B の中に適度に 2 を挿入
($2^B, 1^A, 2^A$ も同様)
- 長い共通部分系列を作ろうとするとき、
その「汚れ」を含めるか、含めないか？
 - 1^A とマッチさせている間は、 1^B 内の 2 とマッチさせても得しない
 - 2^A とマッチさせている間に、 1^B 内の 2 とマッチさせると、 1^B 内の 1 とマッチしなくなる

→ 共通部分系列は広がらず、汚れの分だけスパンが伸びる

$$f_A^\infty = \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \boxed{1^A} \boxed{2^A} \dots$$

$$s = \boxed{1^A} \boxed{} \boxed{1^A} \boxed{} \boxed{1^A} \boxed{} \boxed{1^A}$$

$$f_B^\infty = \boxed{} \boxed{2} \boxed{1^B} \boxed{2} \boxed{} \boxed{2^B} \boxed{} \boxed{1^B} \boxed{} \boxed{2^B} \dots$$

汚れた構成法 (2元の場合)

- $0 < c \leq \sqrt{2} - 1$ を固定
- 長さ M 振幅 a の汚れた 1:

$$\begin{aligned}
 1_{M,a} &= (1^a 2^{ca})^{M/(1+c)a} \\
 &= \underbrace{\boxed{1^a \ 2^{ca} \ 1^a \ 2^{ca} \ 1^a \ 2^{ca}} \dots \boxed{1^a \ 2^{ca}}}_{M}
 \end{aligned}$$

- (参考) きれいな構成: $f_{R^{i-1}} = (1^{R^{i-1}} 2^{R^{i-1}})^{L/R^{i-1}}$
- 汚れた構成: $g_i = (1_{R^{K+1+i}, R^{K-i}} 2_{R^{K+1+i}, R^{K-i}})^{L/R^{K+1+i}}$
 - 整数 $R = (1+c)t$ for integer t , $L = R^{2K+1}$
 - $\text{len}(g_i) = R^{K+1+i} \times 2 \times \frac{L}{R^{K+1+i}} = 2L$

汚れた構成法の分析

補題 8

$w_1 = 1_{\infty, a}$ (or $2_{\infty, a}$) および

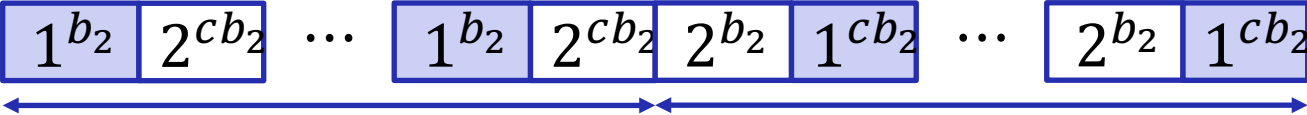
$w_2 = (1_{b_1, b_2} 2_{b_1, b_2})$ の部分系列 s に対し、

$$\text{span}_{w_1} s + 2b_1 \geq (3 + c)\text{len}(s) - \frac{4ab_1}{b_2}$$

■ 証明：

- $w_1 = 1_{\infty, a}$ と仮定

- 対称性より $w_1 = 2_{\infty, a}$ も同様に証明可能

- $w_2 = \boxed{1^{b_2}} \boxed{2^{cb_2}} \cdots \boxed{1^{b_2}} \boxed{2^{cb_2}} \boxed{2^{b_2}} \boxed{1^{cb_2}} \cdots \boxed{2^{b_2}} \boxed{1^{cb_2}}$


- s は $1^{b_2}, 1^{cb_2}, 2^{b_2}, 2^{cb_2}$ の部分語で構成

補題 8 の証明の続き (1/2)

