

# String Algorithms and Data Structures

## Introduction and Pattern Matching

CS 199-225

January 26, 2026

Brad Solomon



UNIVERSITY OF  
**ILLINOIS**  
URBANA - CHAMPAIGN

Department of Computer Science

# Who am I?



## **Brad Solomon**

Teaching Assistant Professor, Computer Science

2233 Siebel Center for Computer Science

Email: [bradsol@illinois.edu](mailto:bradsol@illinois.edu)

## **Office Hours:**

Thursdays, 11:00 AM - 12:00 PM

... or by appointment

<https://courses.engr.illinois.edu/cs225/info/office-hours/>

# Who are you?

Take a moment to introduce yourself to your neighbor!

(Your name, a hobby you enjoy, and one thing you hope to get out of this class)

Piazza Sign up Link: <https://piazza.com/illinois/spring2026/cs199225>





# What will you get out of this class?

Understand fundamental string algorithms

Experience applying data structures, algorithms, and algorithm design principles to real world problems

Justify implementation choices based on theoretical or practical considerations

Build a foundation for future data science projects

6-8  
weeks  
on  
exact  
pattern  
matching

# Course Webpage



<https://courses.grainger.illinois.edu/cs225/sp2026/pages/honors.html>

All course information and links can be found here!

Mediaspace recordings

Piazza

Syllabus

# Syllabus

Please read — many important topics:

Course Goals & Topics

Course Expectations

Grading

Commitments to Diversity, Equity, Inclusion

Commitments to Mental Health

Ethics and Academic Integrity Policies

# Course Expectations

Weekly assignments (11 total):

Small assignments (~ 1-3 hours / week)

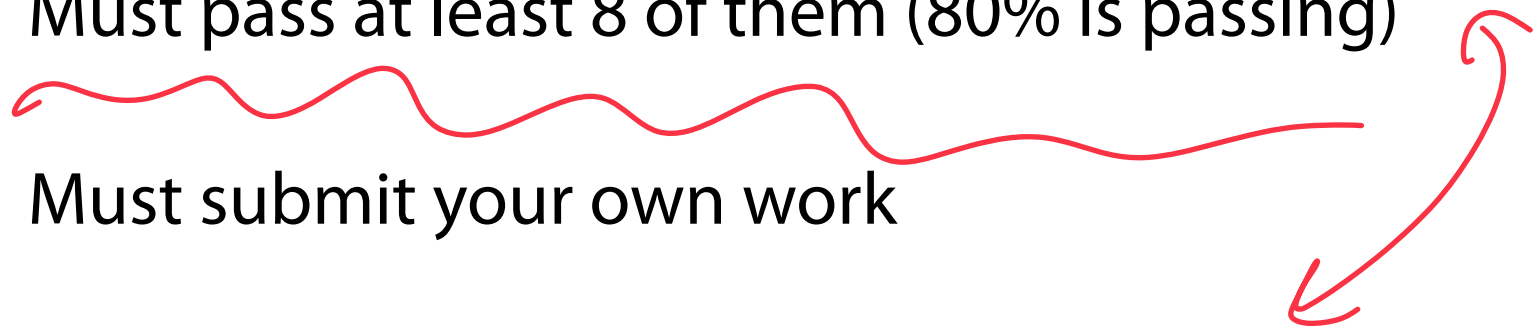
Must pass at least 8 of them (80% is passing)

Must submit your own work

One week extensions for 80% credit

→ ○ credit

2-3 functions



# Course Expectations

Class participation:

No attendance grades

Ask questions (synchronously or asynchronously)

Participate in breakout rooms and polls

# Mental Health

This class should be low-stress, light work-load.

UIUC offers a variety of confidential services:

**Counseling Center:** 217-333-3704

610 East John Street Champaign, IL 61820

**McKinley Health Center:** 217-333-2700

1109 South Lincoln Avenue, Urbana, Illinois 61801

# Diversity, Equity, and Inclusion



“If you witness or experience racism, discrimination, micro-aggressions, or other offensive behavior, you are encouraged to bring this to the attention of...”

