

## 文字列の頻度分布による共通パターン発見

池田, 大輔  
九州大学情報基盤センター

山田, 泰寛  
九州大学大学院システム情報科学府

廣川, 佐千男  
九州大学情報基盤センター

<https://hdl.handle.net/2324/2968>

---

出版情報：情報処理学会研究報告：自然言語処理. 2003 (98), pp.25-32, 2003-09

バージョン：

権利関係：ここに掲載した著作物の利用に関する注意 本著作物の著作権は（社）情報処理学会に帰属します。本著作物は著作権者である情報処理学会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします。

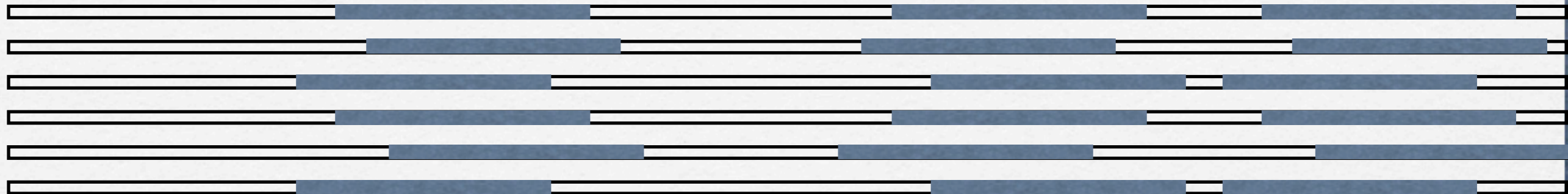
# 文字列の頻度分布による 共通パターン発見

池田 大輔、山田 泰寛\*、廣川佐千男

九州大学情報基盤センター

\*九州大学システム情報学研究院

# 共通パターン発見問題



最長共通部分列問題

NP完全



# 自然言語処理でのパターン発見

- 定型表現や未知語の抽出
  - Nagao & Mori (1994)など
  - 棋譜からの定石発見(中村2002)
- 固定長文字列( $n$ -gram)の頻度を利用
  - 頻度の高いところが抽出したいところ

**長さをどう決めるかが問題**

# 目的

- 共通パターン発見問題の定式化
- 問題を解くアルゴリズムの構築
  - 理論的に「解く」
  - 高速かつ頑健な実装

□ はじめに

□ パタン言語

□ テンプレート発見問題

□ 実験

□ まとめ

# はじめに



FI72&NL154

# 共通/頻出パタンの重要性

- データ/テキスト/Webマイニング
  - 隠れた有用な規則を見つける
- 遺伝子情報学
  - 重要な機能を司る部分は共通
- 圧縮、キャッシュ
- 自然言語処理
  - 定型表現、情報抽出



# 頻度による情報抽出

Ikeda *et. al.*, Discovery Science 2001 (LNCS 2226)

□ テンプレート部分は高頻度であろう

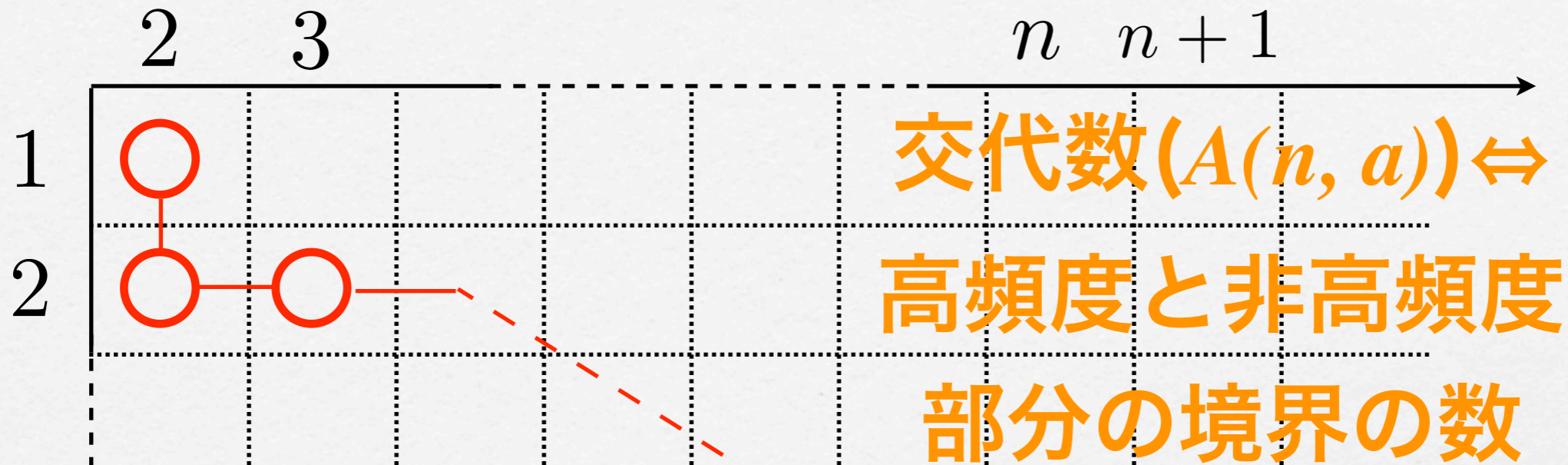
→ 高頻度部分を自動で検出する

□ 高頻度  $\Leftrightarrow$  長さ  $n$  の部分文字列の上位  $a\%$

□ 交代数を用いて  $(n, a)$  を自動的に決定

□ 様々な言語のHTML/XMLファイルに対し  
高精度でテンプレート抽出に成功

# 最適な $(n, a)$ の決定



- $(n+1, a+1)$  との比較は?
- パーセントの妥当性?
- 最小値選択の妥当性?

最も小さい交代数に  
なるように遷移

$a + 1$

# パターン言語



F172&NL154

# パターン言語

□  $\Sigma$  : 定数、  $V$  : 変数

□ **パターン**  $\Leftrightarrow \Sigma \cup V$ 上の文字列

□ 例 :  $p = axbyax (\Sigma = \{a, b\}, V = \{x, y\})$

□ **代入**  $\Leftrightarrow \theta : V \rightarrow \Sigma^*$

□ パターン  $p$  の **言語**  $L(p) = \{w \in \Sigma^* \mid \exists \theta; p\theta = w\}$


# パターン言語の学習

## □ パターン言語の学習

□ 入力：(有限)文字列の集合 $S$

□ 問題：以下を満たすパターン $p$ を見つけよ

$$S \subseteq L(p) \quad \forall S; p = x \rightarrow S \subseteq L(p) = \Sigma^*$$

$$\exists q; L(q) \subseteq L(p)$$


# パターン言語学習の難しさ

- 1変数パターン (Angluin)
- Regular Pattern (Shinohara)
  - 各変数は高々1回しか現われない
- $k$ -variableパターン (Kearns & Pitt)
  - ここまで制限すると多項式時間学習可能

