

DSA-CL III – WINTER TERM 2018-19

DATA STRUCTURES AND ALGORITHMS FOR COMPUTATIONAL LINGUISTICS III

CLAUS ZINN

Çağrı Çöltekin



<https://dsacl3-2018.github.io>

DSA-CL III course overview

What is DSA-CL III?

- Intermediate-level survey course.
- Programming and problem solving, with applications.
 - **Algorithm:** method for solving a problem.
 - **Data structure:** method to store information.
- Second part focused on Computational Linguistics

Prerequisites:

- Data Structures and Algorithms for CL I
- Data Structures and Algorithms for CL II

Lecturers:

- Çağrı Cöltekin
- Claus Zinn

Tutors:

- Marko Lozajic
- Michael Watkins

Slots:

- Mon 12:15 & 18:00
- Wed 14:15 — 18:00 (lab)

Course Materials: <https://dsacl3-2018.github.io>

Why study algorithms?

Their impact is broad and far-reaching.

Internet. Web search, packet routing, distributed file sharing, ...

Biology. Human genome project, protein folding, ...

Computers. Circuit layout, file system, compilers, ...

Computer graphics. Movies, video games, virtual reality, ...

Security. Cell phones, e-commerce, voting machines, ...

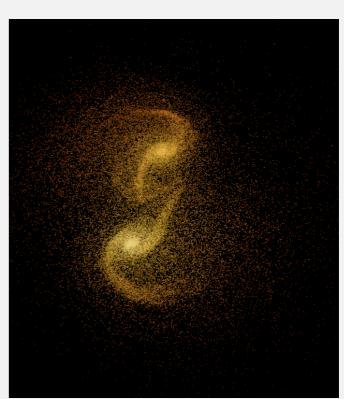
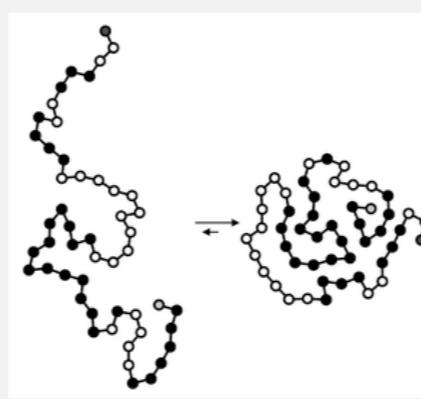
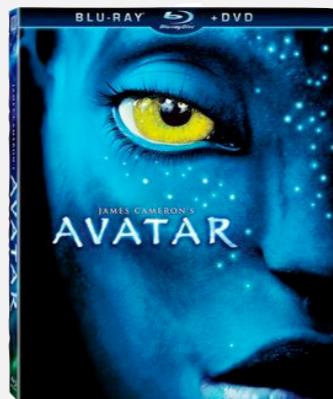
Multimedia. MP3, JPG, DivX, HDTV, face recognition, ...

Social networks. Recommendations, news feeds, advertisements, ...

Physics. N-body simulation, particle collision simulation, ...

⋮

Google™
YAHOO![®]
bing™



Why study algorithms?

Their impact is broad and far-reaching.