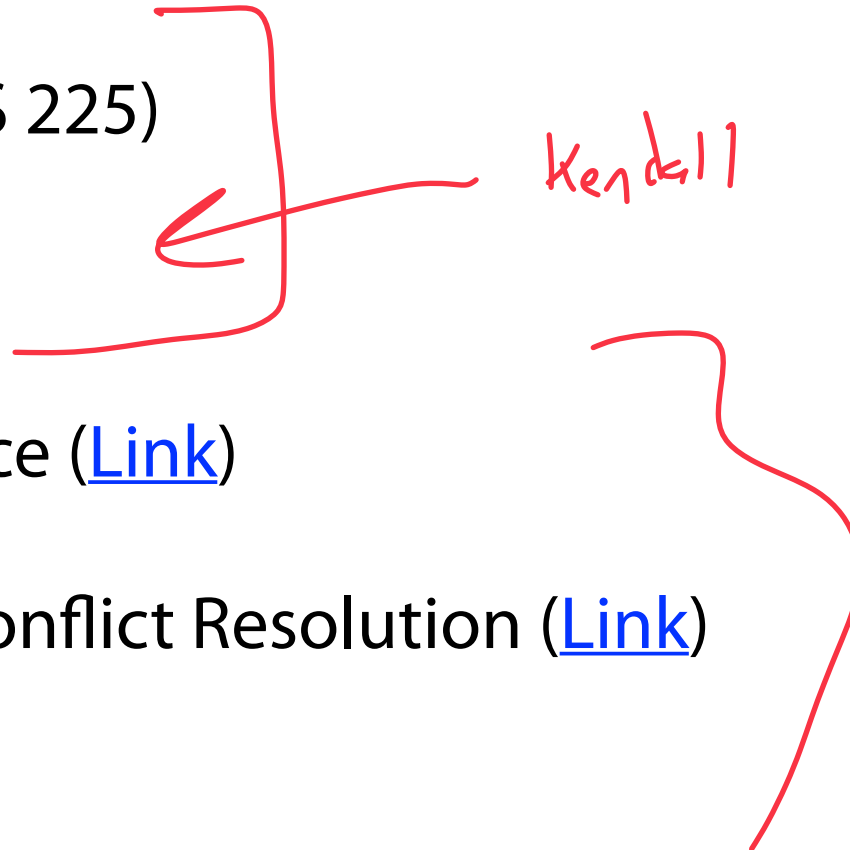
Staff (CAs and TAs for CS 225)

Faculty (Brad Solomon)

Campus Belonging Office ([Link](#))

The Office of Student Conflict Resolution ([Link](#))

CS CARES ([Link](#))



# Learning Objectives

Review fundamentals of strings



Introduce exact pattern matching problem

---

↳ Discussions

# What is a string?

*String*  $S$  is a finite sequence of characters

Characters are drawn from alphabet  $\Sigma$ , usually assumed finite

Nucleic acid alphabet: { A, C, G, T }

English: { A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, Q, R, S, T, U, V, W, X, Y, Z }

What are some other alphabets we could use?

↳ Numeric!

↳ Unicode

↳ ASCII 256 512

↳ Binary Alphabet  $\Leftrightarrow$  hex or any other base values

# What is a string... in C++?

↳ style

**char:** 1-byte (8-bit) character encoding [ASCII 256]

Do this

↳ **std::string:** uses *char* alphabet (by default), has significant operation support

string\_main.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5
6     char c[] = "Hello World";
7
8     std::string str = "Hello World";
9
10    return 0;
11 }
```

# Fundamental operations

## Math

+

-

/

%

\*

## Strings

Append

concatenation

Substring

length

Removal

Find

duplicate

split

Equals

Big list!

Lets discuss a subset

# Fundamental string operations

Big O

“How efficient is my algorithm at searching for a given pattern  $P$ ?”

length of size of a text!

“How much memory do I need to allocate for this text file?”

# Fundamental string operations

**Size** of  $S$ ,  $|S|$ : The number of characters in  $S$ .

$S = \text{"How big?"}$

$|S| = 7$

8

9



Join Code: 225

# Fundamental string operations

**Size** of  $S$ ,  $|S|$ : The number of characters in  $S$ .

$S = \text{"How big?"}$

$|S| = 8$

*9 including null*  
*not we will use C++ convention*

0	1	2	3	4	5	6	7
H	o	w	_	b	i	g	?

*↑  
space*

*↑  
?*

# Fundamental string operations

**Size** of  $S$ ,  $|S|$ : The length of  $S$  (in terms of bytes).

$S.length()$

size.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     std::string S = "Is this a string?";
6     std::string T = "No, this is Patrick.";
7
8     std::cout << S.length() << std::endl;    17
9     std::cout << T.length() << std::endl;    20
10
11     return 0;
12 }
13
14
```