# テンプレート発見問題



# テンプレート

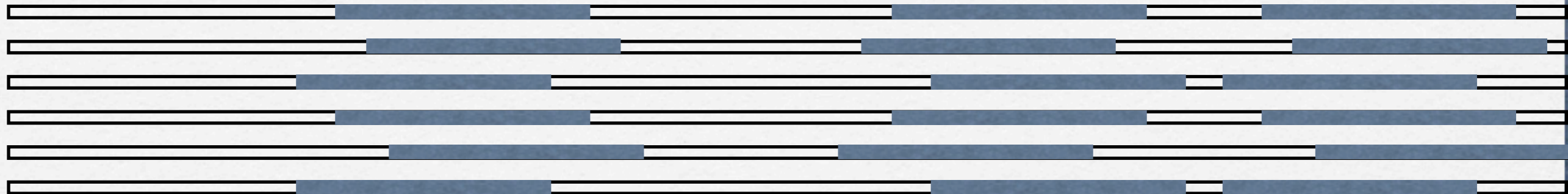
□  $p = w_1 x_1 \cdots w_m x_m$  : パタン

□ テンプレート  $\Leftrightarrow t = (w_1, \dots, w_m)$

□  $t$  の幅  $\Leftrightarrow \min\{|w_i| \mid w_i \in t\}$



# 共通パターン発見問題



—  