■ s をそのような部分語に分解

- $s = 1^{s_1} 1^{s_2} 2^{t_1} 1^{s_3} 2^{t_2} \dots$

- i 番目の 1^* の部分語は s_i 個の 1 で構成

- i 番目の 2^* の部分語は t_i 個の 2 で構成

■ このとき $\text{span}_{w_1} 1^{s_i} \geq \left(\frac{s_i}{a} - 1\right) (1 + c)a$

- 同様に $\text{span}_{w_1} 2^{t_i} \geq \left(\frac{t_i}{ca} - 1\right) (1 + c)a$

■ s に合計 S 個の 1 , \tilde{S} 個の 2 が含まれるとき

$$\text{span}_{w_1} s \geq \sum_i \left(\frac{s_i}{a} - 1\right) (1 + c)a + \sum_i \left(\frac{t_i}{ca} - 1\right) (1 + c)a$$

$$\geq \frac{S}{a} (1 + c)a + \frac{\tilde{S}}{ca} (1 + c)a - \frac{\sum (1 + c)a}{\text{部分語の数} \leq \frac{b_1}{(1+c)b_2} \times 2 \times 2}$$

$$\geq (1 + c)S + \frac{1+c}{c} \tilde{S} - \frac{4ab_1}{b_2}$$

補題 8 の証明の続き (2/2)

- $\text{len}(s) = S + \tilde{S}$ であるため

$$(1 + c)S + \frac{1 + c}{c}\tilde{S} + 2b_1 \geq (3 + c)(S + \tilde{S})$$

を示せばよい

- $S, \tilde{S} \in [0, b_1]$ であり

$0 < c \leq \sqrt{2} - 1$ より $\frac{1+c}{c} > 3 + c$ であるから

$$\begin{aligned} & (1 + c)S + \frac{1+c}{c}\tilde{S} + 2b_1 \\ & \geq (1 + c)S + (3 + c)\tilde{S} + 2S \\ & \geq (3 + c)(S + \tilde{S}) \end{aligned}$$

(証明終)

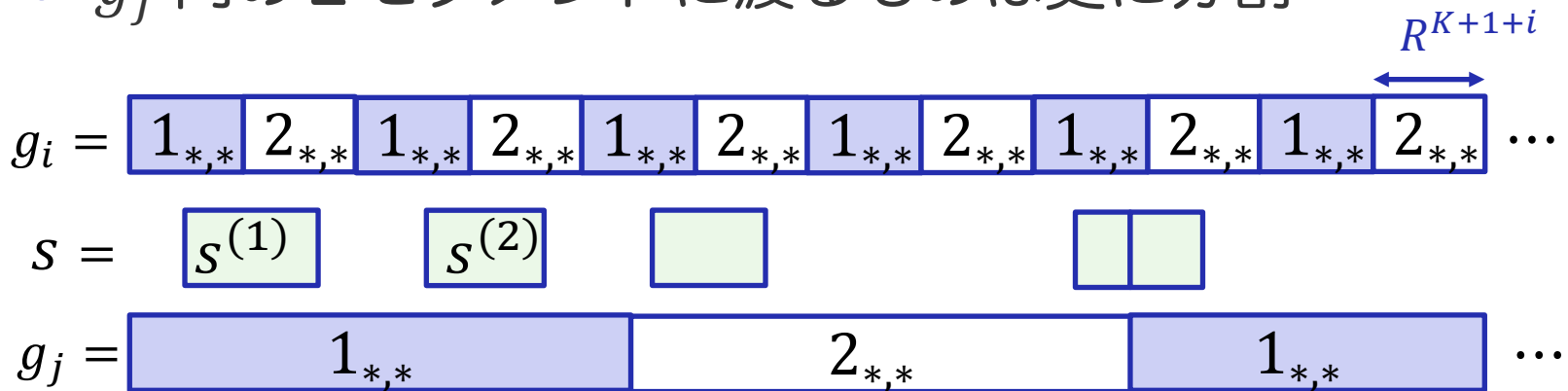
補題 9

g_i, g_j ($i < j$) の共通部分系列 s に対し、 $R \geq 10$ のとき

$$\left(1 + \frac{2}{R}\right) \text{span}_{g_i} s + \text{span}_{g_j} s \geq (3 + c) \text{len}(s) - \frac{10L}{R}$$