# Fundamental string operations

**“Is this book about data structures?”**

*equals / compare*

*find / contain*

**“Is this student enrolled at UIUC?”**

# Fundamental string operations

*S equals* T if each character, in order, is the same

S == T

equals.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     std::string S = "Thing 1";
6     std::string T = "Thing 1";
7
8     if (S == T) {
9         std::cout << "S == T" << std::endl;
10    } else {
11        std::cout << "S != T" << std::endl;
12    }
13    return 0;
14 }
```

# Fundamental string operations



Join Code: 225

*S equals* T if each character, in order, is the same

`S == T`

char\_equals.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     char S[] = "Thing 1";
6     char T[] = "Thing 1";
7
8     if (S == T){
9         std::cout << "S == T" << std::endl;
10    } else {
11        std::cout << "S != T" << std::endl;
12    }
13    return 0;
14 }
```

*(A) 23%*  
*(B) 50%*  
*(C) It crashes 25%*



# Fundamental string operations

*S equals* T if each character, in order, is the same

S == T

char\_equals.cpp

```
1 substring.cpp:8:9: warning: array comparison always evaluates to false [-Wtautological-compare]
2   if (S == T){
3       ^
4   int main() {
5       char S[] = "Thing 1";
6       char T[] = "Thing 1";
7
8       if (S == T){
9           std::cout << "S == T" << std::endl;
10      } else {
11          std::cout << "S != T" << std::endl;
12      }
13      return 0;
14  }
```

*Handwritten notes:*  
A red arrow points from the warning message to the `S == T` comparison in the code.  
To the right of the code, the text `==` is written above `|||`, which is above `.equals`.

# Fundamental string operations

Substring (find)  
concatenation

## Reads

GTATGCACGCGATAG	TATGTCGCAGTATCT	CACCCTATGTCGCAG	GAGACGCTGGAGCCG
TAGCATTGCGAGACG	GGTATGCACGCGATA	TGGAGCCGGAGCACC	CGCTGGAGCCGGAGC
TGTCTTTGATTCTG	CGCGATAGCATTGCG	GCATTGCGAGACGCT	CCTATGTCGCAGTAT
GACGCTGGAGCCGGA	GCACCCTATGTCGCA	GTATCTGTCTTTGAT	CCTCATCCTATTATT
TATCGCACCTACGTT	CAATATTCGATCATG	GATCACAGGTCTATC	ACCCTATTAACCACT
CACGGGAGCTCTCCA	TGCATTTGGTATTTT	CGTCTGGGGGGTATG	CACGCGATAGCATTG
GTATGCACGCGATAG	ACCTACGTTCAATAT	TATTTATCGCACCTA	CCACTCACGGGAGCT
GCGAGACGCTGGAGC	CTATCACCTATTAA	CTGTCTTTGATTCT	ACTCACGGGAGCTCT
CCTACGTTCAATATT	GCACCTACGTTCAAT	GTCTGGGGGGTATGC	AGCCGGAGCACCTA
GACGCTGGAGCCGGA	GCACCCTATGTCGCA	GTATCTGTCTTTGAT	CCTCATCCTATTATT
TATCGCACCTACGTT	CAATATTCGATCATG	GATCACAGGTCTATC	ACCCTATTAACCACT
CACGGGAGCTCTCCA	TGCATTTGGTATTTT	CGTCTGGGGGGTATG	CACGCGATAGCATTG

## Genome

↓ Build one overlap string

CGTCTGGGGGGTATGCACGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCTATGTCGCAGTATCTGTCTTTGATTCTG

↑ concatenation but how do we know what to concat?

# Fundamental string operations

**Concatenation** of  $S$  and  $T$ : characters of  $S$  followed by characters of  $T$

$S = \text{"Beep"}$        $T = \text{"Boop"}$

What is the string  $ST$ ?

Beep Boop

What is the string  $T^f S$ ?

