CS2040S Data Structures and Algorithms

Welcome!

Algorithms Everywhere

· III

Web browser:
Parsing
Substring manipulation (Week 7)

XML trees (Week 5)

Internet

Internet routing:
TCP (congestion control)
IP routing
BGP (Bellman-Ford, Week 9)
Content caching (Week 7)
DNS

Google:

PageRank (Week 10) String matching (Week 7)

Web servers:

Load balancing (PS 2) Scheduling (Week 8) Memory allocation



<u>Database:</u>

B-trees (Week 5) Search (Week 2) Sorting (Week 3)

Data Structures

Problem Solving

Algorithms

What is an *algorithm*?

- Set of instructions for solving a problem
 "First, wash the tomatoes."
 "Second, peel and cut the carrots."
 "Third, mix the olive oil and vinegar."
 "Finally, combine everything in a bowl."
- Finite sequence of steps
- Unambiguous (and precise)
- Designed to accomplish some task

History

- Named for al-Khwārizmī (780-850)
 - Persian mathematician

Many ancient algorithms



History

- Named for al-Khwārizmī (780-850)
 - Persian mathematician

What is the oldest known algorithm?



History

- Named for al-Khwārizmī (780-850)
 - Persian mathematician

- Many ancient algorithms
 - Multiplication: Rhind Papyrus
 - − Babylon and Egypt: ~1800BC
 - Euclidean Algorithm: Elements
 - Greece: ~300BC
 - Sieve of Eratosthenes
 - − Greece: ~200BC



WHAT ARE YOUR GOALS?



BUILD USEFUL SYSTEMS FOR PEOPLE



Jeff Dean (Google)



Ruchi Sanghvi (Facebook)

BUILD FAST SYSTEMS

"If you need your software to run twice as fast, hire better programmers.

But if you need your software to run more than twice as fast, use a better algorithm."

-- Software Lead at Microsoft

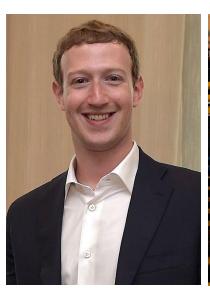
"Bad programmers worry about the code. Good programmers worry about data structures and their relationships."

- Linus Torvalds*



*Creator of git and the linux kernel (in Linux OS, Chrome OS, Android)

START A TECH COMPANY





THINK ABOUT BEAUTIFUL ALGORITHMS



Donald Knuth

"People who analyze algorithms have double happiness. First of all they experience the sheer beauty of elegant mathematical patterns that surround elegant computational procedures. Then they receive a practical payoff when their theories make it possible to get other jobs done more quickly and more economically."

CHANGE THE WORLD

COMPUTER SCIENTISTS HAVE CHANGED THE WORLD...

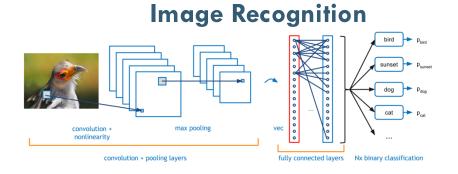


...BY SOLVING PROBLEMS

CS2040S is about solving problems.

It prepares you for what lies ahead...





WHAT ABOUT THE CODE?

WHAT ABOUT THE CODE?

CS2040S is about solving problems in Java.

Why Java?

Language does no

Good aspects:

- Common in industry / real-world / web
- Object-oriented in a deep way
- Modularity / abstraction via OOP
- Avoids memory leak issues of C/C++

Less good aspects:

- Performance (compare to: C++)
- More mental overhead (compare to: Python)

WHAT ABOUT THE CODE?

CS2040S is about solving problems in Java.

CS2040S is not about learning Java.

"If you need your software to run twice as fast, hire better programmers.

But if you need your software to run more than twice as fast, use a better algorithm."

-- Software Lead at Microsoft

CS2030S

"If you need your software to run twice as fast, hire better programmers.

But if you need your software to run more than twice as fast, use a better algorithm."

-- Software Lead at Microsoft

CS2040S

Assumption:

• You are taking CS2030S this semester.

OR

• You already took CS2030S a previous semester.

OR

You already are an expert Java programmer.

OR

You plan to spend some extra time learning Java.

WHAT ABOUT THE CODE?

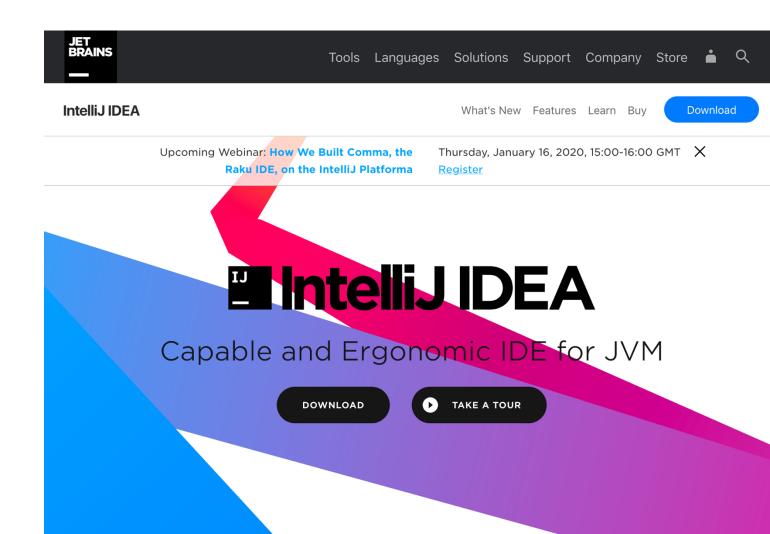
CS2040S is about solving problems in Java.

CS2040S is not about learning Java.

We want you to be great programmers too! Staff are here to help.

Always program using the best development environment you can!

CS2040S Recommended IDE: IntelliJ



Introductions

Teaching Team



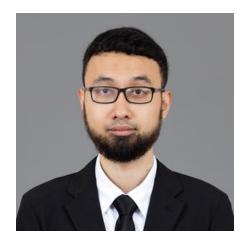
Diptarka Chakraborty



Eldon Chung



Ben Leong



Muhammad Rizki



Lin Shao

What should I call you?

- If I mispronounce your name, please correct me.
- If the name you use is different than the one I have, please tell me.
- If I am addressing you incorrectly, please let me know.

Teaching Staff

Tutors:

~60 tutors

NUS undergrads

Former CS2040S students

Teaching Assistants: (Full Time TAs)

Eldric Liew

Jason Ciu Putra

Enzio Kam

WEEKLY SCHEDULE

2x Lectures. Either:

- 1. Monday: 4pm-6pm (LT11), Wednesday 1pm-2pm (LT11)
- 2. Monday: 2pm-4pm (UT-AUD1), Wednesday 2pm-3pm (LT11)

1x Recitation (room: TBA)

- 18 slots (Thursday/Friday)
- Starts in Week 4

1x Tutorial (room: TBA)

- 14 slots (Tuesday/Wednesday), ~60 groups
- Starts in Week 4

Lectures:

- Two lectures. Either:
 - 1. Monday: 4pm-6pm (LT11), Wednesday 1pm-2pm (LT11)
 - 2. Monday: 2pm-4pm (UT-AUD1), Wednesday 2pm-3pm (LT11)
 - Lecturers: Diptarka Chakraborty, Eldon Chung

Whether live or online:

- Participation expected
- Recorded
- Slides posted after lecture

FAQ: Will you post lectures slides *before* class?

No, only after class.

Three reasons:

- 1. We are preparing the slides before class.
- 2. Reading slides is not a good way to prepare for class. (Read the recommended material instead, if you have time.)
- 3. During class, you will get distracted if you are trying to flip through a separate slide deck. It will NOT help you to catch up.

FAQ: Will you po *before* class? A Respected MIT Professor Had a Simple 4-Word Rule for His Inc. Classroom, and Every Company Should Follow It Patrick Winston knew how to speak. But his even greater skill was getting others to class. BY JUSTIN BARISO, AUTHOR, EQ APPLIED @JUSTINJBARISO listen. 8 to prepare for material instead, if

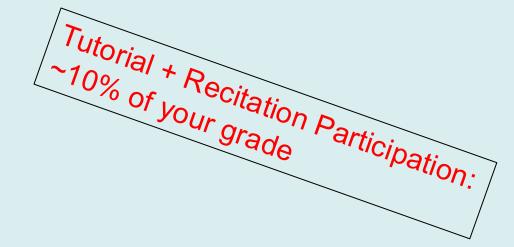
3. Lamg class, you will get distracted if you are trying to flip through a separate slide deck. It will NOT help you to catch up.