■ 証明：

- $g_i = (1_{R^{K+1+i}, R^{K-i}} 2_{R^{K+1+i}, R^{K-i}})^{L/R^{K+1+i}}$
- s を g_i 内のセグメント $(1_{R^{K+1+i}, R^{K-i}} 2_{R^{K+1+i}, R^{K-i}})$ との対応関係により $s^{(1)}, s^{(2)}, \dots$ と分割
- g_j 内の 2 セグメントに渡るものは更に分割



補題 9 の証明 (続き)

- $s = (s^{(1)}, s^{(2)}, \dots, s^{(p)}), p \leq 2 + \frac{\text{span}_{g_i} s}{2R^{K+1+i}} + \frac{2L}{R^{K+1+j}}$
- 各 $s^{(k)}$ は $a = R^{K-j}, b_1 = R^{K+1+i}, b_2 = R^{K-i}$ として補題 8 の仮定を満たす
 - $w_1 \Leftrightarrow g_j$ 内のセグメントの半分
 - $w_2 \Leftrightarrow g_i$ 内のセグメント
- 補題 8 より

$$\text{span}_{g_j} s^{(k)} + 2R^{K+i+1} \geq (3+c)\text{len}(s^{(k)}) - 4R^{i-j}R^{K+i+1}$$
- $\text{span}_{g_j} s \geq \sum_k \text{span}_{g_j} s^{(k)}$ と $\text{len}(s) = \sum_k \text{len}(s^{(k)})$ より

$$\text{span}_{g_j} s + 2pR^{K+i+1} \geq (3+c)\text{len}(s) - 4pR^{i-j}R^{K+i+1}$$
- p の不等式より

$$\begin{aligned} \text{span}_{g_j} s + (1 + 2R^{i-j})\text{span}_{g_i} s \\ \geq (3+c)\text{len}(s) - (1 + 2R^{i-j})(4R^{K+i+1} + 4LR^{i-j}) \end{aligned}$$
 - $R \geq 10, R^{K+i+1} \leq \frac{L}{R}, i < j$ より証明できる (証明終)

汚れた構成法の符号化法

- $[K]$ 上の語 $w = \ell_1 \ell_2 \dots$ に対し $\widehat{w} = g_{\ell_1} g_{\ell_2} \dots$

補題 10

$w_1, w_2 \in [K]^*$ に対し、 \widehat{w}_1 と \widehat{w}_2 の共通部分系列 $s = (w'_1, w'_2)$ が badly-matched のとき

$$\text{span}(w'_1) + \text{span}(w'_2) \geq \left(3 + c - \frac{28}{R}\right) \text{len}(s) - \frac{40L}{R}$$

補題 11

$w_1, w_2 \in [K]^*$, \widehat{w}_1 と \widehat{w}_2 の共通部分系列 $s = (w'_1, w'_2)$ は

$$\text{span}(s) \geq \left(3 + c - \frac{28}{R}\right) \text{len}(s) - 16L \text{LCS}(w_1, w_2) - \frac{40L}{R}$$

汚れた構成法

定理 3

$C_1 \subseteq [K]^n$ with $\text{LCS}(C_1) = \gamma n$, $k \geq 2$ に対し
ある $T = T(K, \gamma, k) \leq O((2k/\gamma)^{2K+2})$ と $\tau: [K] \rightarrow [k]^T$ が存在し、
 C_1 の各シンボルを τ で接続符号化して得られる
 $C_2 \subseteq [k]^N$, $N = nT$ は、
2つの異なる $c, \tilde{c} \in C_2$ の共通部分系列 s に対し以下を満たす:

$$\text{span}(s) \geq (k + \sqrt{k}) \text{len}(s) - 5\gamma k N$$

特に、 $\text{LCS}(C_2) \leq \left(\frac{2+5\gamma k}{k+\sqrt{k}}\right) N < \left(\frac{2}{k+\sqrt{k}} + 5\gamma\right) N$