Mysterious algorithm was 4% of trading activity last week

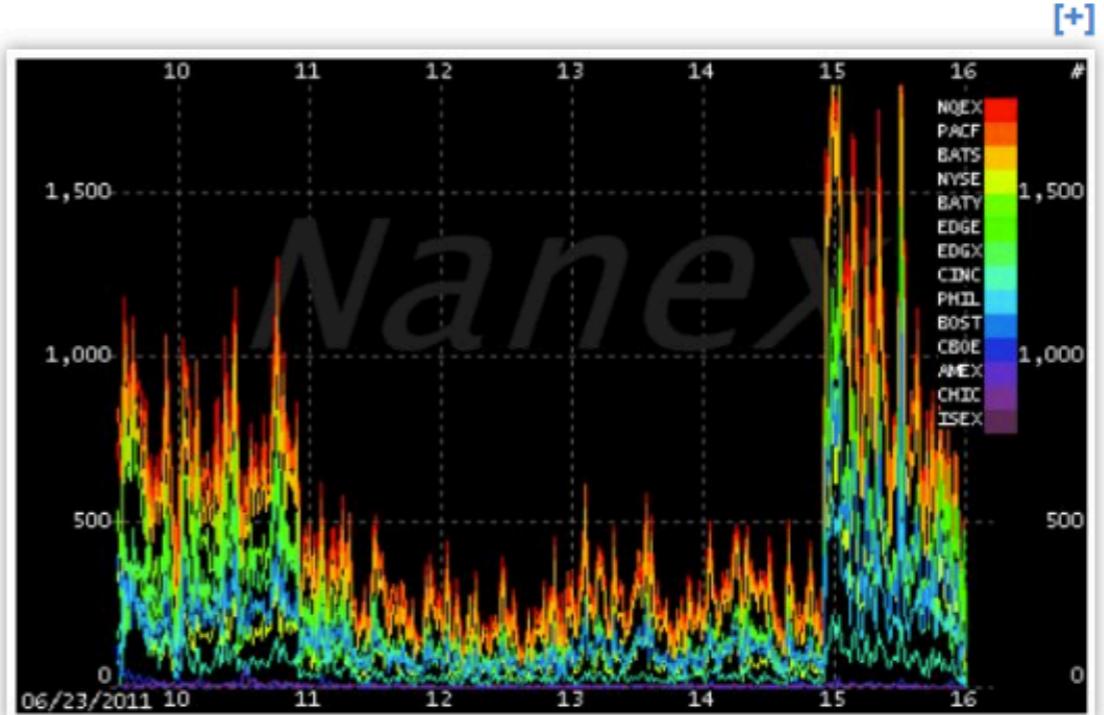
October 11, 2012

A single mysterious computer program that placed orders — and then subsequently canceled them — made up 4 percent of all quote traffic in the U.S. stock market last week, according to the top tracker of [high-frequency trading](#) activity.

The motive of the algorithm is still unclear, [CNBC](#) reports.

The program placed orders in 25-millisecond bursts involving about 500 stocks, according to Nanex, a market data firm. The algorithm never executed a single trade, and it abruptly ended at about 10:30 a.m. ET Friday.

"My guess is that the algo was testing the market, as high-frequency frequently does," says Jon Najarian, co-founder of TradeMonster.com. "As soon as they add bandwidth, the HFT crowd sees how quickly they can top out to create latency." ([Read More: Unclear What Caused Kraft Spike: Nanex Founder.](#))



Generic high frequency trading chart (credit: Nanex)

Why study algorithms?

For intellectual stimulation.

“For me, great algorithms are the poetry of computation. Just like verse, they can be terse, allusive, dense, and even mysterious. But once unlocked, they cast a brilliant new light on some aspect of computing.” — Francis Sullivan

FROM THE
EDITORS

THE JOY OF ALGORITHMS

Francis Sullivan, Associate Editor-In-Chief



THE THEME OF THIS FIRST-OF-THE-CENTURY ISSUE OF COMPUTING IN SCIENCE & ENGINEERING IS ALGORITHMS. IN FACT, WE WERE BOLD ENOUGH—AND PERHAPS FOOLISH ENOUGH—to call the 10 examples we've selected “THE TOP 10 ALGORITHMS OF THE CENTURY.”

Computational algorithms are probably as old as civilization. Sumerian cuneiform, one of the most ancient written records, contains what may be the first algorithm, for calculating the area of a trapezoid. It was probably used for surveying land, which he did. And I suppose we could claim that the Druid algorithm for estimating the start of summer is enshrined in Stonehenge? That's probably not true, but it's a nice thought.

