# Approximate Matching

Ben Langmead

**Department of Computer Science**

# Approximate matching



Strings

*Online* algorithms          *Offline* algorithms

Naive exact matching          Index-assisted

Boyer-Moore          *k*-mer index

We have focused on *exact* matching...

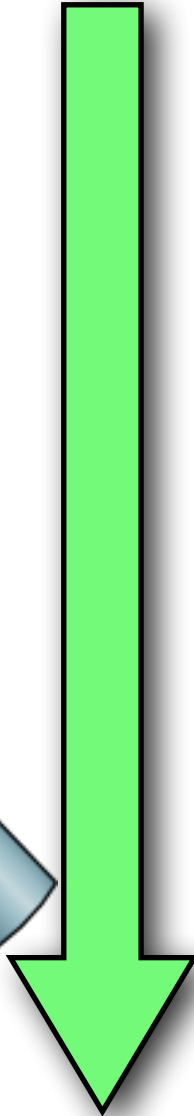... in reality, we have to deal with *differences*

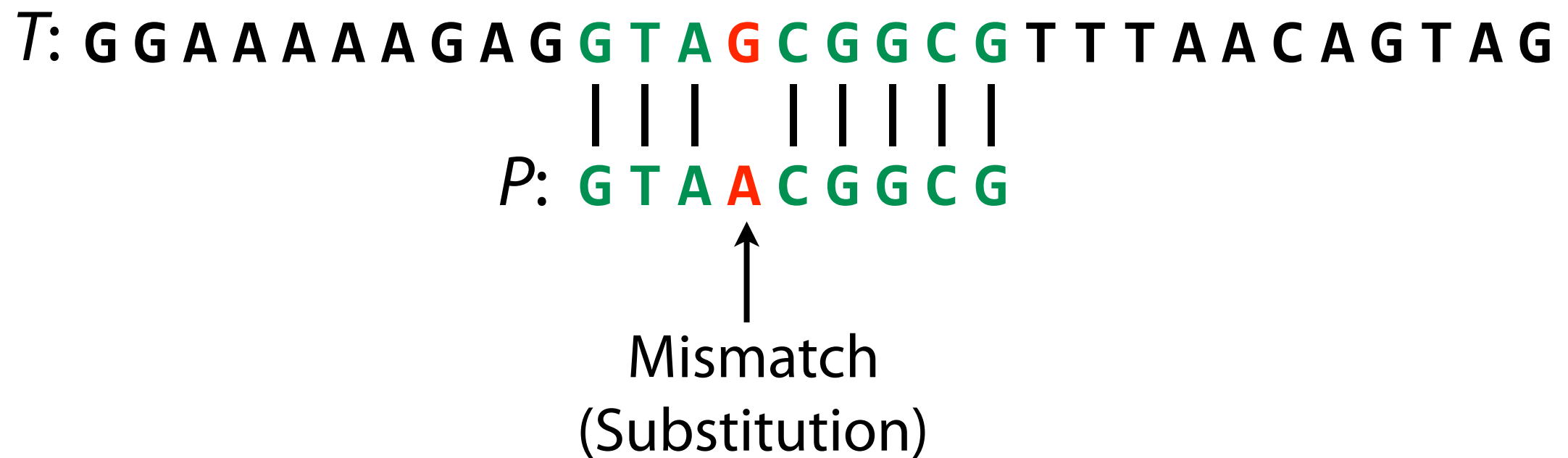# Read

CTCAAACTCCTGACCTTTGGTGATCCACCCGCCTNGGCCTTC

# Reference

GATCACAGGTCTATCACCCTATTAACCACTCACGGGAGCTCTCCATGCATTTGGTATTTT
CGTCTGGGGGGTATGCACGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCCTATGTC
GCAGTATCTGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT
ACAGGCGAACATACTTACTAAAGTGTGTTAATTAATTAATGCTTGTAGGACATAATAATA
ACAATTGAATGTCTGCACAGCCACTTTCCACACAGACATCATAACAAAAAATTTCCACCA
AACCCCCCCTCCCCCGCTTCTGGCCACAGCA────TAAACACA──CTCTGCCAAACCCCAAAA
ACAAAGAACCCTAACACCAGCCTAACCA──ATTTCAAATTT─TAT─CTTTGGCGGTATGCAC
TTTTAACAGTCACCCCCCAACTAACA──ATTATTTTCCCCTCCCAC────CATACTACTAAT
CTCATCAATACAACCCCCGCCCATC──TACCCAGCACACACACACCG─T──CTAACCCCATA
CCCCGAACCAACCAAACCCCAAAG──CACCCCCCACAGTTTATGTAGC──T──ACCTCCTCAAA
GCAATACACTGACCCGCTCAAAC──CCTGGATTTTGGATCCACCCAGCG─G──TTGGCCTAAA
CTAGCCTTTCTATTAGCTCTTAG──AAGATTACACATGCAAGCATCCCCG──TCCAGTGAGT
TCACCCTCTAAATCACCACGATC──AAAGGAACAAGCATCAAGCACGCAG──AATGCAGCTC
AAAACGCTTAGCCTAGCCACACC──CCACGGGAAACAGCAGTGATTAACC──TTAGCAATAA
ACGAAAGTTTAACTAAGCTATACT──ACCCCCAGGGTTGGTCAATTTCGT─CCAGCCACCGC
GGTCACACGATTAACCCAAGTCAAT──GAAGCCGGCGTAAAGAGTGTT─TAGATCACCCCC
TCCCCAATAAAGCTAAAACTCACCTGA──TTGTAAAAAACTCCAGT──ACACAAAATAGAC
TACGAAAGTGGCTTTAACATATCTGAACA──ACAATAGCTAAGACC──GGGATTAGA
TACCCCACTATGCTTAGCCCTAAACCTCAACAGT──TAAACAACC──GCCAGAA
CACTACGAGCCACAGCTTAAAACTCAAAGGACCTGGCGGTGCTTCAT──TAGAGG
AGCCTGTTCTGTAATCGATAAACCCCGATCAACCTCACCACCTCTTGCT──TATA
CCGCCATCTTCAGCAAACCCTGATGAAGGCTACAAAGTAAGCGCAAGTACG──AG
ACGTTAGGTCAAGGTGTAGCCCATGAGGTGGCAAGAAATGGGCTACATTTTC──
AAAACTACGATAGCCCTTATGAAACTTAAGGGTCGAAGGTGGATTTAGCAGTAA─
AGTAGAGTGCTTAGTTGAACAGGGCCCTGAAGCGCGTACACACCGCCCGTCACCC─
AAGTATACTTCAAAGGACATTTAACTAAAACCCCTACGCATTTATATAGAGGAGACAA─
CGTAACCTCAAACTCCTGCCTTTGGTGATCCACCCGCCTTGGCCTACCTGCATAATGAAG
AAGCACCCAACTTACACTTAGGAGATTTCAACTTAACTTGACCGCTCTGAGCTAAACCTA
GCCCCAAACCCACTCCACCTTACTACCAGACAACCTTAGCCAAACCATTTACCCAAATAA