よい削除訂正符号の存在性

補題 13

$\gamma, r > 0$ と整数 $K \geq 2$ が

$$2r \log K + 2h(\gamma) - \gamma \log K < 0$$

を満たすとき、すべての n に対し、

$[K]^n$ 上のサイズ K^{rn} の符号が存在し、

任意の異なる符号語 w, w' に対し、 $\text{LCS}(w, w') < \gamma n$

■ 証明：

- $w_1, \dots, w_{K^{rn}}$ を $[K]^n$ から独立にランダムに、戻すことなく、選んだ系列とする
- $i < j$ に対し、分布 (w_i, w_j) は、互いに異なるという条件下で $[K]^n$ から2つ独立に選ぶ分布に等しい
- 和集合上界より

$$\Pr[\text{LCS}(w_i, w_j) > \gamma n] \leq \binom{n}{\gamma n} K^{(1-\gamma)n} K^{-n} \leq \binom{n}{\gamma n}^2 K^{-\gamma n}$$

■ 証明の続き：

- 和集合上界より、十分大きな n に対し

$$\begin{aligned} \Pr[\exists w_i, w_j \in C, \text{LCS}(w_i, w_j) > \gamma n] &\leq K^{2rn} \binom{n}{\gamma n}^2 K^{-\gamma n} \\ &= 2^{n(2r \log K + 2h(\gamma) - \gamma \log K) + o(n)} < 1 \end{aligned} \quad (\text{証明終})$$

定理 14

整数 k を固定. 任意の $\varepsilon > 0$ に対し $\tilde{r} = (\varepsilon/k)^{O(\varepsilon^{-3})}$ が存在し、無限に多くの N について、レート \tilde{r} 以上の符号 $C \subseteq [k]^N$ が存在し、 $\text{LCS}(C) < \left(\frac{2}{k+\sqrt{k}} + \varepsilon\right)n$

■ 証明：

- 補題 13 で $\gamma = \frac{\varepsilon}{4}, r = \frac{\gamma}{6} = \frac{\varepsilon}{24}$ とすれば $C_1 \subseteq [K]^n, K \leq O\left(\frac{1}{\varepsilon^3}\right), \text{LCS}(C_1) \leq \frac{\varepsilon n}{4}, |C_1| \geq K^{rn}$ となる
- 定理 3 を適用すると $C_2 \subseteq [k]^N, \text{LCS}(C_2) \leq \left(\frac{2}{k+\sqrt{k}} + \varepsilon\right)N$
レートは $\tilde{r} = rn/Tn \geq (\varepsilon/k)^{O(\varepsilon^{-3})}$

(証明終)

明示的な構成法

■ Guruswami, Wang (RANDOM 2015)

- $\forall \gamma > 0, \exists K \leq O\left(\frac{1}{\gamma^5}\right)$

レート $\Omega(\gamma^3)$, $\text{LCS}(C) \leq \gamma n$ の符号 $C \subseteq [K]^n$ を $n(\log n)^{\text{poly}(1/\gamma)}$ 時間で構成できる

- 符号化と復号も $n(\log n)^{\text{poly}(1/\gamma)}$ 時間で可能

■ 定理 3 と組み合わせれば、

$\text{LCS}(C) < \left(\frac{2}{k+\sqrt{k}} + \varepsilon\right) N$ を満たす符号 $C \subseteq [k]^N$ を効率的に構成できる (定理 16)