定数部分



# 位置ごとの頻度

頻度

産経新聞

—HTMLファイル50個

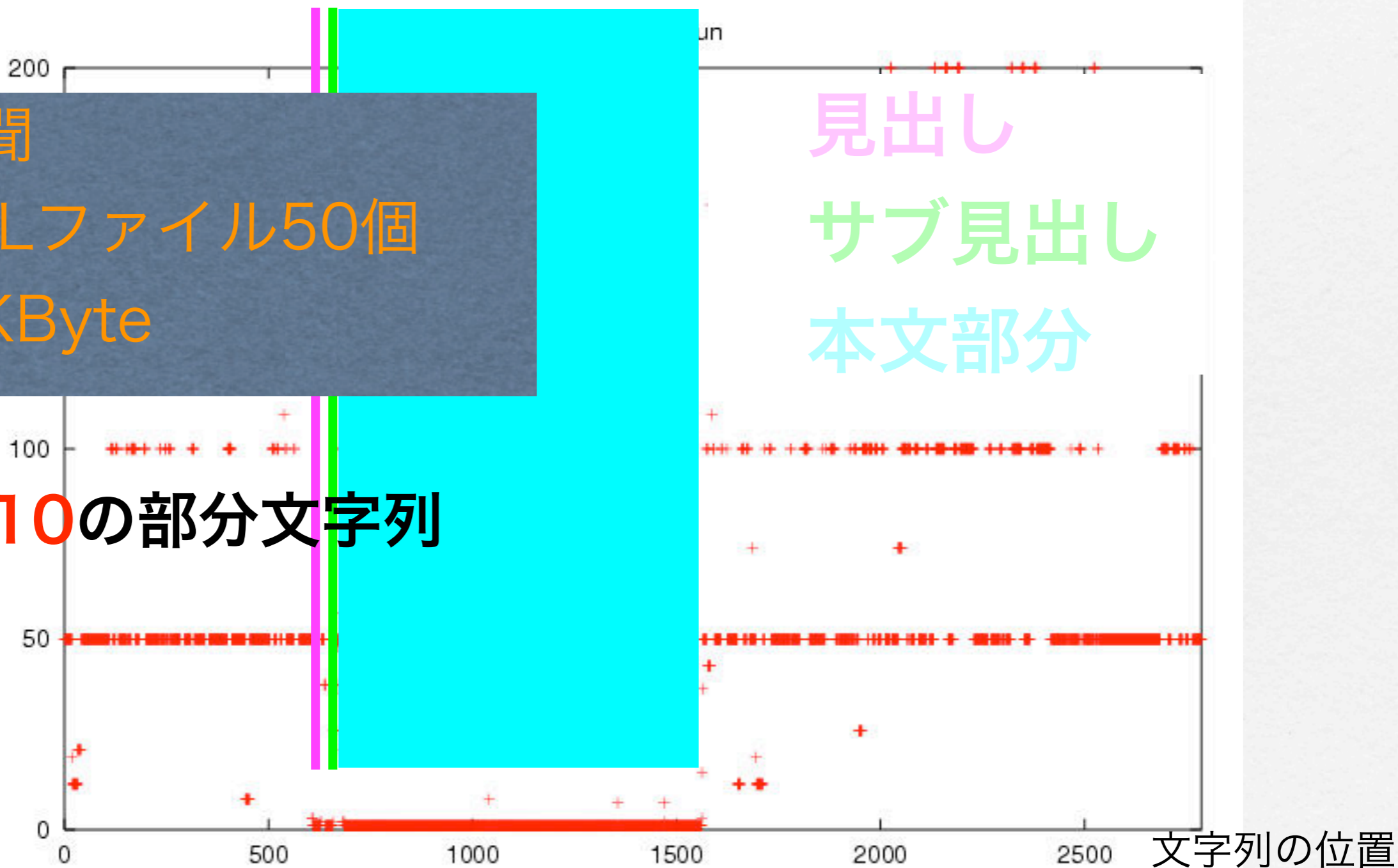
—328KByte

見出し

サブ見出し

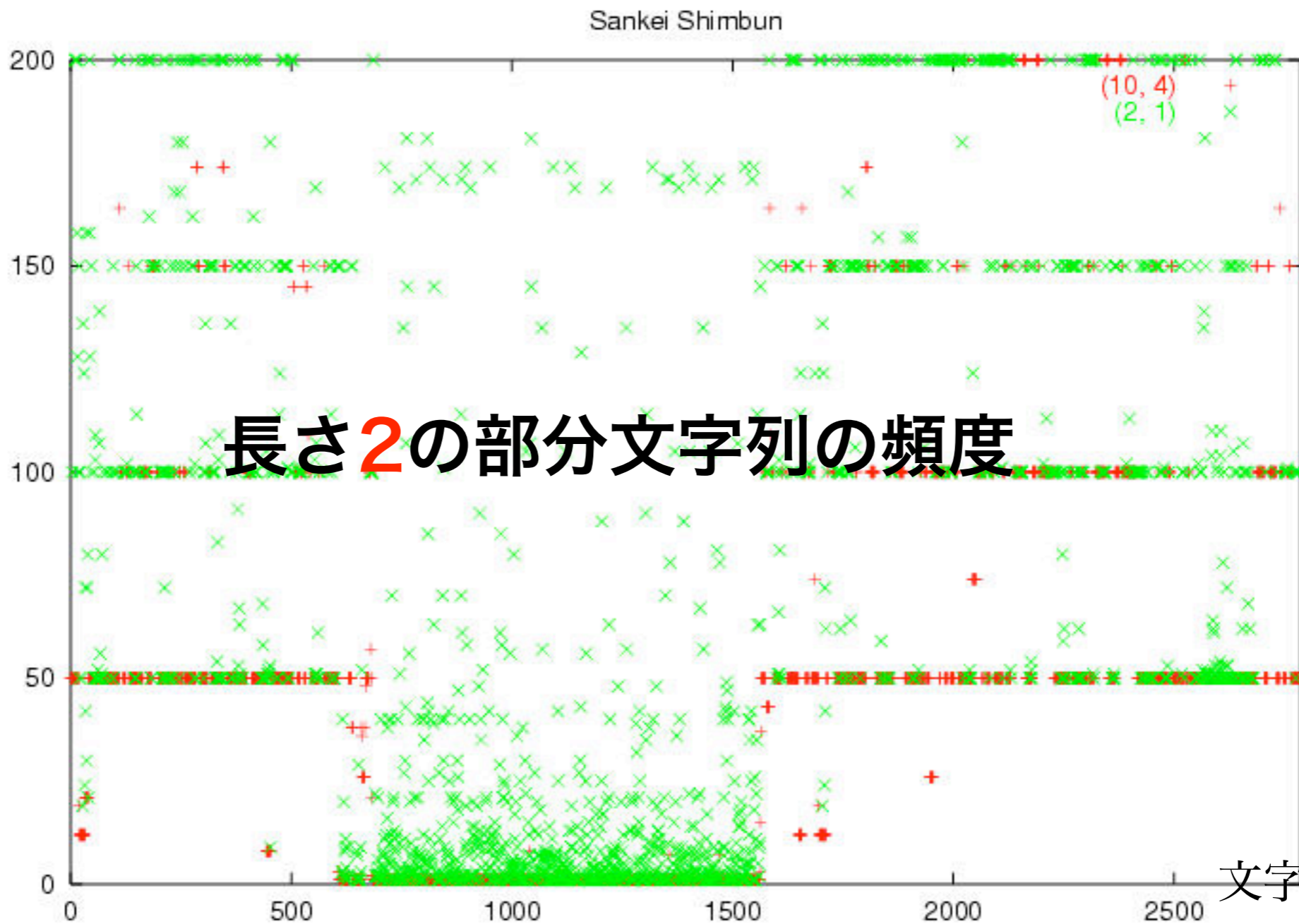
本文部分

長さ10の部分文字列



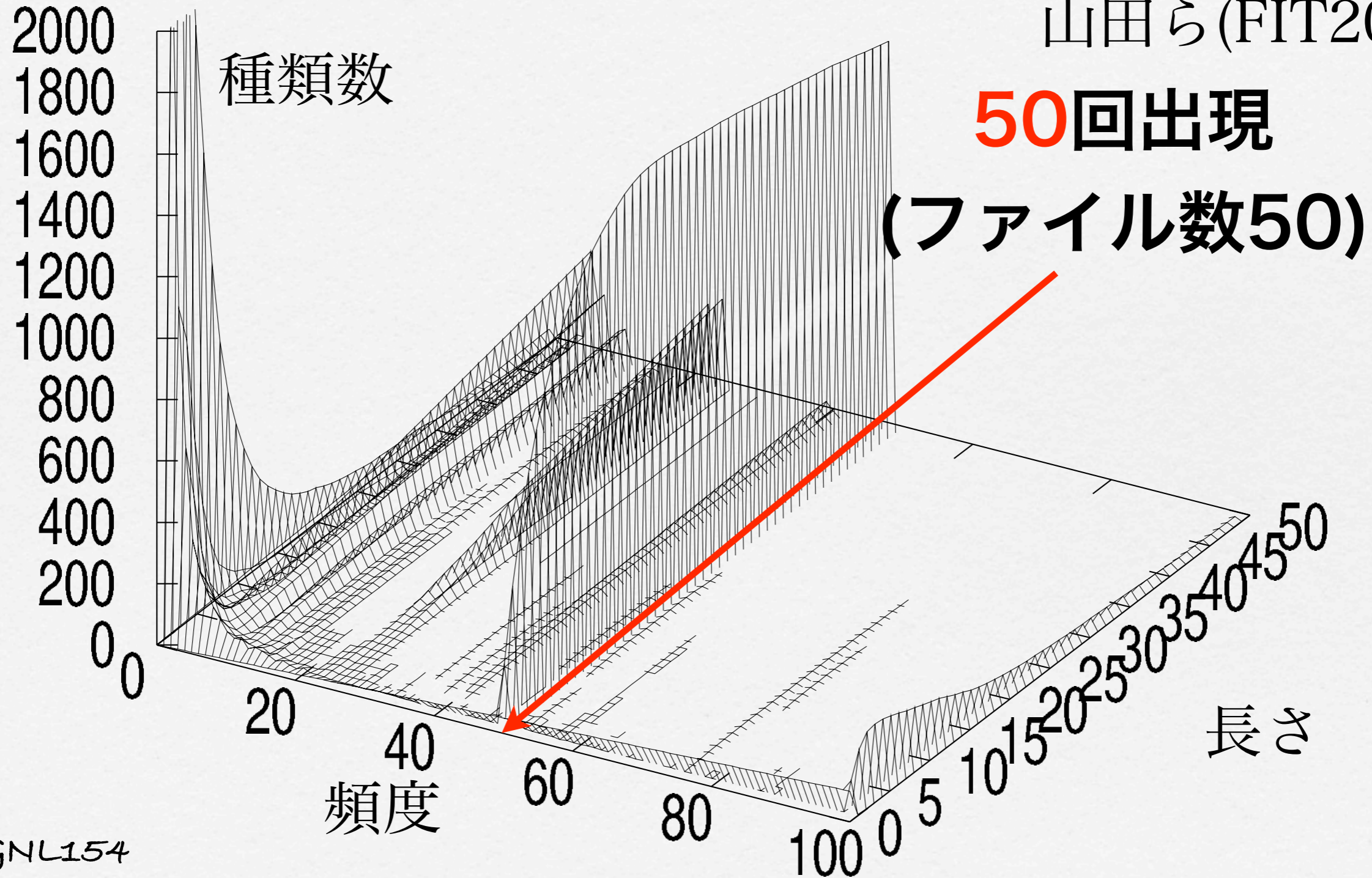
# 位置ごとの頻度 (Cont.)