AGTATAGGCGATAGAAATTGAAACCTGGCGCAATAGATATAGTACCGCAAGGGAAAGATG
AAAAATTATAACCAAGCATAATATAGCAAGGACTAACCCCTATACCTTCTGCATAATGAA
TTAACTAGAAATAACTTTGCAAGGAGAGCCAAAGCTAAGACCCCCGAAACCAGACGAGCT
ACCTAAGAACAGCTAAAAGAGCACACCCGTCTATGTAGCAAAATAGTGGGAAGATTTATA
GGTAGAGGCGACAAACCTACCGAGCCTGGTGATAGCTGGTTGTCCAAGATAGAATCTTAG
TTCAACTTTAAATTTGCCCACAGAACCCTCTAAATCCCCTTGTAAATTTAACTGTTAGTC
CAAAGAGGAACAGCTCTTTGGACACTAGGAAAAAACCTTGTAGAGAGAGTAAAAAATTTA
ACACCCATAGTAGGCCTAAAAGCAGCCACCAATTAAGAAAGCGTTCAAGCTCAACACCCA
CTACCTAAAAAATCCCAAACATATAACTGAACTCCTCACACCCAATTGGACCAATCTATC
ACCCTATAGAAGAACTAATGTTAGTATAAGTAACATGAAAACATTCTCCTCCGCATAAGC
CTGCGTCAGATTAAAACACTGAACTGACAATTAACAGCCCAATATCTACAATCAACCAAC
AAGTCATTATTACCCTCACTGTCAACCCAACACAGGCATGCTCATAAGGAAAGGTTAAAA
AAAGTAAAAGGAACTCGGCAAATCTTACCCCGCCTGTTTACCAAAAACATCACCTCTAGC
ATCACCAGTATTAGAGGCACCGCCTGCCCAGTGACACATGTTTAACGGCCGCGGTACCCT

Sequence differences occur because of...

1. Sequencing error

2. Genetic variation

# Approximate matching

*T*: **G G A A A A A G A G G T A G C G G C G T T T A A C A G T A G**

**| | |   | | | | |**

*P*: **G T A A C G G C G**

Mismatch
(Substitution)

# Approximate matching

*T*: **G G A A A A A G A G G T A G C – G C G T T T A A C A G T A G**

```
          | | | | |   | | |
```

*P*: **G T A G C G G C G**

Insertion

# Approximate matching

$T$:  G G A A A A A G A G G T A G C G G C G T T T A A C A G T A G

$P$:  G T – G C G G C G

Deletion

# Hamming distance

For *X* & *Y* where $|X| = |Y|$, *hamming distance =*
minimum # substitutions needed to turn one into the other

*X*: G A G G T A G C G G C G T T

\|  \| \| \| \|  \| \| \|  \| \| \|

*Y*: G T G G T A A C G G G G T T

*Hamming distance = 3*

# Edit distance

(AKA Levenshtein distance)

For *X* & *Y*, *edit distance* = minimum # edits (substitutions, insertions, deletions) needed to turn one into the other

*X:* T G G C C G C G C A A A A A C A G C

*Edit distance = 2*

*Y:* T G A C C G C G C A A A A – C A G C

*X:* G C G T A T G C G G C T A – A C G C

*Edit distance = 2*

*Y:* G C – T A T G C G G C T A T A C G C

# Approximate matching

Like exact matching, but *pattern P* may be within a certain *distance* (usually Hamming or edit) of *T*.  Each such place is an *approximate match*.