効率的な復号ができる符号

- ハミング距離の大きな符号 (Reed-Solomon) と定理 16 の符号を接続符号化することで実現
 - 外符号の相対ハミング距離 ≈ 1 であれば、 $\text{LCS}(C) < \left(\frac{2}{k+\sqrt{k}} + \varepsilon\right) N$ が保たれる (補題 17)
- 構成法
 - 外符号 : F_q 上の RS 符号
 - 内符号 C_{in} : 符号長 $m = O(\log q)$, $|C_{in}| = q^2$
- 復号法
 - 境界線がわからないため、連続した系列をしらみつぶしに削除訂正 \rightarrow 各位置に対し複数のシンボル候補
 - RS の list recovery によってリスト復号
 - 符号の性質より、一意復号となる

[Bukh, Guruswami, Hastad] のまとめ

- $p^*(k)$: 正レート k 元符号で訂正できる削除割合の上極限
 - $p^*(k) < 1 - 1/k$ が成立
- $\forall k \geq 2, \varepsilon > 0$,
 \exists 効率的に構成可能なレート $r(q, \varepsilon) > 0$ の k 元符号で
 $1 - \frac{2}{k + \sqrt{k}} - \varepsilon$ 割合削除を $n^3 \text{poly}(\log n)$ 時間復号
 - $p^*(k) \geq 1 - \frac{2}{k + \sqrt{k}}$ (よって $p^*(k) = 1 - \Theta\left(\frac{1}{k}\right), k \rightarrow \infty$)
 - 2元符号で $\sqrt{2} - 1 - \varepsilon > 0.414 - \varepsilon$ 割合を訂正可能

削除訂正線形符号の限界

定理 8 [Brakensiek, Guruswami, Zbarsky (2016)]

k 削除訂正できる符号長 n の線形符号 C に対し、

$$\dim(C) \leq \frac{n}{k+1} + (k+1)^2$$

- k 削除訂正できる線形符号で達成可能な漸近的な符号化レートは $1/(k+1)$
 - $(k+1)$ 回繰り返し符号で達成可能

定理 8 の証明

- $C^0 = C, i = 1, \dots, k$ に対し

$$C^i = \{(x_{i+1}, \dots, x_n, x_1, \dots, x_i) \mid (x_1, \dots, x_n) \in C\}$$

補題 9

すべての $0 \leq i < j \leq k$ に対し $\dim(C^i \cap C^j) \leq j - i$

- 補題 9 の証明

- $z \in C^i \cap C^j$ のとき $x, y \in C$ が存在し

$$z = (x_{i+1}, \dots, x_n, x_1, \dots, x_i) = (y_{j+1}, \dots, y_n, y_1, \dots, y_j)$$

- $j > i$ より $(x_{i+1}, \dots, x_{n-j+i}) = (y_{j+1}, \dots, y_n)$

- $k > j$ より x, y は k 削除で同じ結果 $\rightarrow x = y$ である

- $\rightarrow \forall \ell \in \{1, \dots, n\}, x_\ell = x_{\ell+j-i \pmod n}$ ($j - i$ シフトで等価)

- 各 x_ℓ は最大で $j - i$ 種類の値 $\rightarrow x$ は 2^{j-i} 通り以下

$$\dim(C^i \cap C^j) \leq j - i$$

(補題の証明終)

定理 8 の証明

■ 補題 9 より

$$\begin{aligned} n &\geq \dim \left(\bigcup_{i=0}^k C^i \right) && \text{ここが間違い (未解決)} \\ &\geq \sum_{i=0}^k \dim(C^i) - \sum_{i=0}^k \sum_{j=i+1}^k \dim(C^i \cap C^j) \\ &= (k+1) \dim(C) - \sum_{i=0}^k \sum_{j=i+1}^k \dim(C^i \cap C^j) \\ &\geq (k+1) \dim(C) - (k+1)^3 \end{aligned}$$

■ したがって $\dim(C) \leq (k+1)^2 + n/(k+1)$

(証明終)

削除訂正符号に関する既存研究 (1/4)