頻度



# 任意の長さ拡張

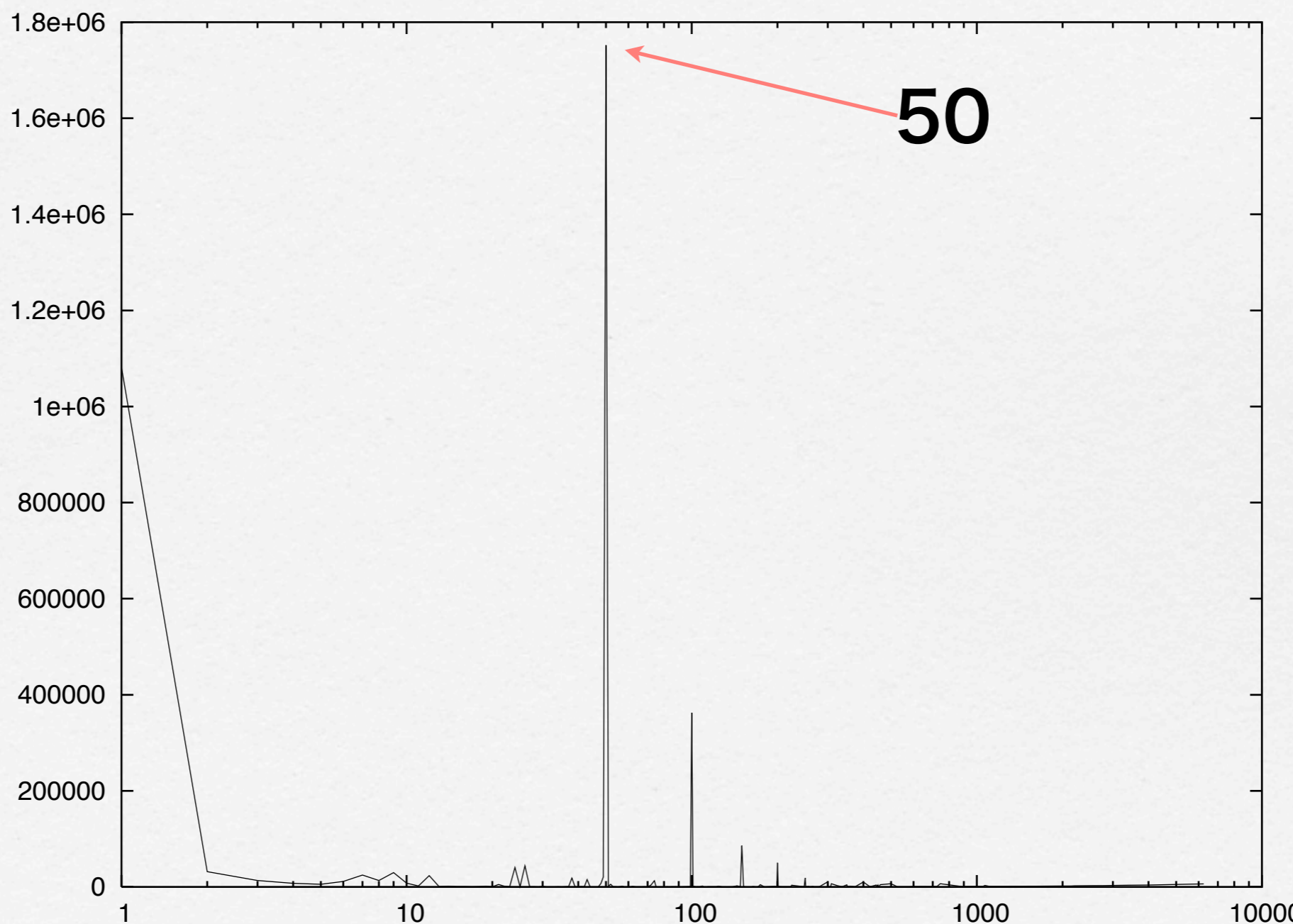
山田ら (FIT2003)



$f$  vs.  $F(f)$

$$F(f) = fV(f)$$

$F(f)$



頻度  $f$

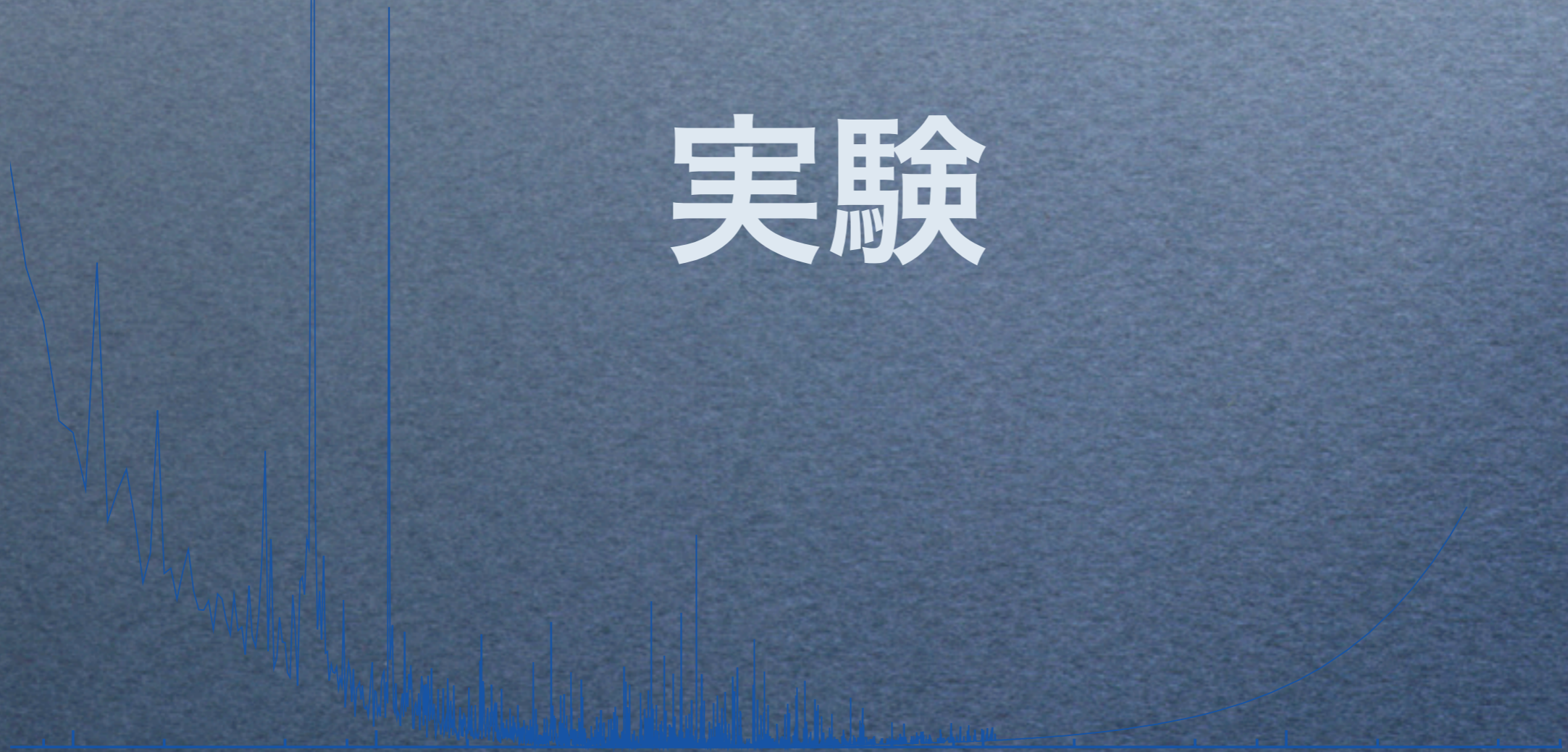
# テンプレート発見問題

- 入力：文字列の有限集合  $S$ 
  - $S$ は未知の正規パターン  $p$  から **自然な** 確率分布に従って生成された
- 問題：以下を満たすテンプレート  $t$ 
  - $t$ は正規パターン  $p$ のテンプレートである
  - $S \subseteq L(p)$  かつ  $\nexists q; L(q) \subseteq L(p)$