Recitations: Thursday/Friday (1 hour)

- Size: ~35 students
- Goal: algorithm design and problem solving

Instructors:

- Lin Shao
- Ben Leong
- Muhammad Rizki
- Eldric Liew
- Enzio Kam



Recitations: Thursday/Friday (1 hour)

- Size: ~35 students
- Goal: algorithm design and problem solving

Register:

- Survey on Coursemology
- Deadline: TONIGHT Monday 11:59pm
- NO ModReg.
- Sessions start Week 4.

Administrative Details

Tutorials:

One tutorial: Tuesday/Wednesday (2 hours)

60 different sessions

Size: ~15 students

Your tutor's job:

- Help you learn the material.
- Tutorial + Recitation Participation: Provide advice on how best to learn.
- Help you catch up if you fall behind.
- Challenge you if you are ahead.
- Show you neat exciting aspects of CS2040S.

Administrative Details

Tutorials:

One tutorial: Tuesday/Wednesday (2 hours)

60 different sessions

Size: ~15 students

How to sign up?

- Survey on Coursemology
- Deadline: TONIGHT Monday 11:59pm
- Please be flexible: fill in as many sessions as you can.
- Please be patient.
- Sessions start on Week 4.

WEEKLY SCHEDULE

2x Lectures. Either:

- 1. Monday: 4pm-6pm (LT11), Wednesd
- 2. Monday: 2pm-4pm (UT-AUD1), We

1x Recitation (room: TBA)

- 18 slots (Thursday/Friday)
- Starts in Week 4

1x Tutorial (room: TBA)

- 14 slots (Tuesday/Wednesday), ~60 groups
- Starts in Week 4

Reminder:

Do not sign for tutorials & recitations on ModReg.

Fill out the Coursemology

survey today.

Administrative Details

Midterm:

Date: Week 7 Monday (3 March, 2025)

Time: 6:30PM to 9PM

Venue: MPSH 2A & 2B

I prefer to have the midterm during regular class time to avoid conflicts, but room schedule for 800+ is difficult!

Administrative Details

Final Exam:

Date: 26 April, 2025

Time: 1PM to 3PM

Venue: TBA

Important note: This year, papers will not be fully MCQ papers.

DROP DATES

Check CourseReg: http://www.nus.edu.sg/CourseReg/

My current impression (BUT PLEASE CHECK):

Drop with:

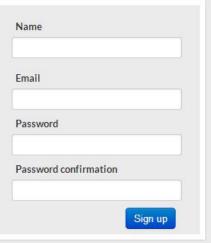
- No record: until 26/01/2024 11:59pm.
- 'W' grade: 27/01/2025 00:00 through 2/03/2025 11:59pm.
- 'F' grade: 3/03/2024 00:00 onwards.



Gamified Online Education Platform:

Making your class a world of games in a universe of fun.







Engaging

Coursemology allows educators to add gamification elements, such as experience points, levels, achievements, to their classroom exercises and assignments.

The gamification elements of Coursemology motivate students to do assignments and trainings.



General

It's built for all subjects. The gamification system of Coursemology doesn't make assumption on the course's subject.

Through Coursemology, any teacher who teaches any subject can turn his course excercies into a online game.



Simple

It's built for all teachers. You don't need to have any programming knowledge to master the platform.

Coursemology is easy and intuitive to use for both teachers and students.

Coursemology

Built here at NUS by Prof. Ben and his team!

Platform for course management:

- Announcements
- Lecture slides
- Problem sets
- Feedback
- Discussion forum
- Etc.

Please check Coursemology, and read your (Coursemology) e-mail.

Replacement for Canvas.

Coursemology

How do you join?

- If you are registered for the class, you have already received an invitation.
- If you are not registered for the class, you will receive an invitation when you do register.

If something doesn't work.... ask on forum!

Canvas

Reasons to check Canvas?

- Almost none...
- Final grade confirmation at the end of the semester...

Least interesting part of CS2040S?

Grades

AND

Talking about grades

AND

Optimizing for grades

AND

Thinking about grades

Solution: Everything gives EXP

ParticipationAND

Assignments

AND

Surveys

AND

Optional exercises

EXP → Levels

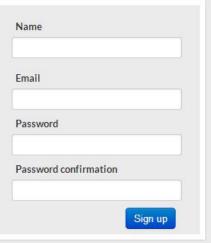
- final level = grade for classwork
- Maximum grade: 35
- Maximum level: > 35
- Bonus / optional work may lead to levels higher than
 35.
- Goal: everyone reaches level 35 (or close)



Gamified Online Education Platform:

Making your class a world of games in a universe of fun.







Engaging

Coursemology allows educators to add gamification elements, such as experience points, levels, achievements, to their classroom exercises and assignments.

The gamification elements of Coursemology motivate students to do assignments and trainings.



General

It's built for all subjects. The gamification system of Coursemology doesn't make assumption on the course's subject.

Through Coursemology, any teacher who teaches any subject can turn his course excercies into a online game.



Simple

It's built for all teachers. You don't need to have any programming knowledge to master the platform.

Coursemology is easy and intuitive to use for both teachers and students.

How to be good at anything CS2040S?

Knowledge

Experience

Talent

Practice, practice, practice...

Lecture Review

Quick review of topics from class

< 15 minutes

Monday/Wednesday:

Due within 24 hours of end of lecture.

(50% reduced EXP after deadline)

Problem Sets	Optional Practice				
	Experience Points	Bonus Experience Points	Requirement For	Start At	Bonus Cut Off
Lecture Training 1: Introduction		125		12 Jan 18:00	15 Jan 23:59
Lecture Training 2: Java OOP		125		14 Jan 17:00	15 Jan 23:59
		Experience Points	Experience Points Experience Points troduction 125 125	Experience Points Experience Points For Requirement For	Experience Points Experience Points Requirement For Start At Start At 125 125 12 12 Jan 18:00

Practice the techniques from class.

Find out what you do and do not understand.

Practice regularly, every week.

Practice the techniques from class.

Find out what you do and do not understand.

Practice regularly, every week.

Problem set 1 is available on Coursemology.... as soon as you complete the Lecture Review today.

Some problems may be hard.

Some problems may involve figuring something out that we haven't done in class!

Most problems will involve writing some Java code. Warning: A tutor is not a compiler.

Your tutor is there to help!

Submit problem sets on Coursemology.

Late submissions:

- 1 week: 25% penalty
- Last day of class: 50% penalty

Hand in problem sets on time!

Even if late, do them anyways!

Submit problem sets on Coursemology.

- Visible (public) testcases.
- Invisible (private) testcases.
- Grading-only testcases.
- Tutor's evaluation and judgement.

Tutors will give feedback on both correctness and how to write better code.

Practice Problems

Ensure that you really know how to use the algorithms & data structures.

Easier than problem sets

Good practice for coding skills (interviews, etc.)

Simple problem solving

Lecture Trainings	Problem Sets	Optional Practice		
Title		Experience Poin	ts Requirement For	Start At
Guess the Number (Binary search)		200	⊕ €	19 Jan 00:00
WiFi (Binary Search)	200	@	19 Jan 00:00

Recitations

Problem solving (Apply ideas from lecture.)

Problem solving techniques (How to go about solving hard problems.)

Missing details
(Fill in missing pieces from lecture.)

Tutorials

Type 1: Review Problems
(Builds directly on ideas from lectures.)

Type 2: Group Challenge Problems (Work together to solve a hard problem.)

Who here plays chess?

Rybka



Rybka (Computer) - Shredder (Computer) 1-0 C67 WCCC Pampiona Openciass (6) 13,05,2009

1.e4 e5 2.包f3 包c6 3.息b5 包f6 4.0-0 包xe4 5.曾e2 包g5 6.包xg5 曾xg5 7.d4 曾e7 8.dxe5 包d4 9.曾d3 曾xe5 10.包c3 息c5 11.曾d1 包e6 12.亘e1 曾d4 13.曾f3 0-0 14.亘e4 曾d6 15.亘h4 曾e5 16.息d2 f5 17.亘e1 曾f6 18.曾h3 曾g6 19.包d5 c6 20.亘xe6 曾xe6 21.包f4 曾xa2 22.亘xh7 cxb5 23.g3 亘f6 24.息c3 查f7 25.曾h4 曾a1+ 26.查g2 曾a6 27.夏xf6 曾xf6 28.粤h5+

1-0

Download PGN

Won four consecutive World Chess Championships (2007-2011).

Today's best chess engines:

- AlphaZero
- Stockfish
- Komodo
- ...

Computer chess reels from 'biggest sporting scandal since Ben Johnson'

Czech mate, Mr Cheat

IT'S a very that has seen puses, and review questing cell aftern bounds across the world.

Rybia, the best aftersoplaying computer on the planet, is a cheat. And stockerologier, Varia Rajhat – one half of a couple dubbed the Posts and Becks of the game – has been sharted as a plaguated, banned from computing, supposed.

of his titles and ordered to hand-

back his trophics and prize money.