■ ランダム削除通信路

- シンボル毎に確率 p で削除
- 通信路容量は未解決
 - $p \rightarrow 0$ のとき, レートは $1 - h(p)$ を達成可能
 - $p \rightarrow 1$ のときでも正レートを達成可能
 - $(1 - p)/9$ を達成 [Mitzenmacher, Drinea 2006]

削除訂正符号に関する既存研究 (2/4)

- Schulman, Zuckerman (SODA '97, IEEE IT '99)
 - 小さな定数割合削除を効率的に訂正する定数レート符号
- Guruswami, Wang (RANDOM 2015)
 - $0 < \varepsilon < 1/2$ に対し、 $(1 - \varepsilon)$ 割合削除を訂正するレート $\Omega(\varepsilon^2)$, アルファベットサイズ $\text{poly}(1/\varepsilon)$ の符号
 - 構成・符号化・復号時間は $N^{\text{poly}(1/\varepsilon)}$
 - $\varepsilon > 0$ に対し、 ε 割合削除を訂正するレート $1 - \tilde{O}(\sqrt{\varepsilon})$ の2元符号
 - 構成・符号化・復号時間は $N^{\text{poly}(1/\varepsilon)}$
 - $0 < \varepsilon < 1/2$, $1/2 - \varepsilon$ 割合削除をサイズ $(1/\varepsilon)^{O(\log \log(1/\varepsilon))}$ でリスト復号するレート $\Omega(\varepsilon^3)$ の2元符号
 - 構成・符号化・復号時間は $N^{\text{poly}(1/\varepsilon)}$

削除訂正符号に関する既存研究 (3/4)

- $\text{del}(N, k)$: 符号長 N の k 削除訂正符号の最大サイズ
 - 定数 k に対し、 k だけに依存する定数 $a_k > 0, A_k < \infty$ が存在し

$$a_k \frac{2^N}{N^{2k}} \leq \text{del}(N, k) \leq A_k \frac{2^N}{N^k}$$

- $\text{del}(N, 1) = \Theta\left(\frac{2^N}{N}\right)$
 - 1 削除訂正の Varshamov-Tenengolts (VT) 符号
 $\{(x_1, \dots, x_n) \in \{0,1\}^N \mid \sum_{i=1}^N ix_i \equiv 0 \pmod{N+1}\}$
- Brakensiek, Guruswami, Zbarsky (SODA 2016)
 - 定数 k に対し、冗長度 $N - \log(\text{del}(N, k)) = O(k^2 \log k \log N)$ の k 削除訂正符号の効率的な構成法
 - 復号時間 $O_k(N(\log N)^4)$

削除訂正符号に関する既存研究 (4/4)

■ Oblivious deletions

- 符号語を見ずに削除パターンを決める通信路
- ランダム削除と最悪ケース削除の間

■ Guruswami, Li (arXiv 2016)

- $\forall p \in (0,1), \exists R > 0, Enc: \{0,1\}^{Rn} \rightarrow \{0,1\}^n, Dec: \{0,1\}^{(1-p)n} \rightarrow \{0,1\}^{Rn} \cup \{\perp\}$ s.t. 任意の pn 個削除のパターン τ と $m \in \{0,1\}^{Rn}$ に対し
$$\Pr[Dec(\tau(Enc(m))) \neq m] \leq o(1)$$
- 明示的構成法で効率的な復号ができる
レート $(1-p)/180$ の符号

任意の削除レートを訂正できる！

未解決問題

- $p^*(k)$ の特定
 - $k = 2$ のとき、 0.414 と 0.5 の間
 - $p^*(k) = 1 - \Theta\left(\frac{1}{k}\right), k \rightarrow \infty$
- p 割合削除を訂正するレート $1 - p - \varepsilon$, アルファベットサイズ $\text{poly}(1/\varepsilon)$ の符号の構成
- ε 割合削除を訂正するレート $1 - \varepsilon \text{poly}(\log(1/\varepsilon))$ の2元符号の構成
 - レート $1 - O(\varepsilon(\log(1/\varepsilon)))$ はランダム符号で存在