Allowing edits is more challenging than just allowing mismatches

We'll return to edits

# Approximate matching

```python
def naive(p, t):
    occurrences = []
    for i in range(len(t) - len(p) + 1):   # Loop over alignments
        match = True
        for j in range(len(p)):            # Loop over characters
            if t[i+j] != p[j]:             # compare characters
                match = False              # mismatch; reject alignment
                break
        if match:
            occurrences.append(i)          # all chars matched; record
    return occurrences
```

# Approximate matching

```python
def naive_approx_hamming(p, t, maxDistance):
    occurrences = []
    for i in range(len(t) - len(p) + 1):   # loop over alignments
        nmm = 0
        for j in range(len(p)):            # loop over characters
            if t[i+j] != p[j]:             # compare characters
                nmm += 1                   # mismatch
                if nmm > maxDistance:
                    break                  # exceeded max hamming dist
        if nmm <= maxDistance:
            occurrences.append(i)          # approximate match
    return occurrences
```

http://bit.ly/CG_NaiveApprox

# Approximate matching

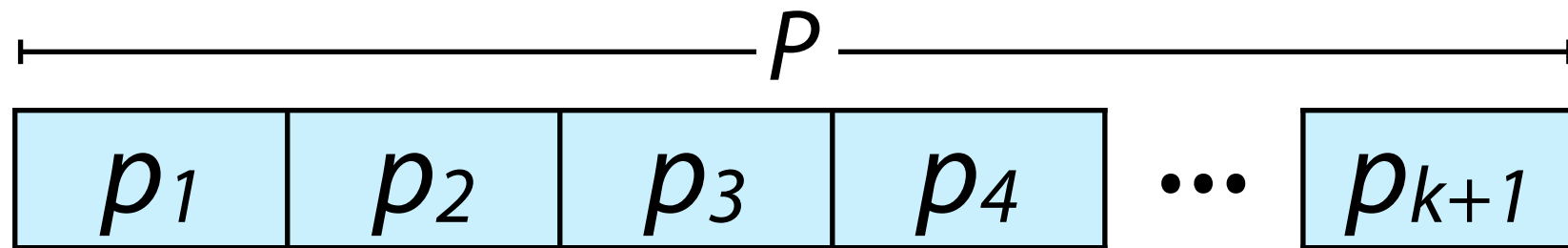**Wanted**: way to apply exact matching algorithms to approximate matching problems

# Approximate matching

# Approximate matching



If *P* occurs in *T* with up to *k* edits…

# Approximate matching

$$P$$

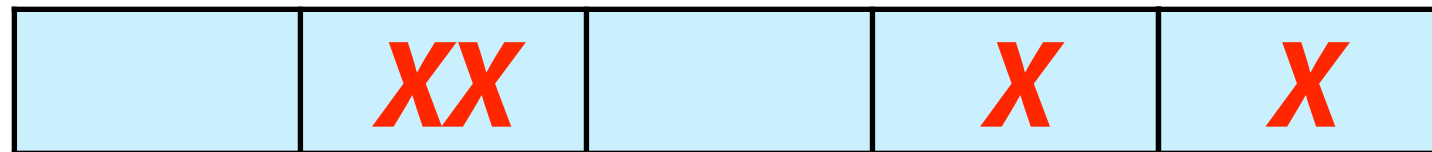| $p_1$ | $p_2$ | $p_3$ | $p_4$ | $\cdots$ | $p_{k+1}$ |

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, …, p_{k+1}$ must appear with 0 edits

# Approximate matching
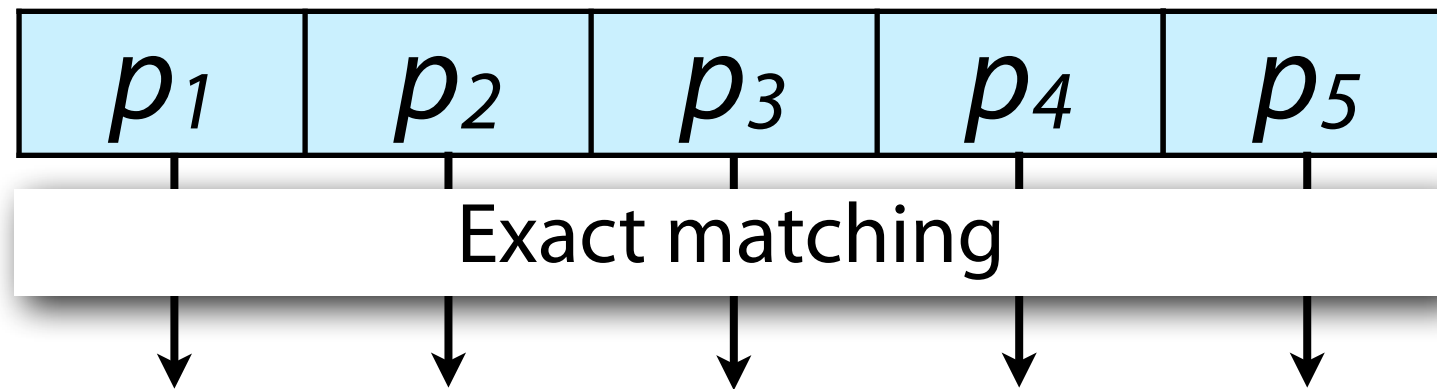


5 partitions
4 edits (*X*)

# Approximate matching



Pigeonhole principle:  *k+1* pigeons, *k* holes.
At least one has >1 pigeon!