Boop \$ Beep

# Fundamental string operations

**Concatenation** of  $S$  and  $T$ : characters of  $S$  followed by characters of  $T$

$S + T$

concat.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     std::string S = "Beep";
6     std::string T = "Boop";
7
8     std::cout << S + T << std::endl;
9     std::cout << T + S << std::endl;
10
11     std::cout << S + '$' + T << std::endl;
12     std::cout << T + '$' + S << std::endl;
13 }
14
```

# Fundamental string operations

**“Is this book about data structures?”**

find()  
substring()

**S: Data Structures**

## 1.1 Why Compact Data Structures?

**T:**

Google’s stated mission, “to organize the world’s information and make it universally accessible and useful,” could not better capture the immense ambition of modern society for gathering all kinds of data and putting them to use to improve our lives. We are collecting not only huge amounts of data from the physical world (astronomical, climatological, geographical, biological), but also human-generated data (voice, pictures, music, video, books, news, Web contents, emails, blogs, tweets) and society-based behavioral data (markets, shopping, traffic, clicks, Web navigation, likes, friendship

# Fundamental string operations

$S$  is a **substring** of  $T$  if there exists (possibly empty) strings  $u$  and  $v$  such that  $T = uSv$

A **substring** is a sequence of characters (a string) contained within another string

**$S$ :** p e p p e r

**$T$ :** I \_ l i k e \_ p e p p e r o n i \_ p i z z a

# Fundamental string operations

A **substring** of  $S$  is a string occurring inside  $S$

$S$ .**substr**(size\_t <sup>start ↓</sup>pos, size\_t len ↓)

substring.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     std::string T = "Hello my name is ";
6
7     std::cout << T.substr(1,4) << std::endl;
8
9     return 0;
10 }
```

ello

# Fundamental string operations

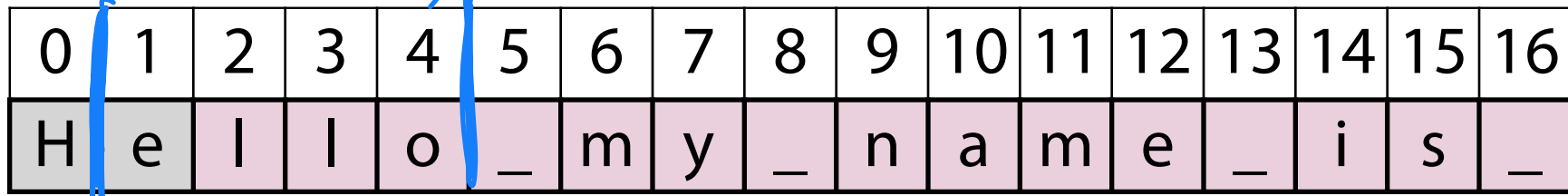


A **substring** of  $S$  is a string occurring inside  $S$

$S$ .**substr**(size\_t pos, size\_t len)

substring.cpp

```
1 #include <string>
2 #include <iostream>
3
4 int main() {
5     std::string T = "Hello my name is ";
6
7     std::cout << T.substr(1,4) << std::endl;
8
9     return 0;
10 }
```