# アルゴリズム

- 任意の長さの部分文字列を数える
  - 実際には10~30程度の長さまで
- すべての頻度 $f$ に対し、ちょうど $f$ 回出現する文字列の数 $V(f)$ を計算
- $f$ を横軸(対数)、 $F(f)$ を縦軸にプロット

# 実験



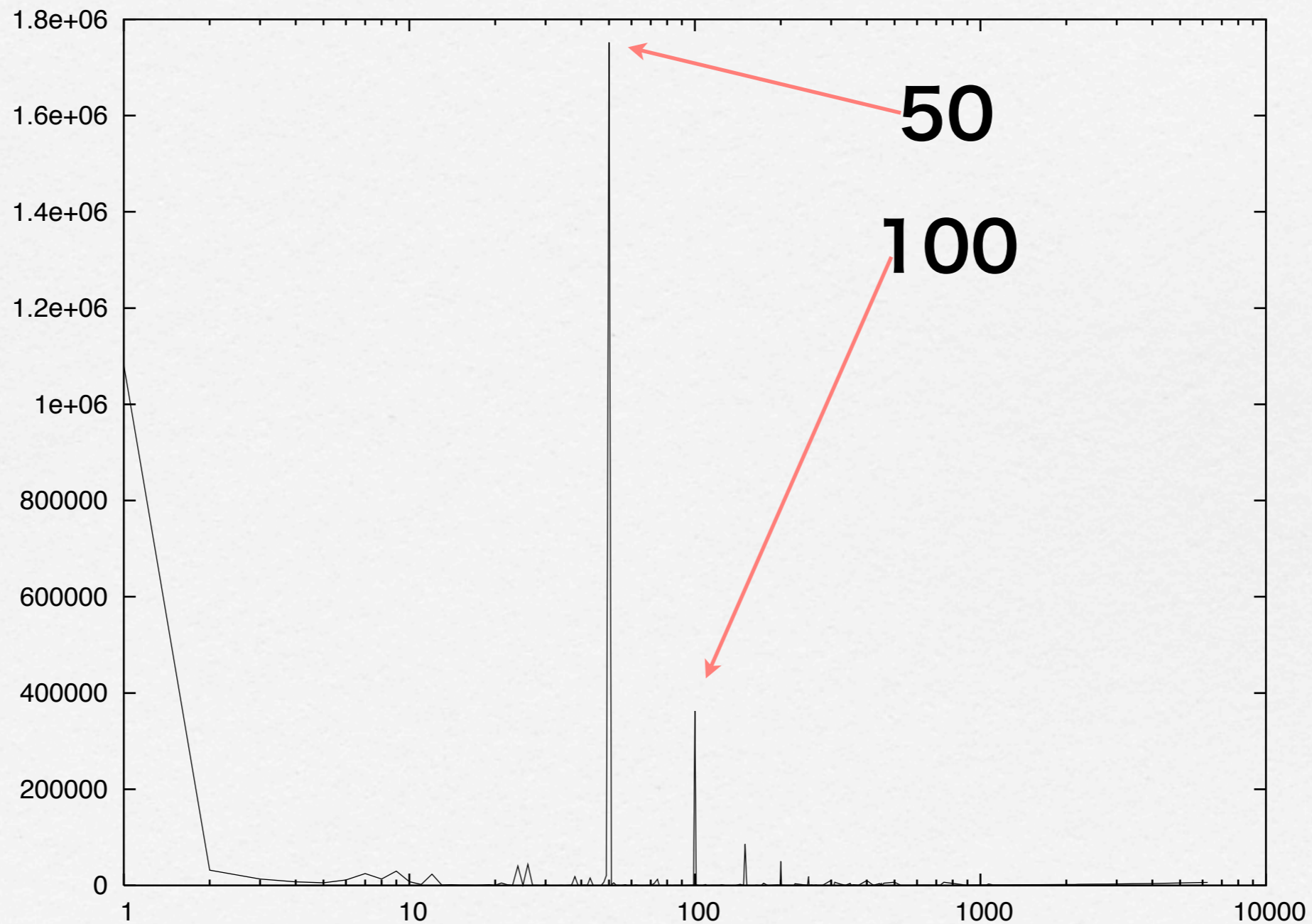
F1729NL154



# 実験データ

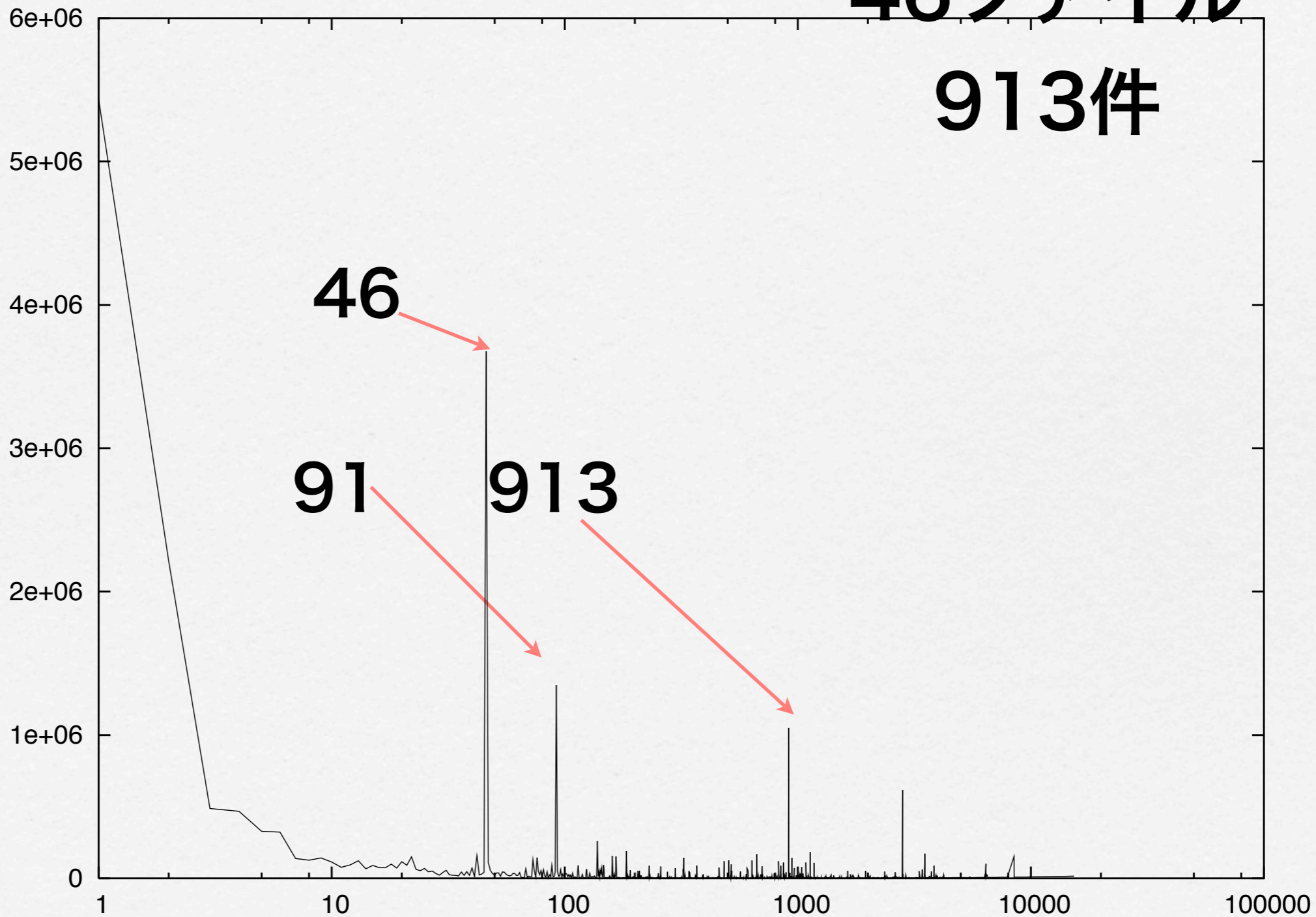
- (複数の)HTMLファイル
  - 前処理は行なわない
- 3種類のデータセット
  - 単一サイトのファイル
  - 単一サイトの周期の異なるテンプレート
  - 複数サイトのファイル