# Approximate matching



We have *k* pigeons, *k*+1 holes, at least one…
…is *empty*

# Pigeonhole principle

| $p_1$ | $p_2$ | $p_3$ | $p_4$ | $p_5$ |
|:-----:|:-----:|:-----:|:-----:|:-----:|

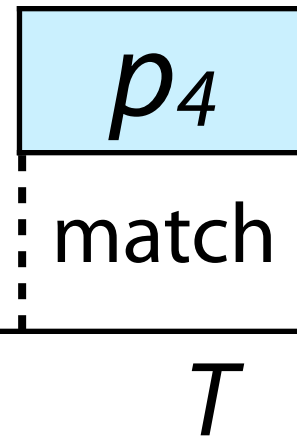Exact matching

$T$

What algorithm can we use?

*Any* exact matching algorithm

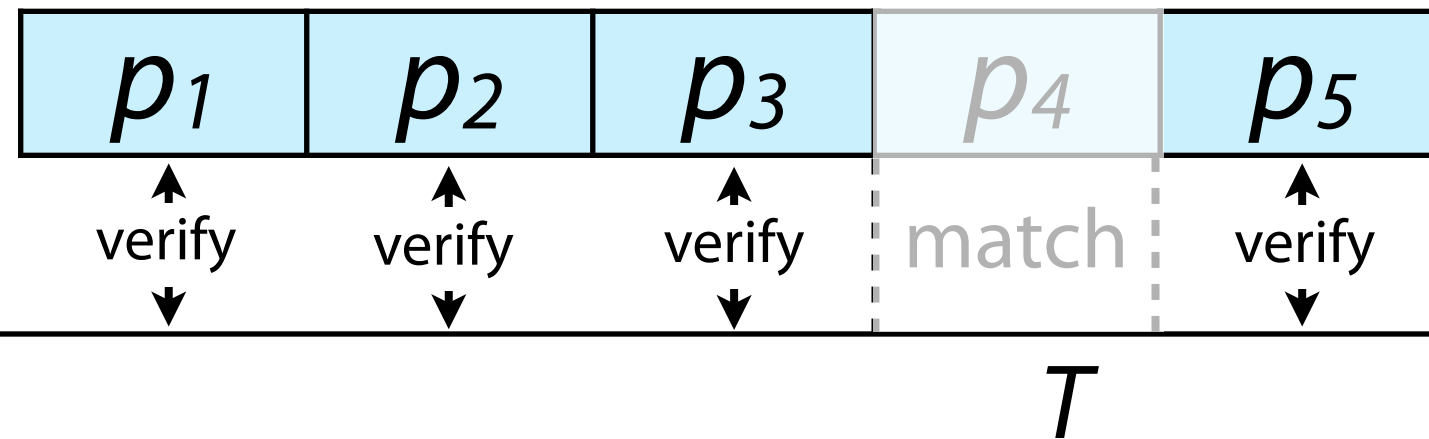If we have a k-mer index, we can use that

Naive exact matching

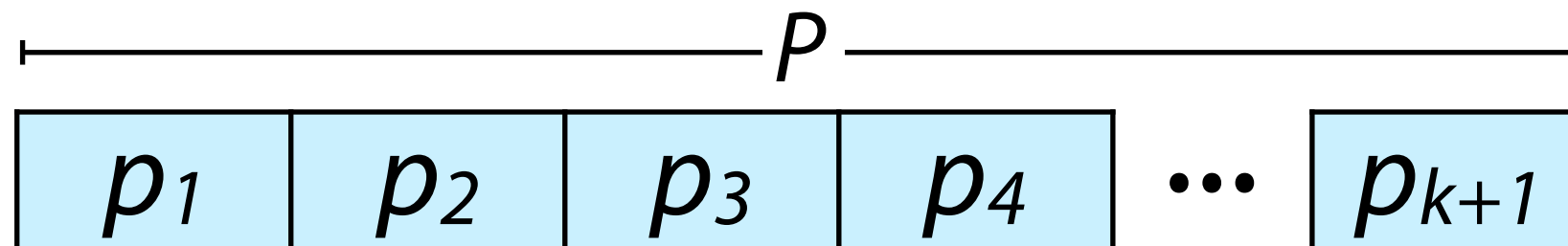Boyer-Moore

# Pigeonhole principle

# Pigeonhole principle



For Hamming distance, verification is essentially just the inner loop of **naive_approx_hamming** from before