Righth, himself an international master of the gatte, was footed guilty by his piece of harically copying autlor afters programs when creating Rytha.

A 34-member panel found the 4th year-old Creek from Ohio, but now living in Hungary, plaguarised two other programs. Crafty and Finis. Their report states. Not a single panel member betweed him innocent. Vasil. Rajlich's claims of complete, originally are contrary to the facts."

Not unce IBM s Doep Blue computer defeated grand master Garry Karpanis in 1997 – and was subsequently accound of cheating – has the world of computer chess been in such upsum. Byblia was been in such upsum. Byblia was

By Tarig Table

the International Computer Games Association world afterprombing from 2000 to 2010.

Peter Deggers, from unfine site. Chew Vibes, said. "The impact in the computer chess world must be comparable to arguably the most famuse example of doping in artilative."—the positive drug testing of Canadian systems. Hen Sobrison."

For his pair, Rapheth has not conmented, were at exteal soft to the association in which he disposed Ryb-La rechalcol cook, written by others.

He never made grand master as a player so turned to programming. If figured there were about 2,000 people in the world storager than me in cheve, he core used, but not one chose player that was stronger than me in programming.

By 2005, Rythia - Crech for Tetle fish' - was ready. The cheft tener is his wife, Iwena, hexalf an international mater. David Levy, providing of the ECGA, said. "We are consisted the evidence against Rights in both overwhelming in its volunte and beyond reasonable conduction in the name."



Posh and checks: Vesik Rajiich and his chief tester, wife lasts

Computer chess reels from 'biggest sporting scandal since Ben Johnson'

Czech mate, Mr Cheat



Relied on code from Crafty and Fruit, other chess engines.

Question: Was it original?

Did not meet the "originality" rules for the tournament.

Arguments for:

Disqualified!

Titles revoked!

<u>Arguments against:</u>

It was much, much better than Crafty, Fruit, or any existing chess player!

Collaboration Policy

- Working together is <u>strongly</u> encouraged!
- List your collaborators / sources when you submit.
- You <u>must</u> write/code/debug your problems sets alone, by yourself---not pair programming!
- Cheating / plagiarism will be dealt with harshly.

(Do the expected valued calculation, taking into account that you do not need to ace all the problem sets to reach Level 35 on Coursemology.)

Should we use ChatGPT to write code?

Should we use ChatGPT to write code?

Arguments for:

- It's pretty good!
- It can save a lot of time.
- Within five years, it will be so good that no one ever needs to write code again.

Should we use ChatGPT to write code?

Arguments for:

- It's pretty good!
- It can save a lot of time.
- Within five years, it will be so good that no one ever needs to write code again.

Arguments against:

- It produces a lot of broken code.
- It produces code with errors that are hard to find.
- Need to learn to code to use ChatGPT well.
- Real goal: learning to think about problems like a computer scientist---requires doing it yourself!

My philosophy (for school):

Rule 1:

Only use ChatGPT to do tasks that you know how to do without ChatGPT.

If you know how to look something up in Wikipedia, can ask ChatGPT!

Rule 2:

Only use ChatGPT for class when it is explicitly allowed in class.

Rule 3:

If it is against the rules to ask your friend, it is against the rule to ask ChatGPT.

Your friend isn't allowed to debug your code for you.

In CS2040S:

Goal: to learn to write code

Do not ask ChatGPT questions about code.

Do not let ChatGPT anywhere NEAR your code.

Do not send your code to your friends.

Do not look at your friends' code.

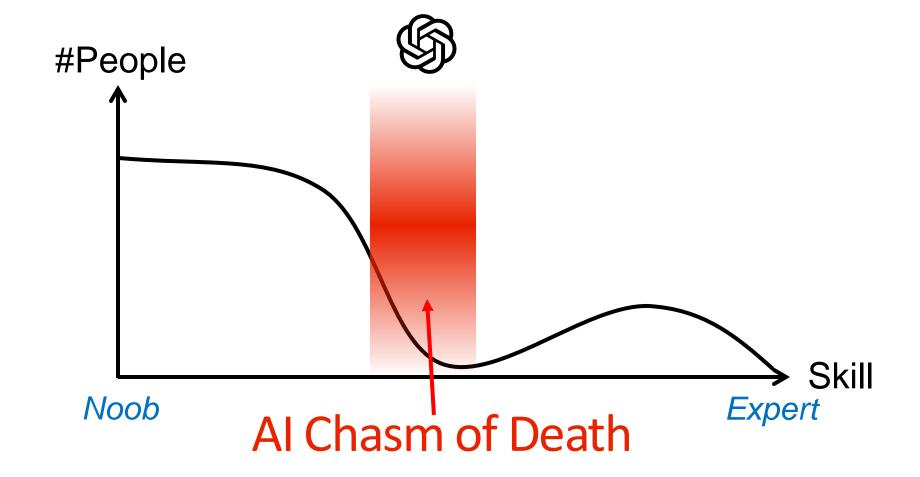
Goal: to learn to solve problems

Lookup ideas / concepts / data structures / algorithms.

Use Wikipedia, Coursera, Youtube, StackOverflow.

Work with your classmates, discuss solutions, ideas, algorithms.

Do not look at / throw away / discard all code : write your own solution.



Should our students use ChatGPT?

Use GPT/GenAl to do what you don't want to do but can do if required

Not what you cannot do

Anything that AI can do are easy problems

The real hard problems are people problems

Doing the thing right

Doing the right thing

WHAT ARE YOUR GOALS?

Learn something.

Become a better person.

- Position yourself to succeed after graduation.
- Get a good grade on a problem set in CS2040S??

WHY DID YOU CHEAT?

Computer chess reels from 'biggest sporting scandal since Ben Johnson'

Czech mate, Mr Cheat

IT's a very that has very turner and marks, sporting off versus and marks, sporting off versus and marks, sporting of the form of the form

dispared Crede from Other, bet 1 Spared better were about 2000 min per Informacy diagnostic programs. Cody and remembers, the row code programs. Cody and remembers, the rowe cand, that is concepted erepointly are contrast, and the fact of the complete erepointly are contrast of fish 2 ways ready. The cheek complete erepointly are contrast of fish 2 ways ready. The cheek contrast of fish 2 ways ready for the contrast of fish 2 ways ready. The cheek contrast of fish 2 ways ready for the contrast of the Cody, and 3 was contrasted the contrast of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the code of the Cody, and 3 was contrasted the cody of the cody



Czech mate, Mr Cheat

IUS a mory that has sere passes and stocks questing off obera-bounds across the world. Rybia, the best obest-playing computer on the planes, is a chear. And stockwoleper, Varol Raphats one half of a couple dalbed the. Chess Vibes, said. The impact is samed from comparing, impeed. In most example of depend in att-or, the trophics and price morey. Caradian springs flow between the Rajlech, henceff un resembled. For his part, Rajlech has not com-

ty by his peers of basecary ying darlier chess programs increasing Rybia.

34-member panel found the year-old Creek from Ohio, but

the International Computer Games Association, world alterpromiting from 2000 to 2010.

a control to the time of the control of the control

nothe facts?

Network(BM) Duep Blue complex defined grand maint facts?

Alegators in 1977 - and was who expected, account of the ICGA, said. "Be Alegators in 1979 - and was who expected, account of the ICGA, said. "Be Alegators in 1979 - and was who expected in the order of the ICGA and "Be Alekto in their own-theliting on an order and they are for the order whether and in solute and they and reasonable boom in such agout." If the ICGA was a second and they are for the ICGA was a recommendation of the ICGA was a recommenda



I thought I was going to fail the class!

If you put in the time, you will not fail the class.

WHY DID YOU CHEAT?

If you cheat, you will learn less, and likely do worse in the end.

Czech mate, Mr Cheat

the built of a couple dubbed the. Chess Vibes, said. The impact is

Not write ISM 1, Dong Dile corrition of delical grand mannet flaming
legions in 1997 – and was subrange of the ICM, and an econtrineed the evolution again
the world of computer choice
on in such agean. Kyhla su is solution and beyond mannet
open on such agean. Kyhla su is nature.



I'm not smart enough...

WHY DID YOU CHEAT?

Everyone else is smarter than I am...

You are smart enough; but you need to keep practicing, keep working at it.

Czech mate, Mr Cheat

one half of a couple dalbed the. Chess Vibes, said. The impact is

arned from competing, stripped if the tribs and ordered to hand ack his trophics and prize money. Raplich, herself un international



- I thought I was going to fail the class! You won't (unless you cheat).
- Everyone else is cheating! No, they aren't.

WHY DID YOU CHEAT?

Czech mate, Mr Cheat

IT'S a every that he and rooks spension bounds across the w

boards across the world.

Rybka, the best aftercuplaying computer on the planet, is a cheat.