# Fundamental string operations

$S$  is a **prefix** of  $T$  if there exists a string  $v$  such that  $T = Sv$

A **prefix** is a substring  $T = uSv$  where  $u = ""$

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

G T T A T A G C T G A T

S

V

Some substring that starts at position 0

# Fundamental string operations

$S$  is a **prefix** of  $T$  if there exists a string  $v$  such that  $T = Sv$

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

# Fundamental string operations

$S$  is a **prefix** of  $T$  if there exists a string  $v$  such that  $T = Sv$

$v = ""$

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

G T T A T A G C T G A T

G T T A T A G C T G A

G T T A T A G C T G

G T T A T A G C T

G T T A T A G C

G T T A T A G

G T T A T A

G T T A T

G T T A

G T T

G T

G

← prefix can be full length



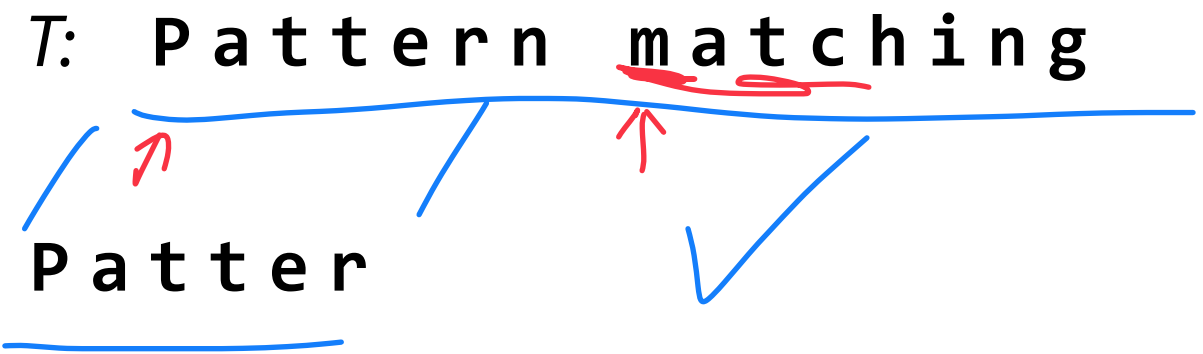
# Fundamental string operations



Join Code: 225

$S$  is a **prefix** of  $T$  if there exists a string  $v$  such that  $T = Sv$

|| 0 1 2 3 4 5 6 7 8



matching 20%

greatest common prefix

Patrick

# Fundamental string operations

$S$  is a **prefix** of  $T$  if there exists a string  $v$  such that  $T = Sv$

$T$ : P a t t e r n m a t c h i n g

P a t t e r



m a t c h i n g



P a t r i c k



# Fundamental string operations

$S$  is a **suffix** of  $T$  if there exists a string  $u$  such that  $T = uS$  

A **suffix** is a substring  $T = uSv$  where  $v = ""$  

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

G T T A T A G C T G A T

$u$

$S$

$\uparrow$

can be ""

# Fundamental string operations

$S$  is a ***suffix*** of  $T$  if there exists a string  $u$  such that  $T = uS$