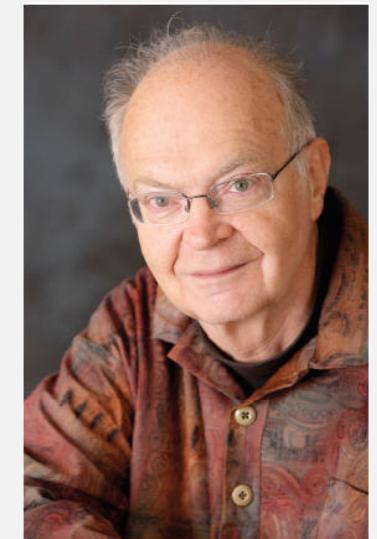
Like so many other things that technology affects, algorithms have advanced in startling and unexpected ways in the last few decades. In fact, the 10 algorithms of the century we chose for this issue have been essential for progress in communications, health care, manufacturing, economics, and more. The search for better algorithms, and the variety of progress in these areas has stimulated the search for even more algorithms. Recall one late-night session on the March 1998 cover when we asked “What's the next craze? After all, they don't look very appealing.” After the usual speculation about the future behavior of stock gulls, someone gave what must be the right answer: “A very smart person first ate a crab.”

The reason for this is the modern “law of Moore,” “The invention creates its own necessity.” Our need for powerful machines always exceeds their availability. Each significant computer advance creates a demand for more power, and that creates larger, competition to be done. New algorithms are an attempt to bridge the gap between the demand for cycles and the available data is still almost untouched. There are still very big challenges coming from more “traditional” tasks, too. For example, progress in speech recognition has been slow, but the recent breakthroughs in handwriting recognition are impressive. A large turning-point calculation is likely to be correct. Think of the way that Moore's Law has transformed the computer industry. It is not small, but the added confidence in its power is huge.

Is there an analog for things such as large, multidisciplinary teams? I suspect there is, but I don't know what it is. Solvable methods for solving specific cases of “impossible” problems always exceed their availability. Each significant computational advance creates a demand for more power, and that creates larger, competition to be done. New algorithms are an attempt to bridge the gap between the demand for cycles and the available data is still almost untouched. There are still very big challenges coming from more “traditional” tasks, too. For example, progress in speech recognition has been slow, but the recent breakthroughs in handwriting recognition are impressive. A large turning-point calculation is likely to be correct. Think of the way that Moore's Law has transformed the computer industry. It is not small, but the added confidence in its power is huge.

The next century is going to be very fruitful for us, but it is not going to be dull either. ■

“An algorithm must be seen to be believed.” — Donald Knuth



Why study algorithms?

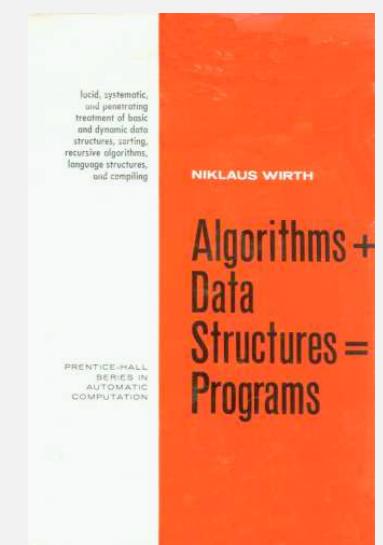
To become a proficient programmer.

“I will, in fact, claim that the difference between a bad programmer and a good one is whether he considers his code or his data structures more important. Bad programmers worry about the code. Good programmers worry about data structures and their relationships. ”

— Linus Torvalds (creator of Linux)



“Algorithms + Data Structures = Programs. ” — Niklaus Wirth



Why study algorithms?

They may unlock the secrets of life and of the universe.

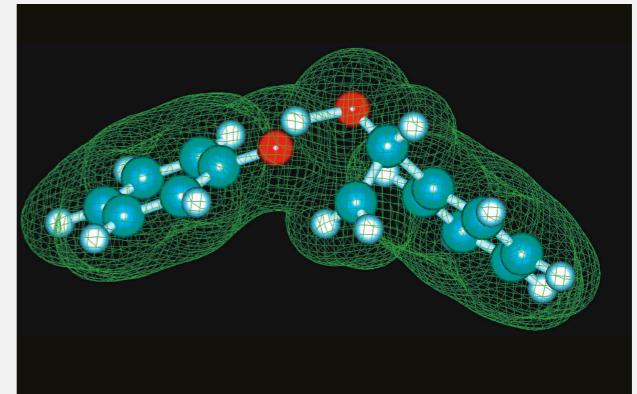
“ Computer models mirroring real life have become crucial for most advances made in chemistry today.... Today the computer is just as important a tool for chemists as the test tube. ”