And its skeeleper, Vanit Rajhati – one half of a comple dubbed the Posh and Becks of the game – has been sharined as a plagurist, beautiful processing the process of the planet.

based from computing, strapped of the titles and ordered to hand back bit recipites and prior morey. Righth, herself an recruational staster of the game, was found gutly by him poers of hoseably copying satter where programs which creating Rights. A 34-morber panel found the

A 34-monther parel found the filty-que full Creat foun Office, but now living in Hungary, plaguarised tous other programs. Creaty and Finit. Their appart states. Not a single parel monther believed him incorest. Vesil. Rajacht-Valians of complete originality are contrary to the fairs.

to the facts?

Not once IBM's Deep Bloe or point defeated grand matter that Kasparov in 1997 - and was weaparely account of chating his the world of computer of boos in such upone. Rybka w

By Tariq Tahir

from 2000 to 2016. Peter Doggors, from ordine site Check Whee, said. The impact of the computer check model must be comparable to agreedy the most famuse example of doping in athbrace – the prostore drug testing of Canadian systems Hers Solmoon?

I mental, see at entail sent is the a sociation in which the dispited Ryl is a reclassificable within by others. He never made grand mester as player on turned to programming. If figured their waits about 2,00 people in the world storage the mean clean, he cover used, but is see a clean, he cover used, but is see a clean, he cover used.

By 2005, Rythia - Creck for Ittle field - was ready. The chiteeper is he wide, Iwas, Iwasil a stirrrational mater. David Levpresident of the ICCIA, said. "Baare commend the evidence again Eighth is both overwhelming in as soluting and beyond transmall conduction in a storie."



Posh and checks: Vasik Rajich and his chief tester,

- I thought I was going to fail the class! You won't (unless you cheat).
- Everyone else is cheating! No, they aren't.

WHY DID YOU CHEAT?

It wasn't really cheating... The rules are the rules.

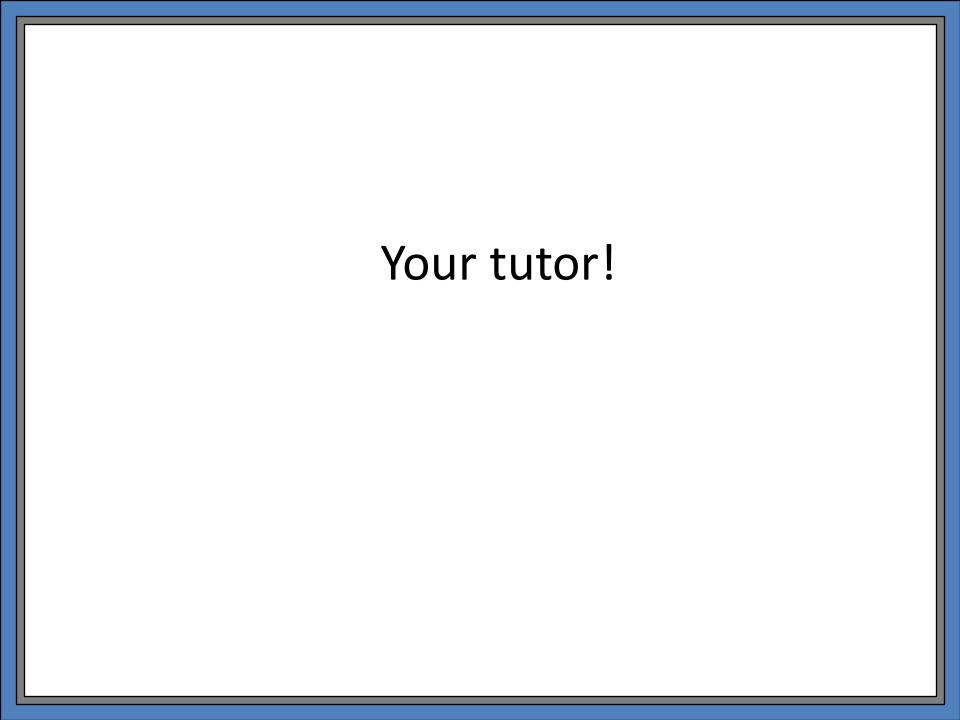
WHAT ARE YOUR GOALS?

Learn something.

Become a better person.

Position yourself to succeed after graduation.

Where can I get help?



Discussion Forums

Coursemology forums:

- Lectures
- Problem Sets
- Java
- Recitations & Tutorials
- and more....

Ask questions on Coursemology

(Don't e-mail me unless it is private.)

With 750+ students, I can't respond immediately to all emails (and I apologize if sometimes my reply is terse!), but someone on the teaching staff is always watching the forums.

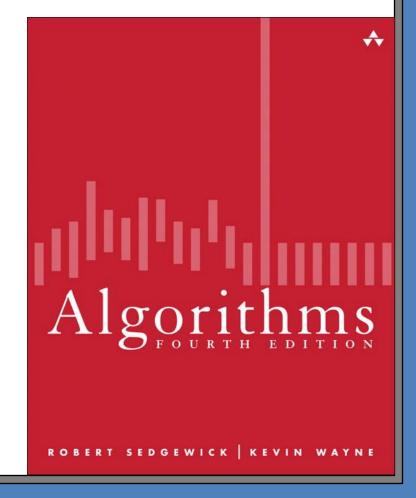
Other questions

Options:

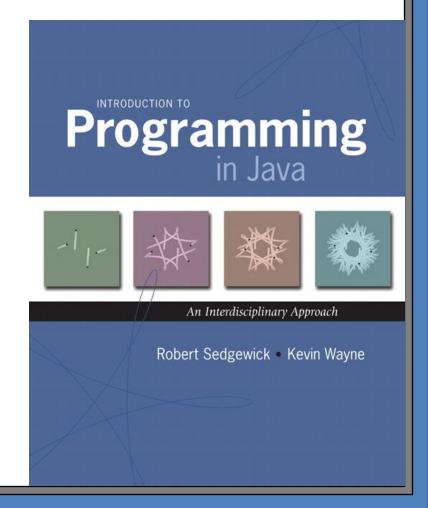
- Chat with me after lecture.
- Talk after recitation (with instructor).
- Schedule an appointment with me.

Textbook: Algorithms

Robert Sedgewick and Kevin Wayne



Optional: Introduction to Programming in Java Robert Sedgewick and Kevin Wayne



VISUALGO.NET



PSA: What if I hate all this?

Think about whether CS is for you!

- I think CS is really neat!
- I also think many other fields are really neat!
- Many philosophies for a happy life...
 - Philosophy 1: Do what you are good at.
 - Philosophy 2: Do what you love.
- Either way, you have many options!

Admin note:

It is easier to apply for a departmental transfer before the end of Year 1.

Stress?

If stress is getting high:

- Chat with University Counselling Services.
- They have lots of good advice on stress management, organization, how to achieve your goals!

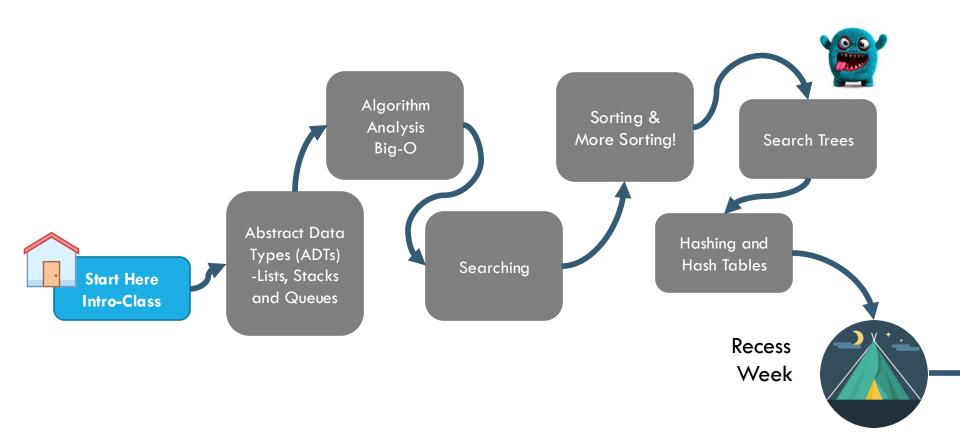
(Counselling is *not* for sick people---it's for people who need advice.)

Take care of your mental health!

What is CS2040S about?