# Pigeonhole principle

$$P$$

| $p_1$ | $p_2$ | $p_3$ | $p_4$ | $\cdots$ | $p_{k+1}$ |

**Advantages**

Reuse favorite exact matching algos; fast and easy

Flexible; works for Hamming and edit distance*

**Disadvantages**

Large $k$ yields small partitions matching many times by chance; lots of verification work

$k+1$ exact matching problems, one per partition

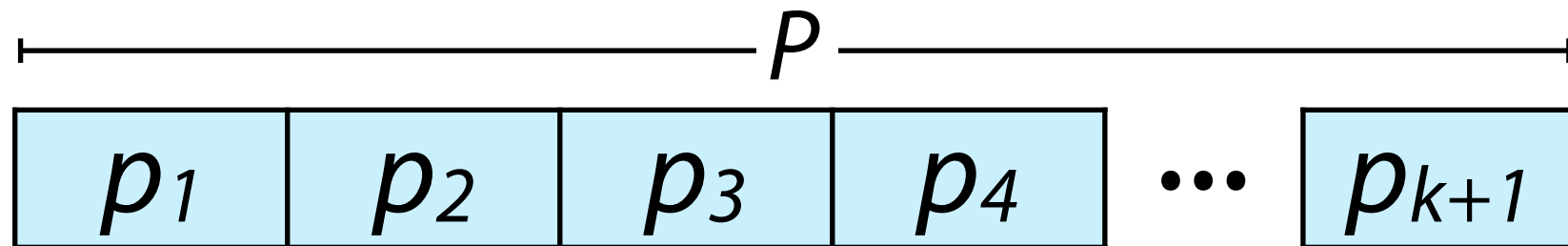* we don't know how to do edit distance verification yet

# Implementation of pigeonhole principle with Boyer-Moore as exact matching algorithm: http://j.mp/CG_ApproxBM

| | Boyer-Moore, exact | | | Boyer-Moore, ≤1 mismatch with pigeonhole | | | Boyer-Moore, ≤2 mismatches with pigeonhole | | |
|---|---|---|---|---|---|---|---|---|---|
| | # character comparisons | wall clock time | # matches | # character comparisons | wall clock time | # matches | # character comparisons | wall clock time | # matches |
| **P**: "tomorrow" <br><br> **T**: Shakespeare's complete works | 786 K | 1.91s | 17 | 3.05 M | 7.73 s | 24 | 6.98 M | 16.83 s | 382 |
| **P**: 50 nt string from Alu repeat* <br><br> **T**: Human reference (hg19) chromosome 1 | 32.5 M | 67.21 s | 336 | 107 M | 209 s | 1,045 | 171 M | 328 s | 2,798 |

\* GCGCGGTGGCTCACGCCTGTAATCCCAGCACTTTGGGAGGCCGAGGCGGG

# Generalizing pigeonhole, part 1

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, ..., p_{k+1}$ must appear with 0 edits
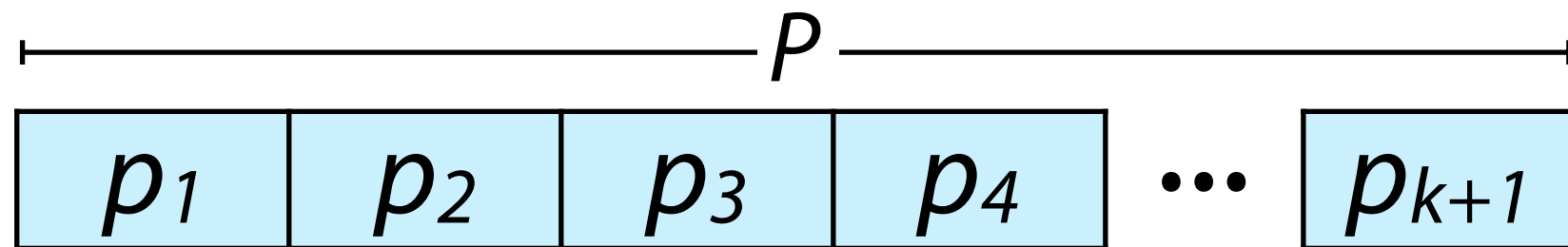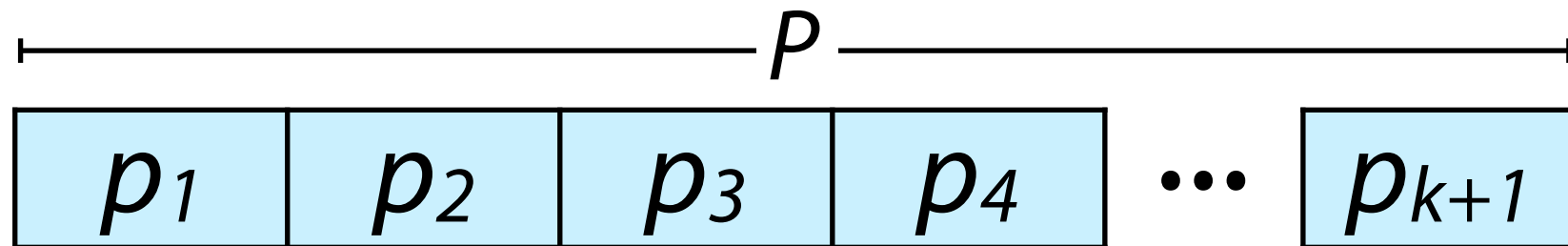


But doesn't *have to* be "at least one of" …

what would we have to change for "at least two of"?