— Royal Swedish Academy of Sciences

(Nobel Prize in Chemistry 2013)

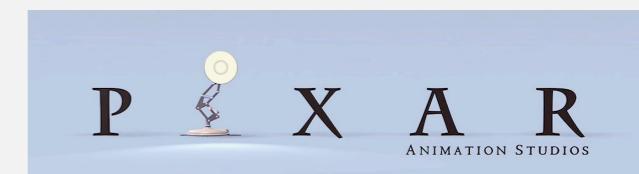
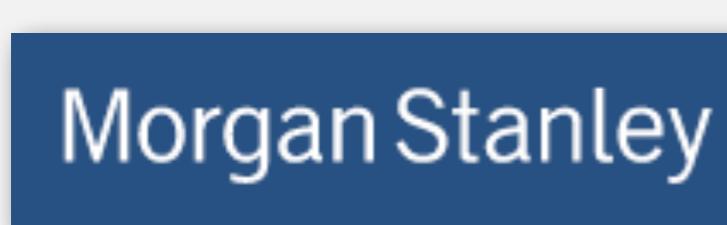
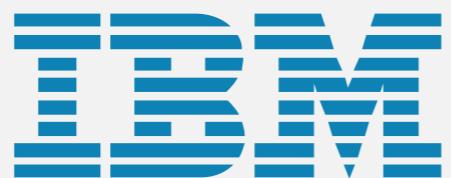


Martin Karplus, Michael Levitt, and Arieh Warshel



Why study algorithms?

For fun and profit.



Why study algorithms?

- Their impact is broad and far-reaching.
- Old roots, new opportunities.
- For intellectual stimulation.
- To become a proficient programmer.
- They may unlock the secrets of life and of the universe.
- To solve problems that could not otherwise be addressed.
- Everybody else is doing it.
- For fun and profit.

Why study anything else?



Coursework and grading

Reading material for most lectures

Weekly programming assignments

Four graded assignments. 60%

- Due on Tuesdays at 11pm via electronic submission (Github Classroom)
- Collaboration/lateness policies: see web.

Written exam. 40%

- Midterm practice exam 0%
- Final exam 40%

Honesty Statement

Honesty statement:

- Feel free to cooperate on assignments that are not graded.
- Assignments that are graded must be your own work. Do **not**:
 - Copy a program (in whole or in part).
 - Give your solution to a classmate (in whole or in part).
 - Get so much help that you cannot honestly call it your own work.
 - Receive or use outside help.
- Sign your work with the honesty statement (provided on the website).
- Above all: You are here for yourself, practice makes perfection.

Organisational issues

Presence:

- A presence sheet is circulated **purely** for statistics.
- Experience: those who do not attend lectures or do not make the assignments usually fail the course.
- Do not expect us to answer your questions if you were not at the lectures.

Office hours:

- Office hour: **Monday, 14:00-15:00**, please make an **appointment!**
- Please ask questions about the material presented in the lectures during the lectures — Everyone benefits
- We will discuss each assignment that is not graded during the next lab.