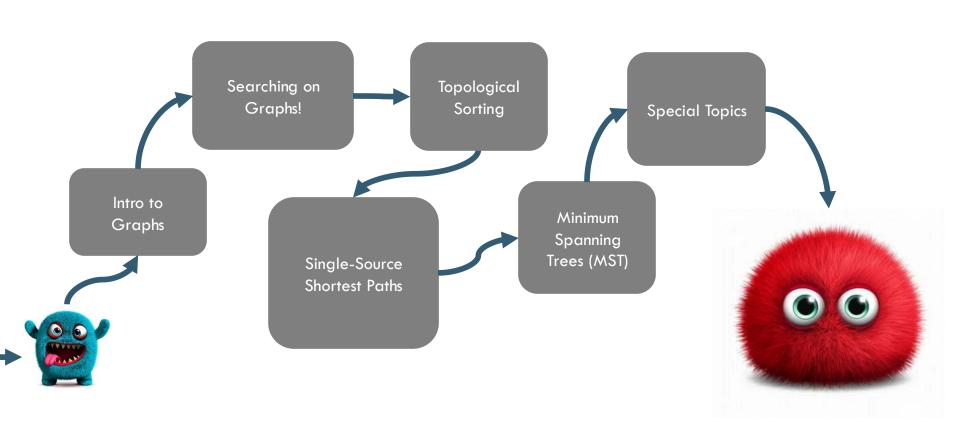
COURSE STRUCTURE



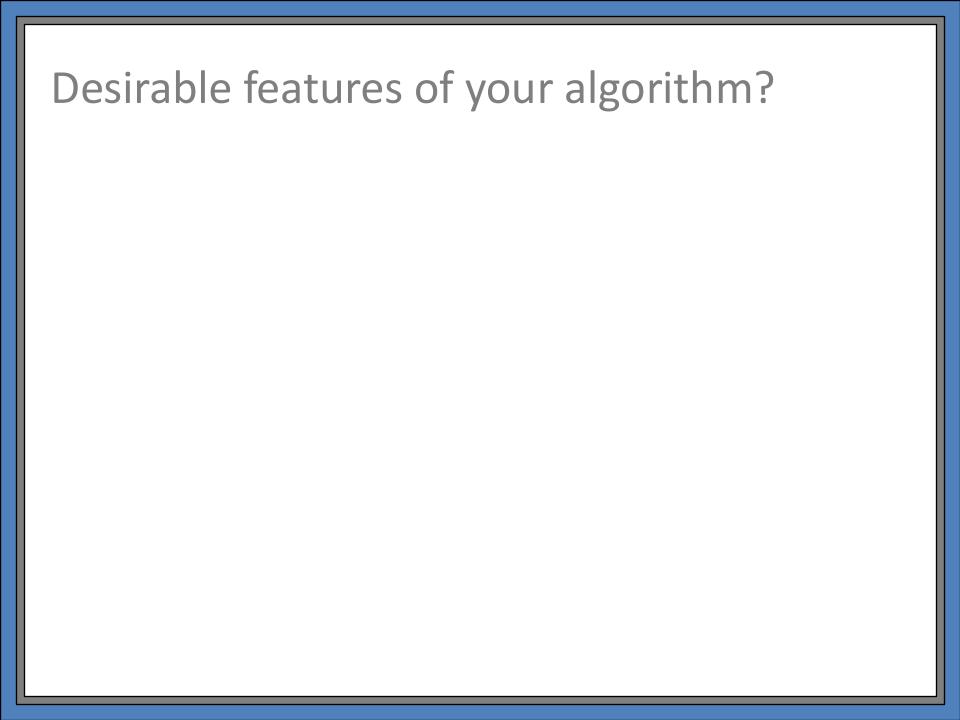
PART I: ORGANIZING YOUR DATA



COURSE STRUCTURE



PART II: MODELLING AND SOLVING PROBLEMS



Desirable features of your algorithm:



Desirable features of your algorithm:



How do you choose the right algorithm for the right problem?

How do you design new algorithms for new problems?

Algorithms

Goals of this course:

- How to organize and manipulate data?
 - Efficiency
 - -Time: How long does it take?
 - -Space: How much memory? How much disk?
 - Others: Energy, parallelism, etc.
 - Scalability
 - Inputs are large: e.g., the internet.
 - Bigger problems consume more resources.
- Solve real (fun!) problems

FRAMEWORK: ALGORITHM & DATA STRUCTURE

What problem does it solve?

How does it work?

How to implement it?

What is its asymptotic performance?

What is its real world performance?

What are the trade-offs?



GOALS

By the end of this course, you should be able to:

- Apply algorithmic thinking and techniques for solving computational problems.
- Describe the structure and operation of different data structures and algorithms under the standard computational model.
- Assess the suitability of different data-structures and algorithms for a specific computational problem.
- Adapt existing data-structures and algorithms to solve specific computational problems.

A little algorithmic thinking...

Searching

Finding a big item.

Finding many big items.

Searching an array.

Searching faster?

Similarity Test

Comparing two texts

An algorithm

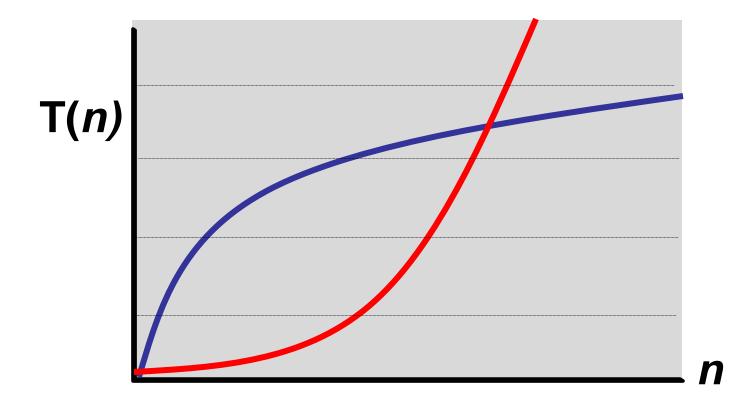
Performance profiling

Fixing a few mistakes

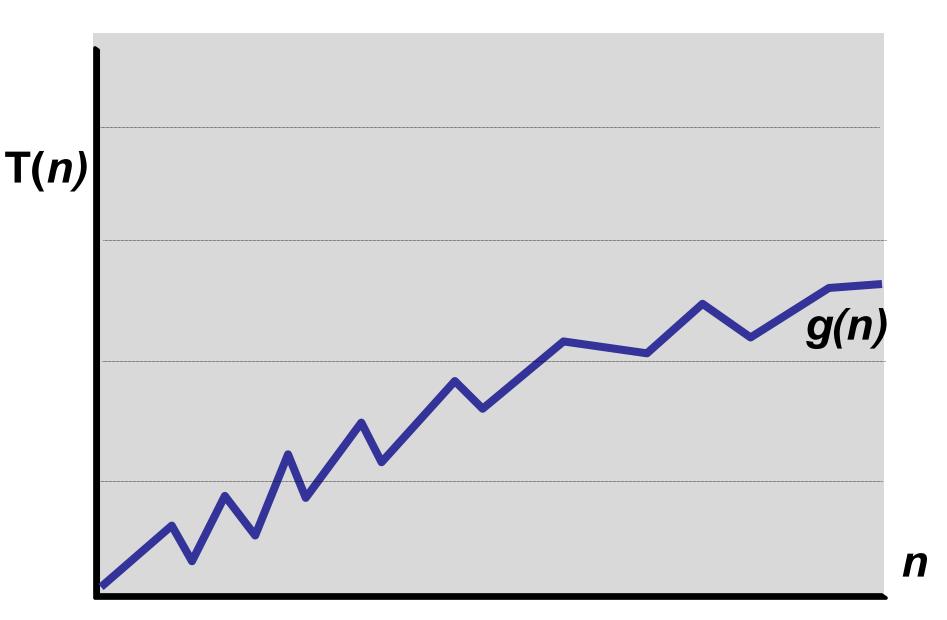
Big-O Notation

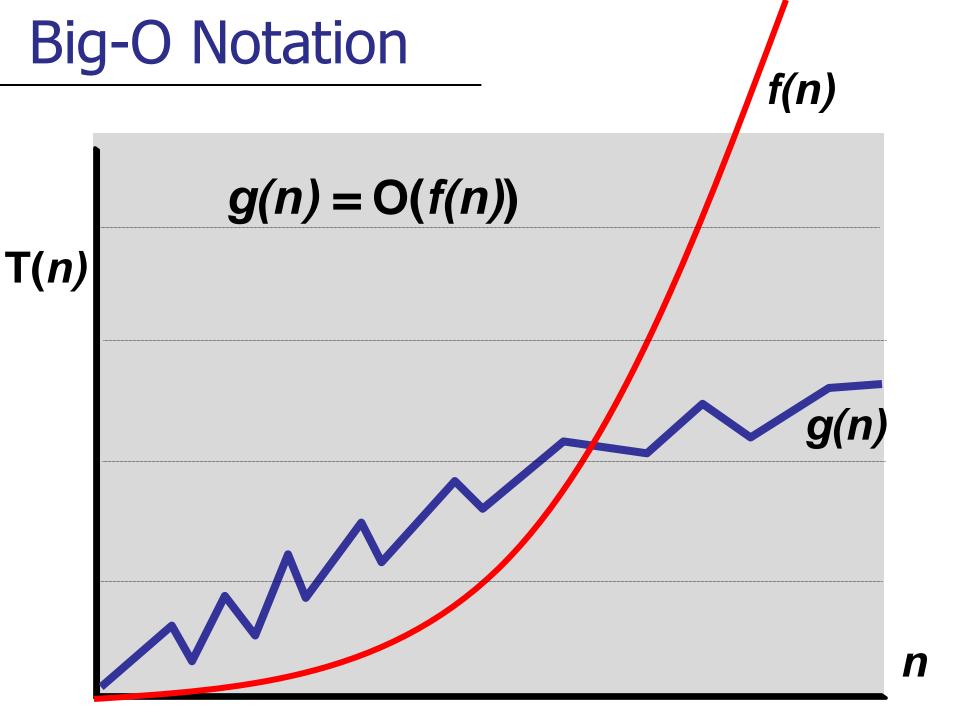
How does an algorithm scale?