# 産経新聞の記事ファイル

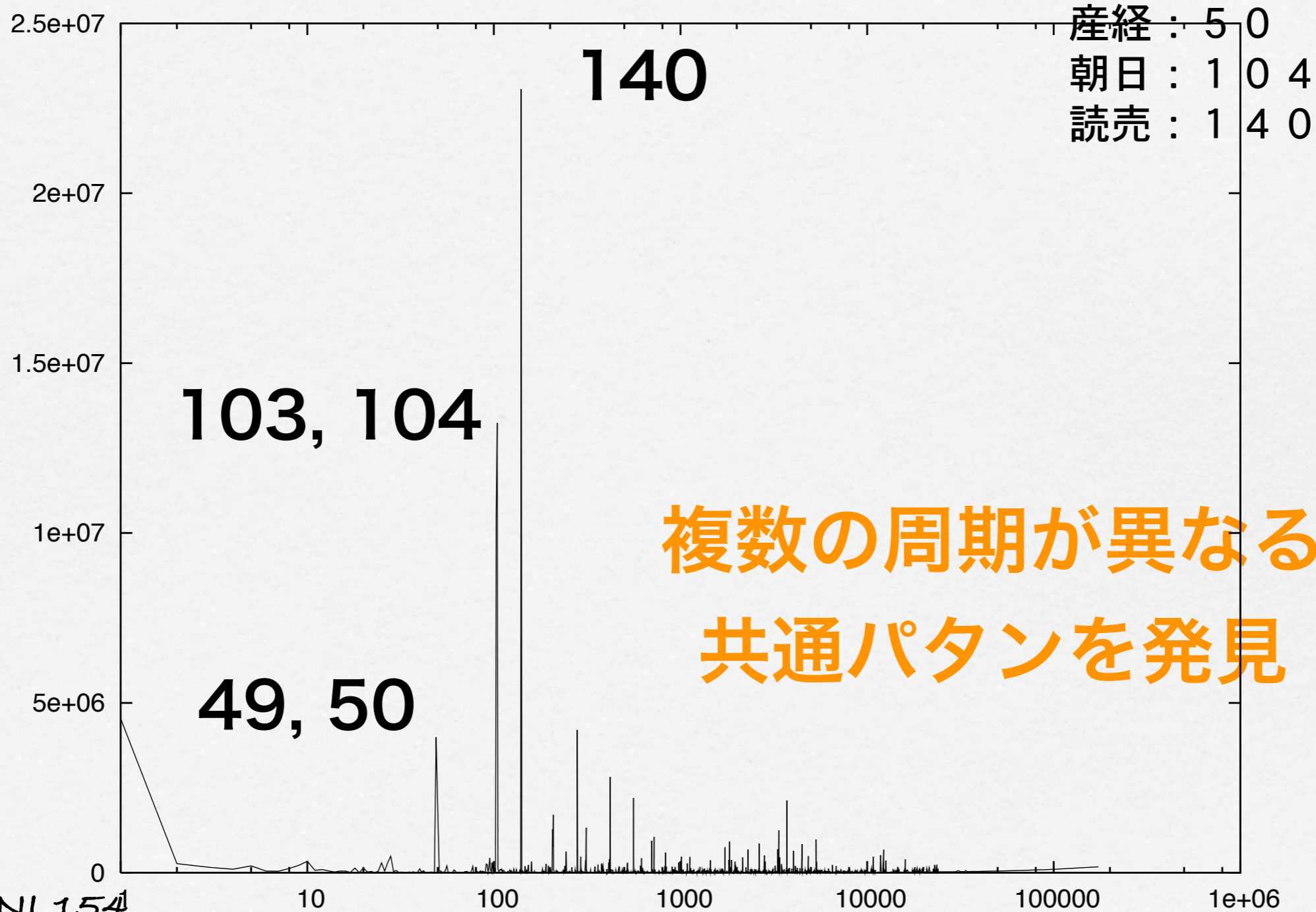


# Yahoo!検索結果 46ファイル

913件



# 3サイトのニュース記事



# 九大サイト内

九州大学\_\_Kyushu University\_\_

九州大学\_\_Kyushu Universit...

九州大学  
Kyushu University

日本語 English

f 学部・大学院  
faculty graduate school

C 教育・学生生活  
campus life

r 研究  
research

S 社会との連携  
society

行事・イベント

入試情報

国際交流・留学

教官情報

図書館

博物館

病院

同窓会・後援会

改革・評価

新キャンパス

広報誌

教職員向け情報

マップ

九州大学ニュース  
更新日: 2003年9月26日

- ▶ 中村哲氏講演会が10月14日に開催
- ▶ 農学研究院が「科学技術振興機構」から研究費を定締結(8/9)
- ▶ 細田大 来学(8/6)
- ▶ オープンキャンパスに8,000人(8/6)
- ▶ 日産のカルロス・ゴーン社長へ 名誉博士号を授与(8/1)
- ▶ 平成15年度「21世紀COEプログラム」に4拠点採択(7/17)
- ▶ ロボカップ世界大会で九大・福大連合チーム優勝!(7/11)
- ▶ 第42回七大戦 結団式開催(7/7)
- ▶ 一年生に 梶山総長が「講義」(7/2)
- ▶ 大島造船所と包活連携を推進(7/1)