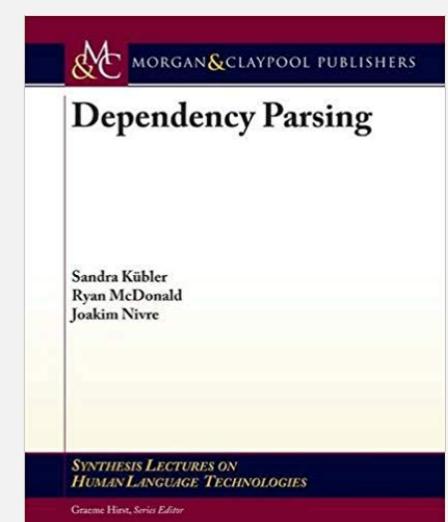
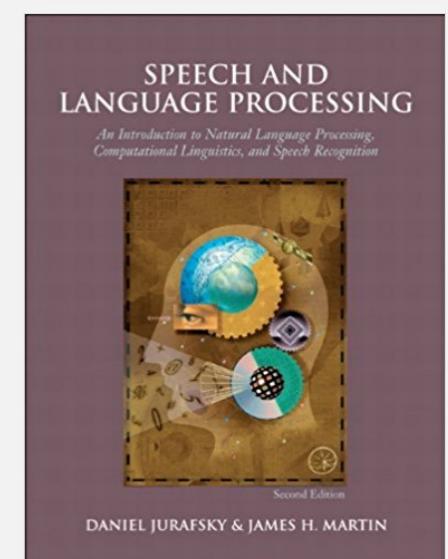
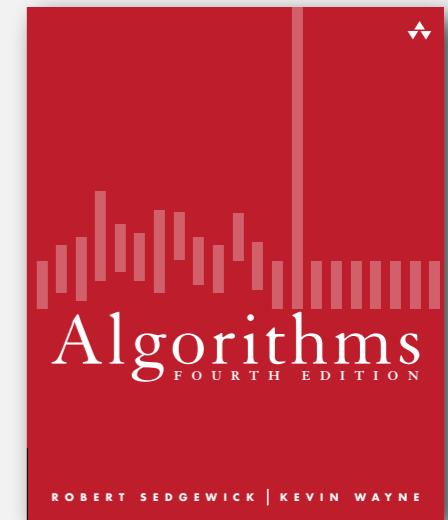
Registration:

- Do the first assignment, A0.

Resources (textbook)

Required reading.

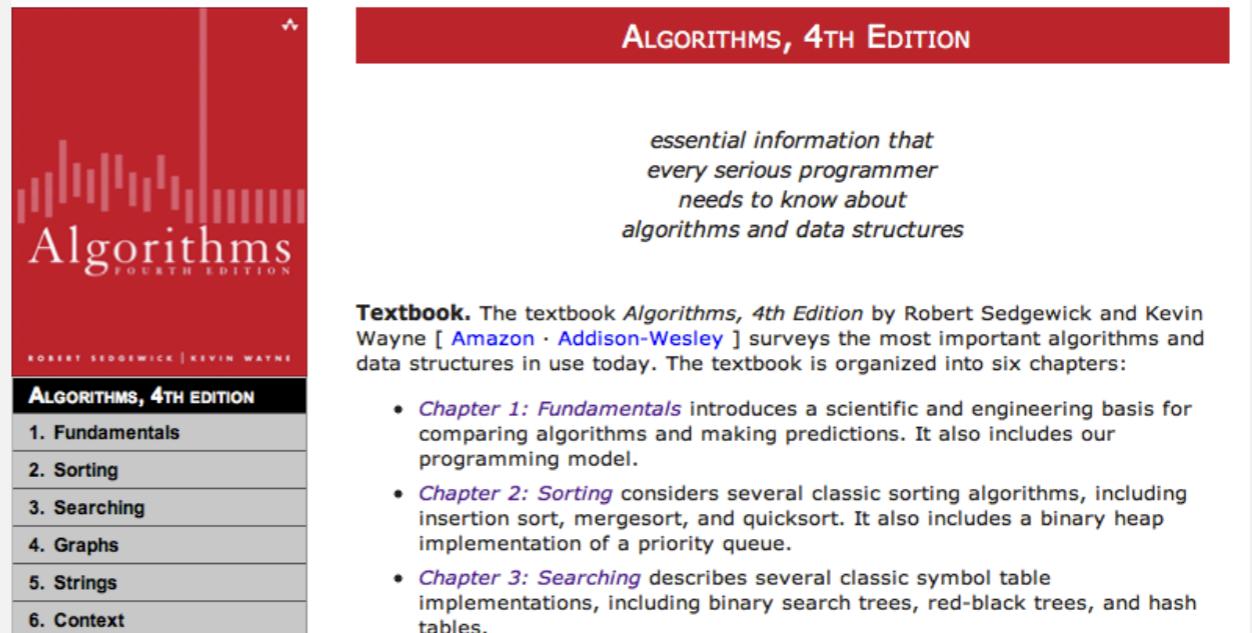
- Algorithms 4th edition by R. Sedgewick and K. Wayne, Addison-Wesley Professional, 2011, ISBN 0-321-57351-X.
 - Readable from university network thru Safari books:
 - see [proquest.tech.safaribooksonline.de/
9780132762571](http://proquest.tech.safaribooksonline.de/9780132762571)
- Speech and Language Processing, Jurafsky & Martin, 2nd Edition, Prentice Hall
 - Draft chapters of 3rd. edition available
 - see web.stanford.edu/~jurafsky/slp3/
- Dependency Parsing, Kübler, McDonald & Nivre, Morgan & Claypool