- For large inputs, what is the running time?
- -T(n) = running time on inputs of size n

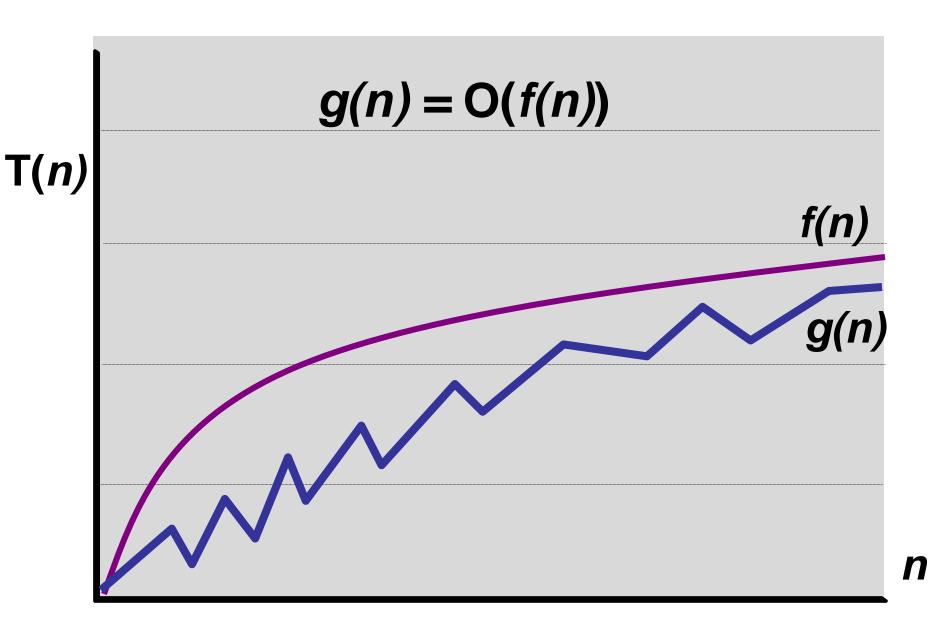


Big-O Notation





Big-O Notation



Example

Γ(n) big-O			•
	Γ(n)	f(n)	big-O

f(n) = nT(n) = 1000nT(n) = O(n)

 $f(n) = n^2$ T(n) = 1000n $T(n) = O(n^2)$

 $T(n) = n^2$ f(n) = n $T(n) \neq O(n)$ $T(n) = 13n^2 + n$ $f(n) = n^2$ $T(n) = O(n^2)$

Do asymptotics matter?

Algorithms are more important than language:

Fact: C can be 20x as fast as Python!

Algorithm	Language	Time	10,000 elements
Fast (MergeSort)	Slow (Python)	2n log(n) μs	0.266s
Slow (InsertionSort)	Fast (C)	0.01n ² μs	1s

(Source: MIT 6.006)

Given:

a collection of n items

Find:

the largest item in the collection

- 1. Sort the collection.
- 2. Return the largest item.

Given:

a collection of n items

Find:

the largest item in the collection

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest item.

- 1. Iterate through the collection.
- 2. Return the largest item.

Given:

a collection of n items

Find:

the largest item in the collection

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest item.

 $O(n \log n)$

- 1. Iterate through the collection.
- 2. Return the largest item.

Given:

a collection of n items

Find:

the largest item in the collection

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest item.

 $O(n \log n)$

Algorithm 2.

- 1. Iterate through the collection.
- 2. Return the largest item.

O(n)

Given:

a collection of n items

Find:

the largest k items in the collection

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

- Iterate through the collection. Store the largest k items seen.
- 2. Return the largest k items.

Extra memory: 1

Time: $O(n \log n + k)$

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

Extra memory: k

Time: ??

- Iterate through the collection. Store in order the largest k items seen.
- 2. Return the largest k items.

Extra memory: 1

Time: $O(n \log n + k)$

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

Extra memory: k

Time: ??

Store k largest items in a sorted array. Insert new big items into the array.

- Iterate through the collection. Store in order
 - the largest k items seen.
- 2. Return the largest k items.

Example: k = 3

 Collection:
 1
 3
 5
 6
 2
 7
 1
 3
 1

1

Big items: 1

Cost of processing first item: 1

Example: k = 3

 Collection:
 1
 3
 5
 6
 2
 7
 1
 3
 1
 1

1

Big items: 1 3

Cost of processing second item: 1

Example: k = 3

 Collection:
 1
 3
 5
 6
 2
 7
 1
 3
 1
 1

Big items: 1 3 5

Cost of processing third item: 1

Example: k = 3

Collection:



ı

Big items:

1 | 3 | 5

Example: k = 3

Big items: 3 5

Example: k = 3

Collection:



Big items: 3

Example: k = 3

Collection:



Big items: 3

Example: k = 3

Collection:



Big items:

 $3 \mid 5$

Example: k = 3

Collection: $\begin{bmatrix} 1 & 3 & 5 & 6 & 2 & 7 & 1 & 3 & 1 & 1 \\ & & & & & & & \end{bmatrix}$

Big items: $3 \quad 5 \quad 6$

Cost of processing fourth item: k

Example: k = 3

 Collection:
 1
 3
 5
 6
 2
 7
 1
 3
 1
 1

Big items: 3 5 6

Cost of processing fifth item: 0

Example: k = 3

Collection: $\begin{bmatrix} 1 & 3 & 5 & 6 & 2 & 7 & 1 & 3 & 1 & 1 \\ & & & & & & & & \end{bmatrix}$

Big items: $3 \quad 5 \quad 6$

Cost of processing: k

Example: k = 3

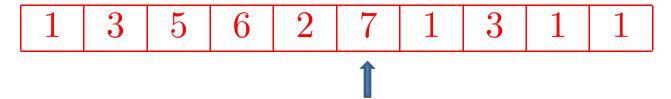
 Collection:
 1
 3
 5
 6
 2
 7
 1
 3
 1
 1

Big items: 5 6 7

Cost of processing: k

Example: k = 3

Collection:



Big items:

5 | 6 | 7

Ask yourself:

What is the worst-case time to process the entire collection?

Example: k = 3

Collection: $\begin{bmatrix} 1 & 3 & 5 & 6 & 2 & 7 & 1 & 3 & 1 & 1 \\ & & & & & & & & & & & \end{bmatrix}$

Big items: $5 \quad 6 \quad 7$

Worst case cost to process collection: O(nk)

Extra memory: 1

Time: $O(n \log n + k)$

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

Extra memory: k

Time: O(nk)

- Iterate through the collection. Store in order the largest k items seen.
- 2. Return the largest k items.

Extra memory: 1

Time: $O(n \log n + k)$

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

Extra memory: k

Time: O(nk)

Better solution in Week 5.

- Iterate through the collection. Store the largest k items seen.
- 2. Return the largest k items.

Extra memory: 1

Time: $O(n \log n + k)$

Algorithm 1.

- 1. Sort the collection.
- 2. Return the largest k items.

Extra memory: k

Time: O(nk)

Ask Yourself:

Between these two,

Can we say one algorithm is always better than the other?

- Iterate through the collection. Store the largest k items seen.
- 2. Return the largest k items.

Given:

a sorted array of n items

Find:

find smallest item at least as big as value x?

Algorithm 1. Linear search

Given:

a sorted array of n items

Find:

find smallest item at least as big as value x?

Week 2:

Check midpoint and recurse into one of the two halves.

Algorithm 1. Linear search

Algorithm 2. Binary search

Given:

a sorted array of n items

Find:

find smallest item at least as big as value x?

Algorithm 1. Linear search

O(n)

Algorithm 2. Binary search

 $O(\log n)$

Are we done?

Is binary search always the best searching algorithm on real machines?

Algorithm 1. Linear search

O(n)

Algorithm 2. Binary search

 $O(\log n)$

Are we done?