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

# Fundamental string operations

$S$  is a **suffix** of  $T$  if there exists a string  $u$  such that  $T = uS$

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

G T T A T A G C T G A T 

T T A T A G C T G A T

T A T A G C T G A T

A T A G C T G A T

T A G C T G A T

A G C T G A T

G C T G A T

C T G A T

T G A T

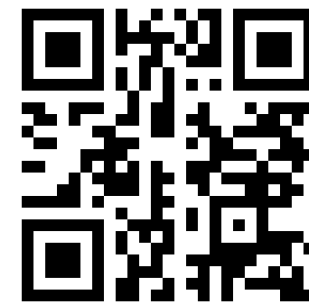
G A T

A T

T

# Fundamental string operations

$S$  is a **suffix** of  $T$  if there exists a string  $u$  such that  $T = uS$



Join Code: 225

$T$ : Pattern matching

ing

tern

ring

# Fundamental string operations

$S$  is a **suffix** of  $T$  if there exists a string  $u$  such that  $T = uS$

$T$ : Pattern matching

ing



tern



ring



# Fundamental string operations



Size,  $|S|$

`S.length()`

Equals,  $S == T$

`S == T`

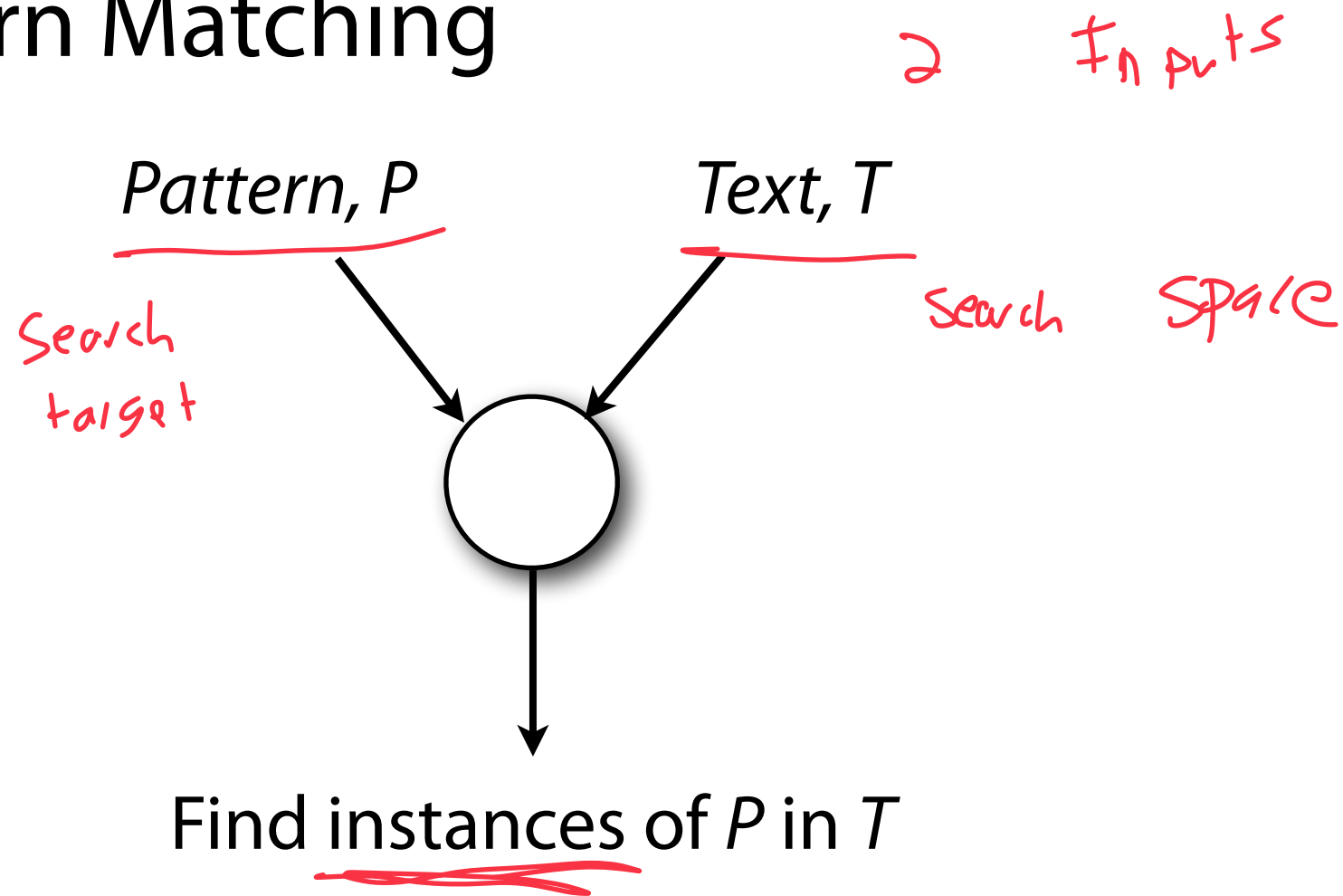
Concatenation,  $ST$

`S + T`

Substring,  $uSv$

`S.substr(pos, len)`

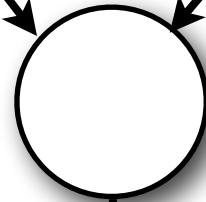
# Exact Pattern Matching



# Exact Pattern Matching

*Pattern, P*

*Text, T*



Find instances of  $P$  in  $T$

'instances': An exact, full length copy

# Exact Pattern Matching

Find places where *pattern*  $P$  occurs as a substring of *text*  $T$ . Each such place is an occurrence or match.

$P$ : word

$T$ : There would have been a time for such a word

# Exact Pattern Matching

Find places where *pattern*  $P$  occurs as a substring of *text*  $T$ . Each such place is an *occurrence* or *match*.

$P$ : word

$T$ : There would have been a time for such a word  
                                  word                                  word

# Exact Pattern Matching

Find places where *pattern*  $P$  occurs as a substring of *text*  $T$ . Each such place is an *occurrence* or *match*.

$P$ : word

$T$ : There would have been a time for such a word

Alignment 1: word 

Alignment 2: word 

**Alignment:** a way of putting  $P$ 's characters opposite  $T$ 's. May or may not correspond to a match.

# Exact Pattern Matching

Find places where *pattern*  $P$  occurs as a substring of *text*  $T$ . Each such place is an *occurrence* or *match*.

$P$ : word

$T$ : There would have been a time for such a word

Alignment 1: word

Alignment 2: word

Not a match!

Match!

**Alignment:** a way of putting  $P$ 's characters opposite  $T$ 's. May or may not correspond to a match.

# Exact Pattern Matching

What's a simple algorithm for exact matching?

P: word

T: There would have been a time for such a word

for each alignment in text

↳ check if pattern equals text's substring

# Exact Pattern Matching

What's a simple algorithm for exact matching?

P: word

T: There would have been a time for such a word

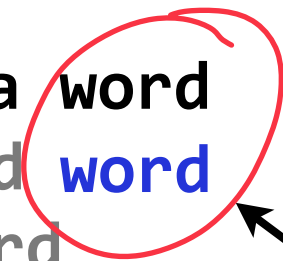
word word word word word word word word word word

word word word word word word word word word

word word word word word word word word word

word word word word word word word word word

word word word word word word word word word



One occurrence

Try all possible alignments. For each, check if it matches. This is the *naïve algorithm*.