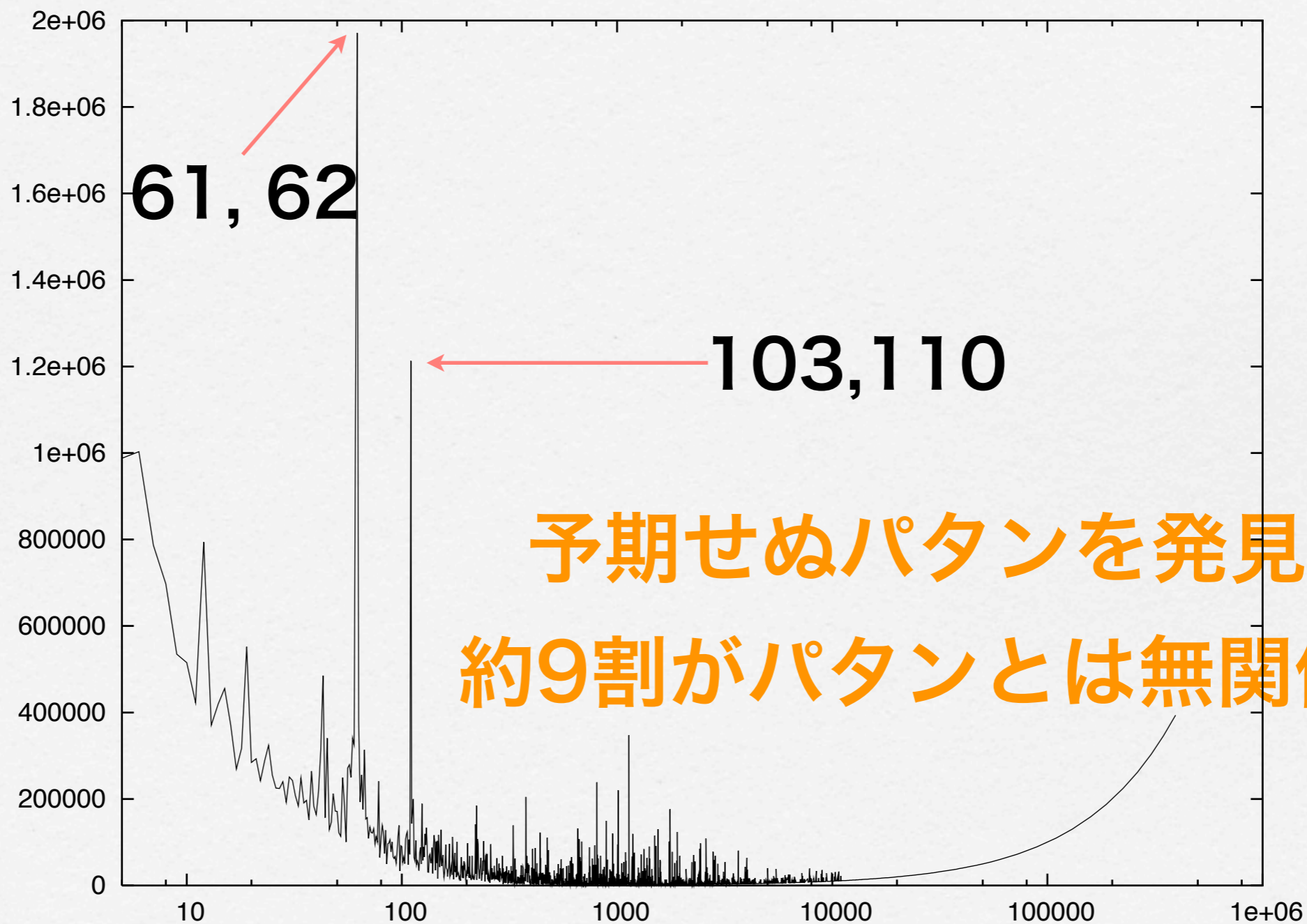
1 最初のページ 前のページ 次のページ 最新のページ

Copyright 2003 kyushu University. All rights reserved.