Information theory says that this problem requires at least $\Omega(\log n)$ ops.

But that's not the whole story...

Algorithm 1. Linear search

O(n)

Algorithm 2. Binary search

 $O(\log n)$

Predicting Performance

Example: 100 TB of data

- 1) Store data sorted in an array
 - \Rightarrow Scan all the data: O(n)
 - ⇒ (Binary) search: O(log n)
- 2) Store data in a linked list
 - \Rightarrow Scan all the data: O(n)
 - ⇒ Search: O(n)
- 3) Store data in a balanced tree
 - \Rightarrow Scan all the data: O(n)
 - ⇒ Search: O(log n)

Predicting Performance

Example: 100 TB of data

- 1) Store data sorted in an array
 - ⇒ Scan all the data: O(n) ◆
 - ⇒ (Binary) search: O(log n)
- 2) Store data in a linked list
 - ⇒ Scan all the data: O(n) ◆
 - ⇒ Search: O(n)
- 3) Store data in a balanced tree
 - ⇒ Scan all the data: O(n) ←
 - ⇒ Search: O(log n)

Which is fastest?

Predicting Performance

Example: 100 TB of data

- 1) Store data sorted in an array
 - ⇒ Scan all the data: O(n) ◆
 - ⇒ (Binary) search: O(log n)
- 2) Store data in a linked list
 - ⇒ Scan all the data: O(n) ◆
 - ⇒ Search: O(n)
- 3) Store data in a balanced tree
 - ⇒ Scan all the data: O(n) ←
 - ⇒ Search: O(log n)

Same theoretical performance!

Predicting Performance

Example: 100 TB of data

- 1) Store data sorted in an array
 - ⇒ Scan all the data: O(n)
 - ⇒ (Binary) search: O(log n)

This is MUCH faster in practice by an order of magnitude.

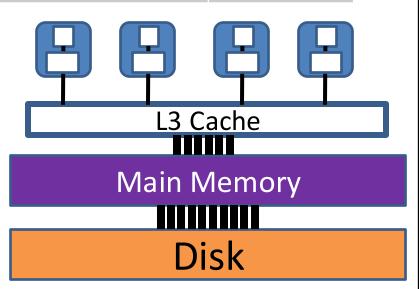
- 2) Store data in a linked list
 - \Rightarrow Scan all the data: O(n)
 - ⇒ Search: O(n)
- 3) Store data in a red-black tree
 - \Rightarrow Scan all the data: O(n)
 - ⇒ Search: O(log n)

Haswell Architecture (2-18 cores)

Memory Type	size	line size	clock cycles
L1 cache	64 KB	64 B	~4
L2 cache	256 KB	64 B	~10
L3 cache	2-40 MB	64 B	40-74
L4 (optional)	128 MB		
Main Memory	< 128 GB	16 KB	~200-350
SSD Disk	BIG	Variable (e.g., 16KB)	~20,000
Disk	BIGGER	Variable (e.g., 16KB)	~20,000,000

Notes:

- Several other "caches" e.g., TLB, micro-op cache, instruction cache, etc.
- L1/L2 caches are per core.
- L3/L4 cache are shared per socket.
- Main memory shared cross socket.



Sneak Peek: More Searching

Faster to keep your data stored in a B-tree than in a sorted array!

Algorithm 1. Linear search

Array: $O(\log(n/B))$

Algorithm 2. Binary search

B-tree: $O(\log_B n)$

Algorithm 3.
Store data in B-tree

B is a parameter of the cache.

A little algorithmic thinking...

Searching

Finding a big item.

Finding many big items.

Searching an array.

Searching faster?

Similarity Test

Comparing two texts

An algorithm

Performance profiling

Fixing a few mistakes

A little algorithmic thinking...

Searching

Similarity Test

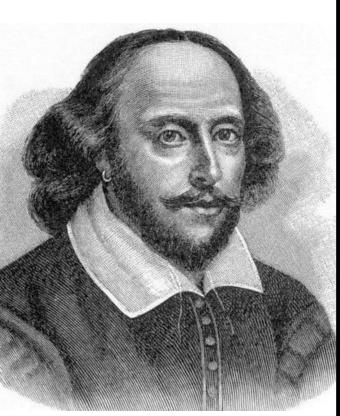
Fii

Fi

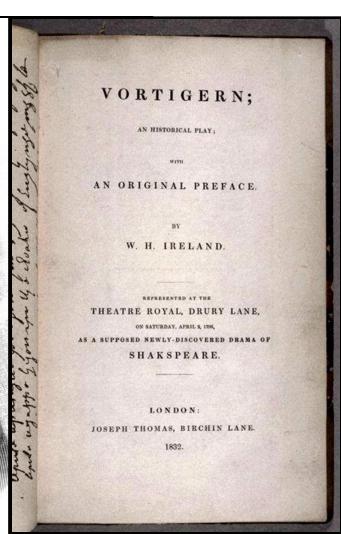
Se

Let's see how we might apply algorithmic thinking in our problem solving.

Who wrote this?



William Shakespeare??



mystery play "found" in 1796



William Henry Ireland??

Document distance

- How similar are two documents?
 - Are two documents written by the same author?
 - Detect forgeries
 - Find plagiarism / cheating
 - Was Homer one author or many?

What does "similar" mean?

Metrics of similarity

What are good ways to quantify similarity?

Metrics of similarity

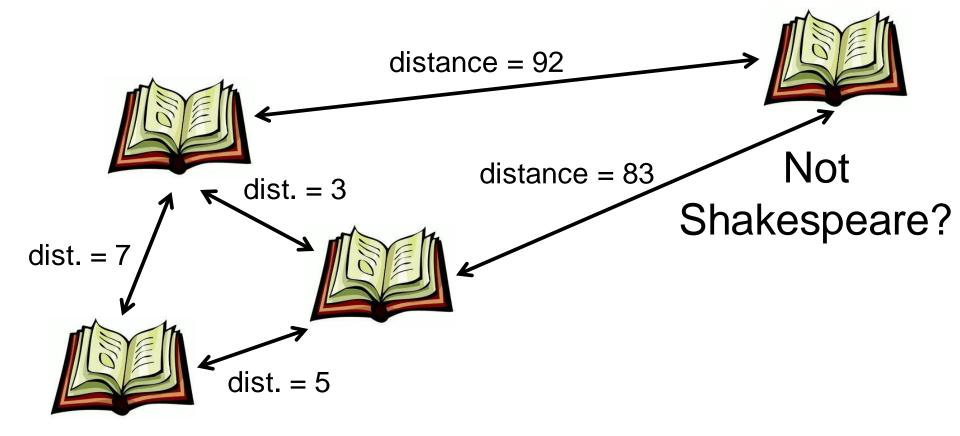
- Binary: (e.g., detect plagiarism)
 - Exactly same words in same order

Scalar:

- Number of words in the same order
- Number of shared uncommon words
- Same # of words per sentence
- Same ratio of adjectives / nouns
- Written on similar paper / using similar ink

Document distance

How similar are two documents?



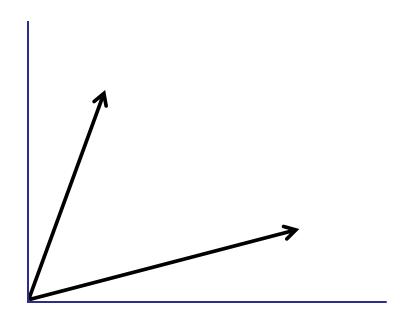
Shakespeare

Vector Space Model

Strategy:

View each document as a high-dimensional vector.

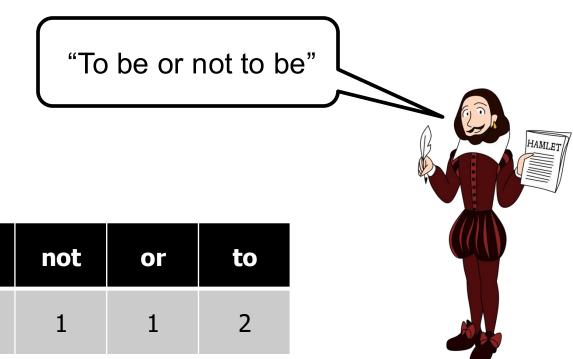
[Salton, Wang, Yang '75]



Vector Space Model

Strategy:

Each dimension is # times each word appears.



4d-vector:

be

Vector Space Model

Strategy:

- View each document as a high-dimensional vector.
- The metric of similarity is the angle between the two vectors.