# Assignment 1: a\_naive

Learning Objective:


Conceptualize exact pattern matching w/ naïve search



Demonstrate understanding of fundamental operations

Think about as you code: is naïve search a good solution?

Prefix  
&  
Suffix



# End-of-class brainstorm

How can we improve the naïve algorithm?

# End-of-class brainstorm

How can we improve the naïve algorithm?

... if you have infinite space?

$O(1)$   
↑

could make tree as faster encoding of text

↳ use hashing

create a look up table for every string in the text & its location

↳ preprocess data to speed up search

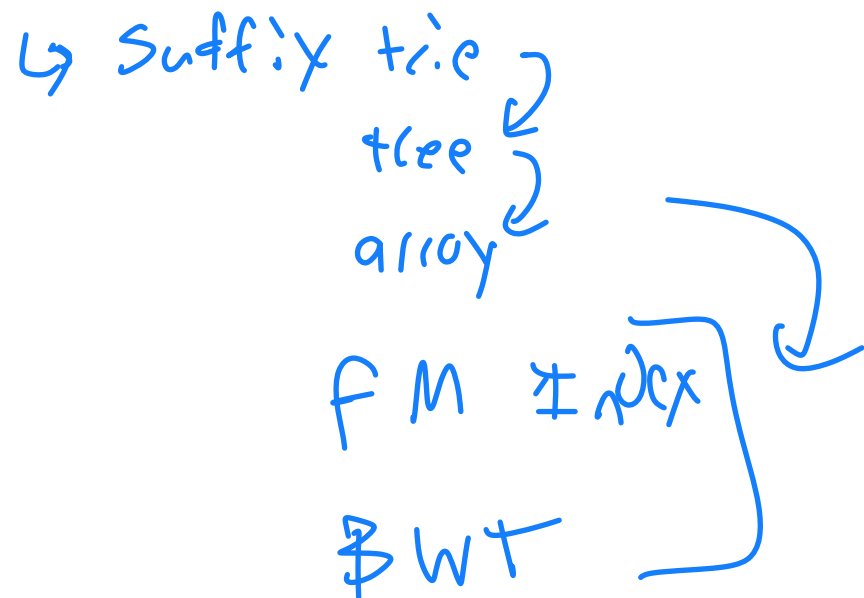


# End-of-class brainstorm

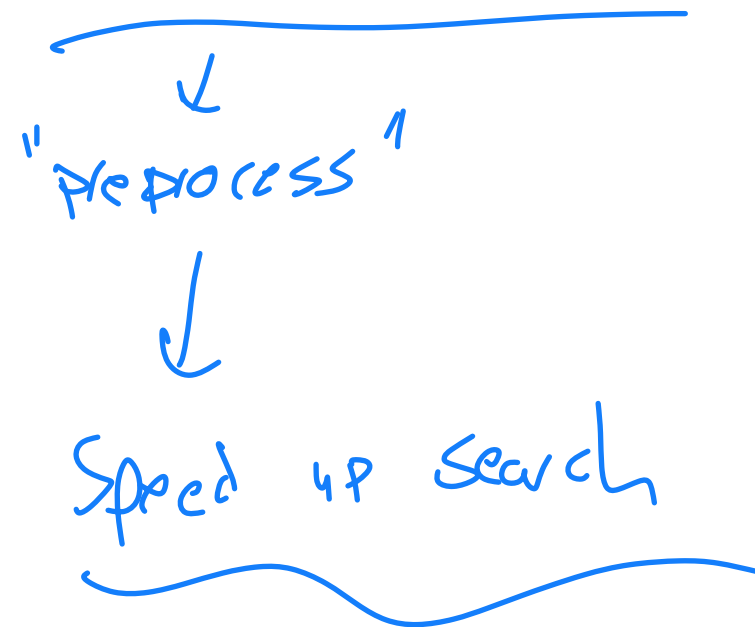
How can we improve the naïve algorithm?

... if I tell you the text ahead of time?

↳ Suffix trie  
tree  
array  
FM Index  
BWT



Preprocess  
↓  
Speed up search



# End-of-class brainstorm

How can we improve the naïve algorithm?

... if you have infinite space?

... if I tell you the pattern ahead of time?

... if I tell you the text ahead of time?