九大トップページから深さ3まで  
リンクをたどり598ファイルを収集

<http://www.kyushu-u.ac.jp/>

# 九大サイト内(Cont.)



予期せぬパターンを発見  
約9割がパターンとは無関係

# まとめ



FI72&NL154

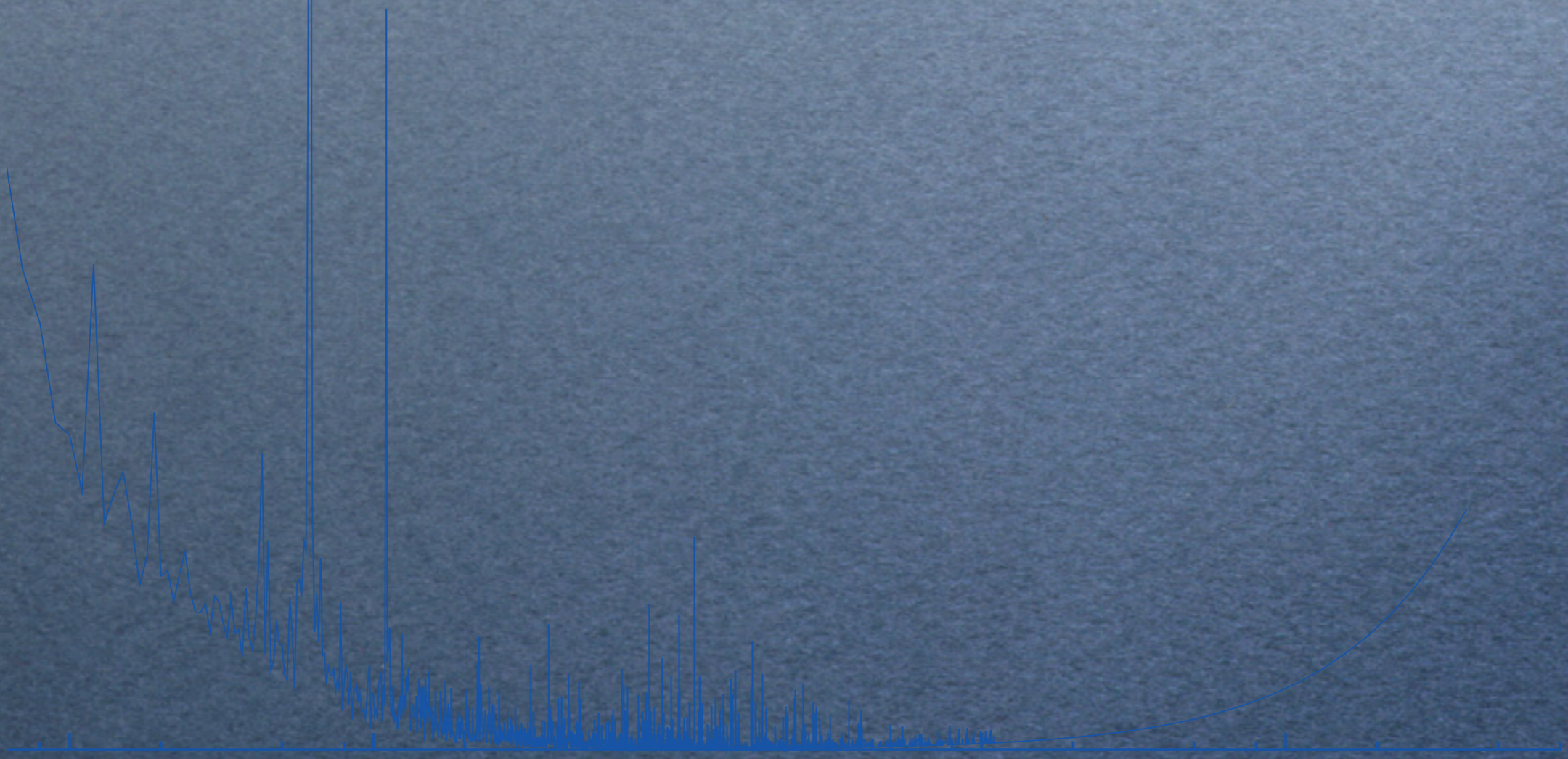
# まとめ

- テンプレート発見問題の定式化
- 共通パターン発見手法の提案
  - 正例のみで動く
  - 入力に対する前提や背景知識が不要
- HTMLファイルによる実験
  - 複数テンプレート
    - 周期の異なるテンプレートの混在も可能
  - ノイズに対し頑健



# 今後の課題

- パタン抽出の自動化
  - 複数の候補出力や対話的操作
- 「パタン」の拡張
  - 近似文字列や異周期テンプレートを包含
    - 後者は実験的にうまくいくことが示されている
- 精度の理論的な評価



FI72&NL154

# 機械学習

- パタン言語の学習
  - Angluin (1980)
  - 変数の出現を制限する
    - 共通部分よりも変数の対応
  - 負例(背景)をつかうものもあり

# 文字列処理

"共通問題"は部分列

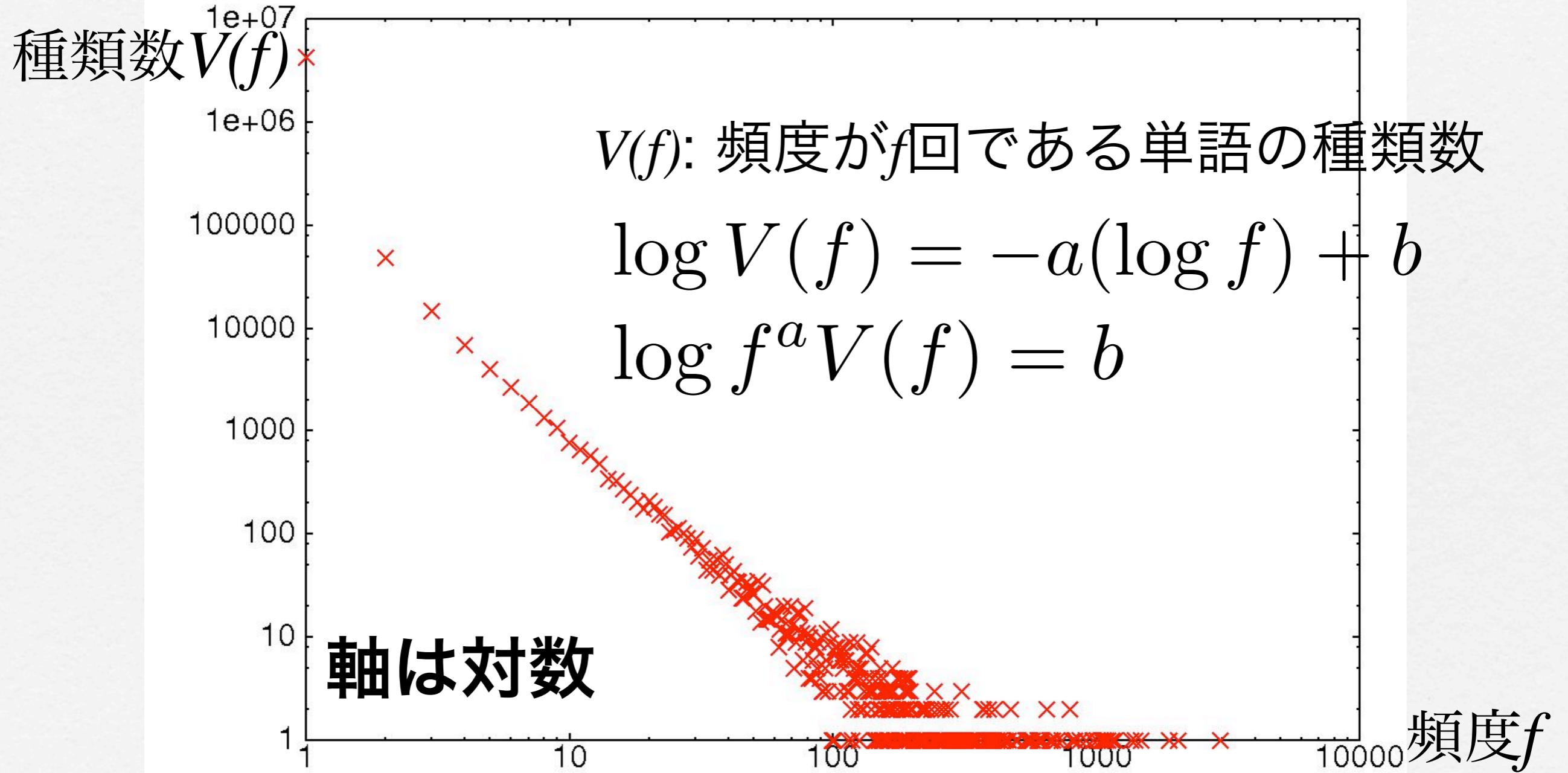
- 最長**共通**部分列**問題**

- 入力：(複数の)文字列

- 出力：入力の共通部分列で最長のもの

NP完全(Maier 1978)

# Zipfの法則との関連





[Mail to Web master](#)

[サイトマップ](#)

キーワードを  
入力してください

[日本語](#)

[English](#)

[ホーム > 病院](#)

**f** 学部・大学院  
faculty graduate school

**C** 教育・学生生活  
campus life

**r** 研究  
research

**S** 社会との連携  
society

[行事・イベント](#)

[入試情報](#)

[国際交流・留学](#)

[教官情報](#)

[図書館](#)

[博物館](#)

[病院](#)

[同窓会・後援会](#)

[改革・評価](#)

[新キャンパス](#)

[広報誌](#)

[教職員向け情報](#)

[マップ](#)

## 病 院

[トップページに戻る](#)

- 九州大学の病院
  - ▶ 医学部附属病院
  - ▶ 歯学部附属病院
  - ▶ 生体防御医学研究所附属病院

[ページのトップへ](#)

[トップページに戻る](#)