- Identical: $\Phi = 0$
- No words in common: $\Phi = \pi/2$

Compare Two Documents

Given: documents A and B

- 1. Create vectors v_A and v_B
 - count number of times each word appears
- 2. Calculate vector norms
- 3. Calculate dot product: $(v_A \cdot v_B)$
- 4. Calculate angle $\Phi(v_A, v_B)$

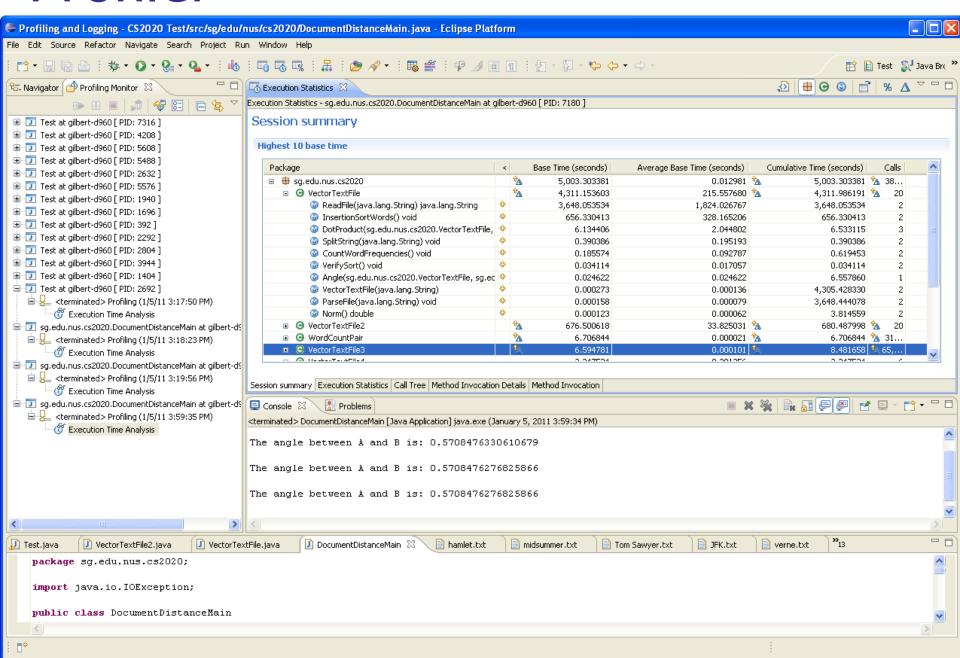
Performance Profiling

(Dracula vs. Lewis & Clark)

Step	Function
Create vectors:	Read each file
	Parse each file
	Sort words in each file
	Count word frequencies
Dot product:	
Norm:	
Angle:	
Total:	

72minutes ≈ **4,311.00s**

Profiler



Performance Profiling

(Dracula vs. Lewis & Clark)

Step	Function	Running Time
Create vectors:	Read each file	1,824.00s
	Parse each file	0.20s
	Sort words in each file	328.00s
	Count word frequencies	0.31s
Dot product:		6.12s
Norm:		3.81s
Angle:		6.56s
Total:		72minutes ≈ 4,311.00s

Problem 1: Why does it take SO long to read the file?

ReadFile (excerpt)

```
// Open the file as a stream and find its size
inputStream = new FileInputStream(fileName);
iSize = inputStream.available();
// Read in the file, one character at a time, normalizing as we go.
for (int i=0; i<iSize; i++)</pre>
    // Read a character
    char c = (char)inputStream.read();
    // Ensure that the character is lower-case
    c = Character.toLowerCase(c):
    // Check if the character is a letter
    if (Character.isLetter(c))
        strTextFile = strTextFile + c:
    // Check if the character is a space or an end-of-line marker
    else if (((c == ' ') || (c == ' \setminus n')) \&\& (!strTextFile.endsWith(" ")))
                strTextFile = strTextFile + ' ':
```

What happens when:

textFile = textFile + c

- 1. Creates new temporary string.
- 2. Copies textFile to the new string.
- 3. Adds the new character *c*.
- 4. Reassigns textFile to point to the new string.

What happens when:

textFile = textFile + c

- 1. Creates new temporary string.
- 2. Copies textFile to the new string.
- 3. Adds the new character *c*.
- 4. Reassigns textFile to point to the new string.

Copying a string of k characters takes time k!

How long to read in a file of n characters?.

How long to read in a file of n characters?.

$$1 + 2 + 3 + 4 + ... + n = n(n+1)/2 = \Theta(n^2)$$

Very, very, very slow!

Fix the string problem!

```
// Open the file as a stream and find its size
inputStream = new FileInputStream(fileName);
iSize = inputStream.available();
// Initialize the char buffer to be arrays of the appropriate size.
charBuffer = new char[iSize];
// Read in the file, one character at a time, normalizing as we go.
for (int i=0; i<iSize; i++)</pre>
    // Read a character
    char c = (char)inputStream.read();
    // Ensure that the character is lower-case
    c = Character.toLowerCase(c);
    // Check if the character is a letter
    if (Character.isLetter(c))
        charBuffer[iCharCount] = c;
        iCharCount++:
    // Check if the character is a space or an end-of-line marker
    else if (((c == ' ') || (c == '\n')) && (!strTextFile.endsWith(" ")))
        charBuffer[iCharCount] = ' ';
        iCharCount++:
```

Performance Profiling, V2

(Dracula vs. Lewis & Clark)

Step	Function	Running Time
Create vectors:	Read each file	1.09s
	Parse each file	3.68s
	Sort words in each file	332.13s
	Count word frequencies	0.30s
Dot product:		6.06s
Norm:		3.80s
Angle:		6.06s
Total:		11 minutes ≈ 680.49s

Problem 2: Can we sort faster?

Performance Profiling

(Dracula vs. Lewis & Clark)

Version	Change	Running Time
Version 1		4,311.00s
Version 2	Better file handling	676.50s
Version 3	Faster sorting	6.59s

Use O(n log n) sorting algorithm (MergeSort) instead of $O(n^2)$ sorting algorithm (InsertionSort).

Problem 3: Do we need to sort at all?

Document Distance

"If you need your software to run twice as fast, hire better programmers."

But if you need your software to run more than twice as fast, use a better algorithm."

-- Software Lead at Microsoft

Version	Change	Running Time
Version 1		4,311.00s
Version 2	Better file handling	676.50s
Version 3	Faster sorting	6.59s
Version 4	No sorting!	2.35s

Use O(n) hashing approach instead of O(n log n) sorting algorithm (MergeSort).

Goals for the Semester

Speed up your code by a factor of 2040!

Algorithms:

- Design of efficient algorithms
- Analysis of algorithms

Implementation:

- Solve real problems
- Analyze and measure performance
- Improve performance via better algorithms

For next time...

Wednesday lecture:

- Java introduction, OOP
- Stacks, Queues, Lists

Lecture 1 Review: Released at 6pm. Due tomorrow, 6pm.

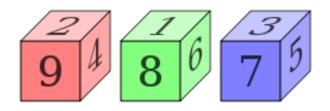
Problem Set 1: Released at 6pm. Due Monday morning.

Important: Deadline TODAY 11:59pm.

- Fill out recitation survey (Coursemology)
- Fill out tutorial survey (Coursemology)

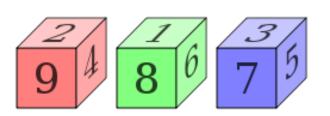
Imagine three dice:

- A has six sides: 2, 2, 4, 4, 9, 9
- B has six sides: 1, 1, 6, 6, 8, 8
- C has six sides: 3, 3, 5, 5, 7, 7



Imagine three dice:

- A has six sides: 2, 2, 4, 4, 9, 9
- B has six sides: 1, 1, 6, 6, 8, 8
- C has six sides: 3, 3, 5, 5, 7, 7

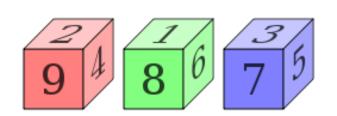


Game:

- Alice chooses a die.
- Bob chooses a die.
- Alice and Bob both roll.
- The higher value wins.

Imagine three dice:

- A has six sides: 2, 2, 4, 4, 9, 9
- B has six sides: 1, 1, 6, 6, 8, 8
- C has six sides: 3, 3, 5, 5, 7, 7



Game:

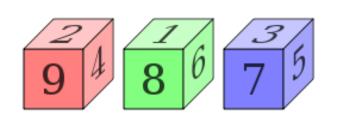
- Alice chooses a die.
- Bob chooses a die.
- Alice and Bob both roll.
- The higher value wins.

Questions:

- Which die should Alice choose?
- Which die should Bob choose?
- Who is more likely to win?

Imagine three dice:

- A has six sides: 2, 2, 4, 4, 9, 9
- B has six sides: 1, 1, 6, 6, 8, 8
- C has six sides: 3, 3, 5, 5, 7, 7



Game:

- Alice chooses a die.
- Bob chooses a die.
- Alice and Bob both roll.
- The higher value wins.

Questions:

- Which die should Alice choose?
- Which die should Bob choose?
- Who is more likely to win?