Resources (web)

Book site for first part of class

- Brief summary of content.
- Download code from book.
- APIs and Javadoc.



<http://algs4.cs.princeton.edu>

What's ahead

Week	Monday (lectures)	Wednesday (lab)	
01	Oct 22 A: <i>Introduction</i> B: <i>Complexity theory</i>	Oct 24 lab: <i>Language Guessing</i>	
02	Oct 29 A: <i>Elementary sorts</i> B: <i>Quicksort</i>	Nov 1 <i>No class</i>	
03	Nov 05 A: <i>Undirected graphs</i> B: <i>Undirected graphs</i>	Nov 07 lab: <i>Sorting</i>	
04	Nov 12 A: <i>Directed graphs</i> B: <i>Directed graphs</i>	Nov 14 lab: <i>Undirected graphs</i>	
05	Nov 19 A: <i>Distance measures</i> B: <i>Binary heaps & heapsort</i>	Nov 21 lab: <i>Directed graphs</i>	
06	Nov 26 A: <i>Binary heaps & heapsort</i> B: <i>Exam practice</i>	Nov 28 lab: <i>Burkhard-Keller trees</i>	
07	Dec 03 A: <i>Formal languages and automata</i> B: <i>Formal languages and automata</i>	Dec 05 lab: <i>TBA</i>	
08	Dec 10 A: <i>Regular grammars and finite state automata</i> B: <i>Regular grammars and finite state automata</i>	Dec 12 lab: <i>TBA</i>	

What's Ahead

	Dec 17	Dec 19
09	A: <i>Finite-state transducers</i> B: <i>Finite-state transducers and computational morphology</i>	lab: TBA
Sem. break	<i>No class</i>	<i>No class</i>
10	Jan 7 A: <i>Context-free languages and constituency parsing</i> B: <i>Context-free languages and constituency parsing</i>	Jan 9 lab: TBA
11	Jan 14 A: <i>Dependency grammars and treebanks</i> B: <i>Dependency grammars and treebanks</i>	Jan 16 lab: TBA
12	Jan 21 A: <i>Dependency parsing</i> B: <i>Dependency parsing</i>	Jan 23 lab: TBA
13	Jan 28 A: <i>A gentle introduction to classification</i> B: <i>Transition-based dependency parsing</i>	Jan 30 lab: TBA
14	Feb 04 A: <i>Exam review / practice</i> B: <i>Exam review / practice</i>	Feb 06 lab: TBA
15	Feb 11 A: <i>Exam</i> B:	Feb 13 lab: TBA