If $P$ occurs in $T$ with up to $k$ edits, then at least <span style="color:red">two</span> of _____ must appear with 0 edits

# Generalizing pigeonhole, part 1

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, ..., p_{k+1}$ must appear with 0 edits



But doesn't *have to* be "at least one of" …
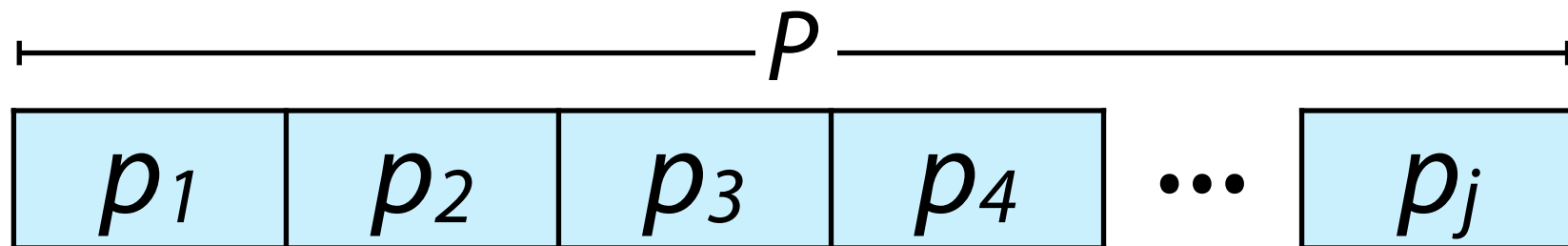
what would we have to change for "at least two of"?

If $P$ occurs in $T$ with up to $k$ edits, then at least <span style="color:red">two</span> of $p_1, p_2, ..., p_{k+2}$ must appear with 0 edits

# Generalizing pigeonhole, part 2

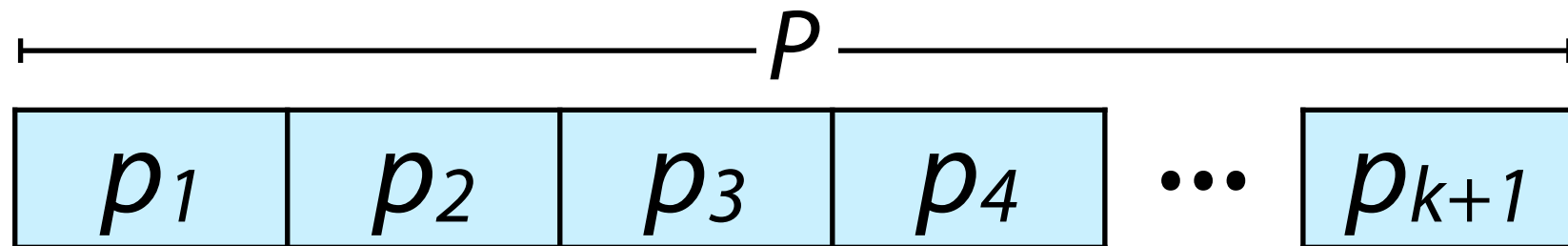If *P* occurs in *T* with up to *k* edits, then at least one of $p_1, p_2, \ldots, p_{k+1}$ must appear with 0 edits



Let $p_1, p_2, \ldots, p_j$ be a partitioning of *P*. If *P* occurs with up to *k* edits, then at least one of $p_1, p_2, \ldots, p_j$ must occur with ≤ **???** edits.
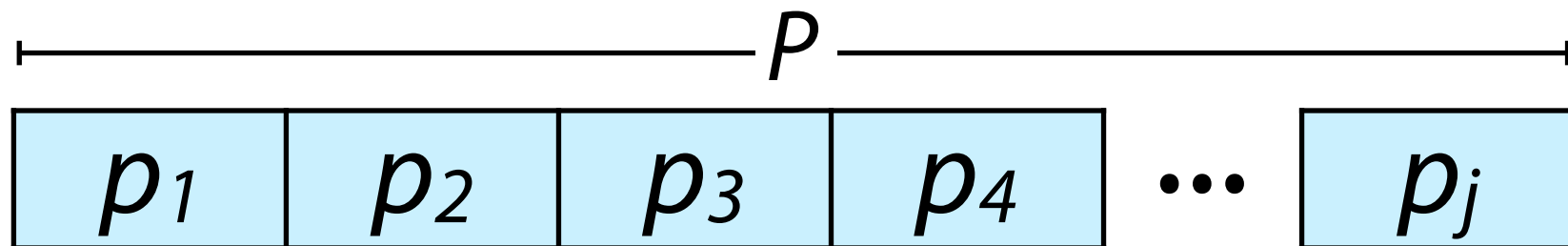
# Generalizing pigeonhole, part 2

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, ..., p_{k+1}$ must appear with 0 edits
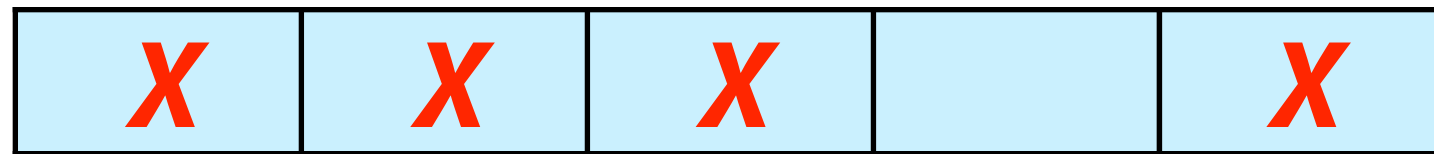


Let $p_1, p_2, ..., p_j$ be a partitioning of $P$. If $P$ occurs with up to $k$ edits, then at least one of $p_1, p_2, ..., p_j$ must occur with $\leq floor(k / j)$ edits.
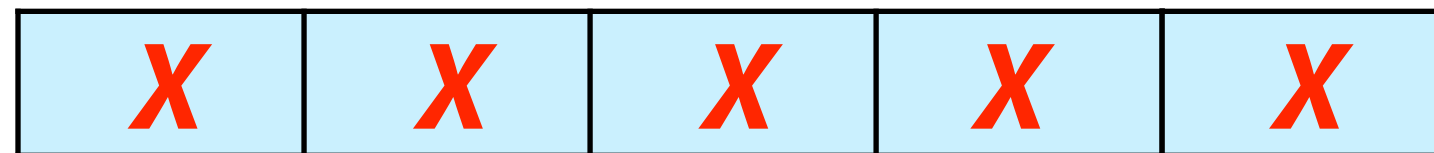
# Generalizing pigeonhole, part 2

At least one of $p_1, p_2, ..., p_5$ occurs with...

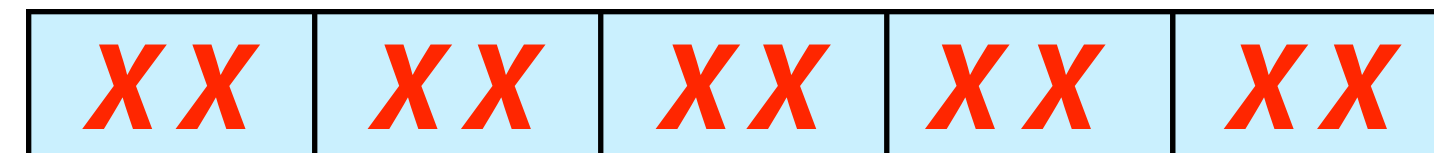| | | | | | | |
|---|---|---|---|---|---|---|
| $X$ | $X$ | $X$ | | $X$ | $k = 4$ edits | $\leq 0$ edits |
| $X$ | $X$ | $X$ | $X$ | $X$ | $k = 5$ edits | $\leq 1$ edits |
| $X\,X$ | $X\,X$ | $X$ | $X\,X$ | $X\,X$ | $k = 9$ edits | $\leq 1$ edits |
| $X\,X$ | $X\,X$ | $X\,X$ | $X\,X$ | $X\,X$ | $k = 10$ edits | $\leq 2$ edits |

$P$

etc

# Generalizing pigeonhole, part 2

**Let $j = k + 1$**

*Why?*
Smallest value s.t. $floor(k / j) = 0$

*Why make $floor(k / j) = 0$?*
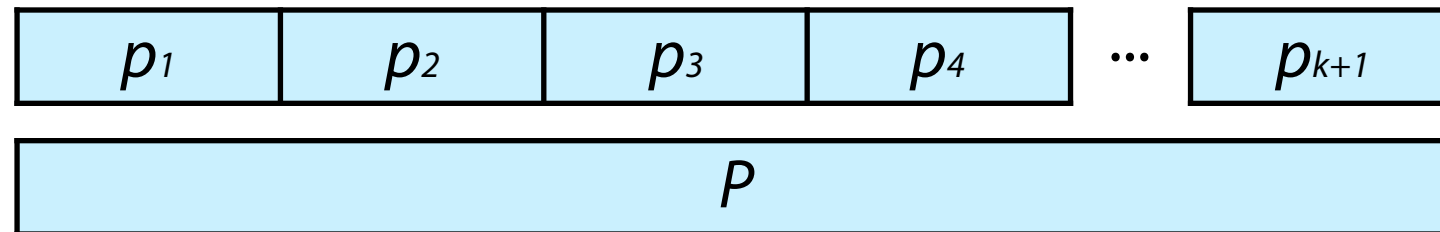So we can use exact matching

*Why is smaller j good?*
Yields fewer, longer partitions

*Why are long partitions good?*
Makes exact-matching filter more specific, minimizing # candidates

**General**

*Pigeonhole principle*

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, ..., p_j$ must appear with $floor(k / j)$ edits

**Specific**

*Pigeonhole principle with $j = k + 1$*

If $P$ occurs in $T$ with up to $k$ edits, then at least one of $p_1, p_2, ..., p_{k+1}$ must appear with 0 edits

# A different principle

We partitioned *P* into non-overlapping substrings

| $p_1$ | $p_2$ | $p_3$ | $p_4$ | ... | $p_{k+1}$ |

$P$

Now consider *overlapping* substrings

$P_3$   $P_6$   $P_{n-l+1}$

$P_2$   $P_5$   ...   $P_{n-l}$

$P_1$   $P_4$   $P_{n-l-1}$

$P$

# Approximate string matching: more principles



Say substrings are length $q$.  There are $n - q + 1$ such substrings.

1 edit to $P$ changes *at most q* substrings

Minimum # of length-$q$ substrings unedited after $k$ edits?    $n - q + 1 - kq$

$q$-*gram* lemma: if $P$ occurs in $T$ with up to $k$ edits, alignment must contain $t$ exact matches of length $q$, where $t \geq n - q + 1 - kq$

kq is *worst case; could be < kq*

# Approximate string matching: more principles

If *P* occurs in *T* with up to *k* edits, alignment contains an exact match of length *q,* where $q \geq floor(n / (k + 1))$

Obtained by solving for *q*: $\quad n - q + 1 - kq \geq 1$

Exact matching filter: find matches of length $floor(n / (k + 1))$ between *T* and *any* substring of *P*. Check vicinity for full match.

# Approximate matching principles

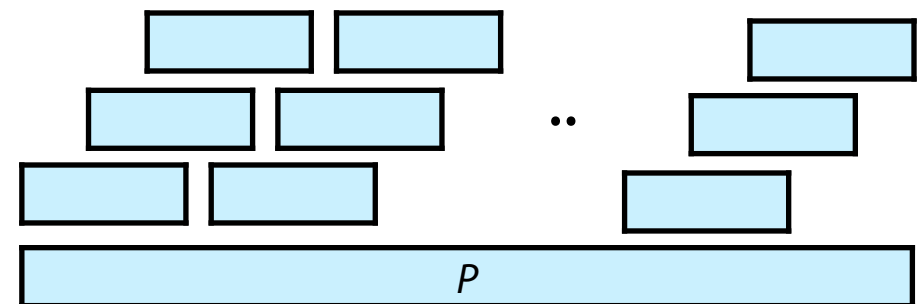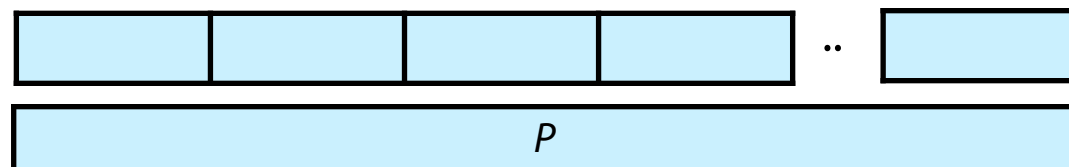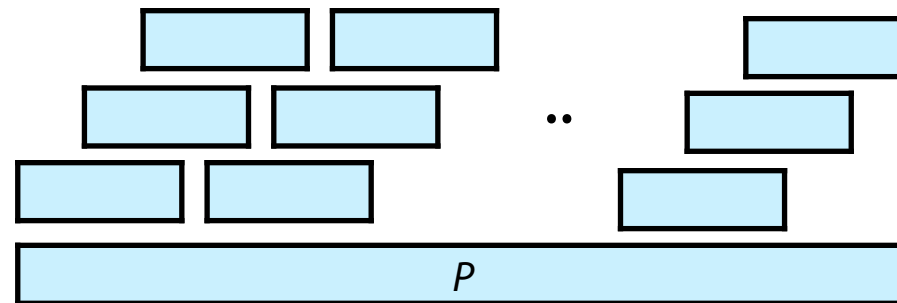|  | Non-overlapping substrings | Overlapping substrings |
|---|---|---|
| **General** | *Pigeonhole principle*<br><br>$p_1, p_2, ..., p_j$ is a partitioning of $P$. If $P$ occurs with $\leq k$ edits, at least one partition matches with $\leq floor(k / j)$ edits. | *q-gram lemma*<br><br>If $P$ occurs with $\leq k$ edits, alignment contains $t$ exact matches of length $q$, where $t \geq n - q + 1 - kq$ |
| **Specific** | *Pigeonhole principle with $j = k + 1$*<br><br>$p_1, p_2, ..., p_{k+1}$ is a partitioning of $P$. If $P$ occurrs in $T$ with $\leq k$ edits, at least one partition matches exactly. | *q-gram lemma with $t = 1$*<br><br>If $P$ occurs with $\leq k$ edits, alignment contains an exact match of length $q$ where $q \geq floor(n / (k + 1))$ |

# Sensitivity

Sensitivity = fraction of "true" approximate matches discovered by the algorithm

*Lossless* algorithm finds all of them, *lossy* algorithm doesn't necessarily

We've seen *lossless* algorithms.  Most everyday tools are *lossy*.  Lossy algorithms are usually much speedier & still acceptably sensitive.



Example lossy algorithm: pick $q > floor(n / (k + 1))$