Sorting

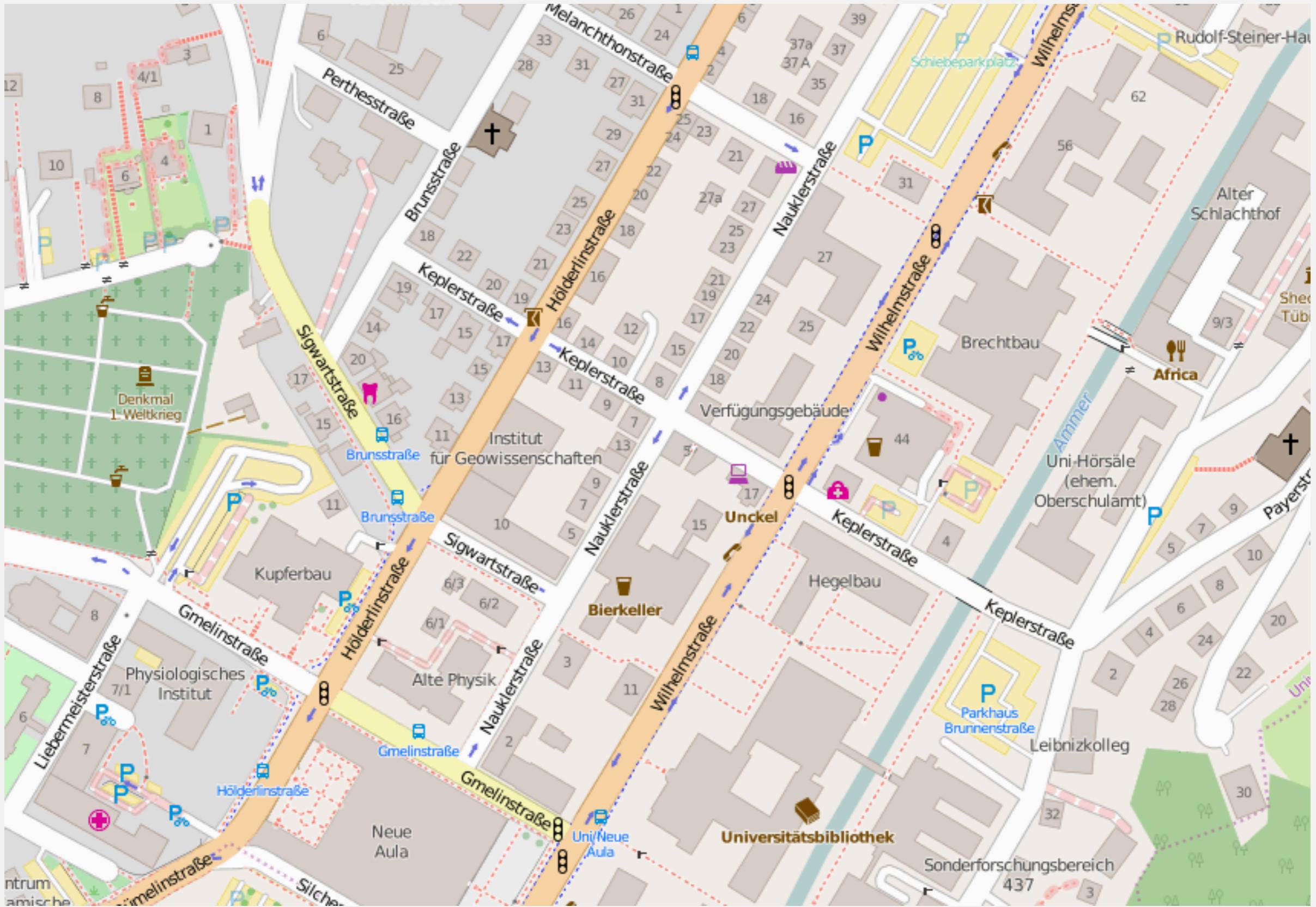


Bundesarchiv, Bild 183-22350-0001
Foto: Junge, Peter Heinz | 20. November 1953

Undirected Graphs



Directed Graphs



String Distance

The screenshot shows a search results page from a search engine. At the top left is a red logo consisting of two curved lines forming a stylized 'e'. To its right, the search term 'malorca' is displayed in a large, dark font. Below the search bar, there is a horizontal navigation menu with five items: 'All' (highlighted in blue), 'Images', 'Maps', 'News', and 'Shopp'. A thin blue horizontal bar is positioned below the menu line. The main content area displays the search results summary: 'About 114.000.000 results (0,64 seconds)'. Below this, a large bold text block says 'Showing results for **mallorca**' followed by 'Search instead for **malorca**'.

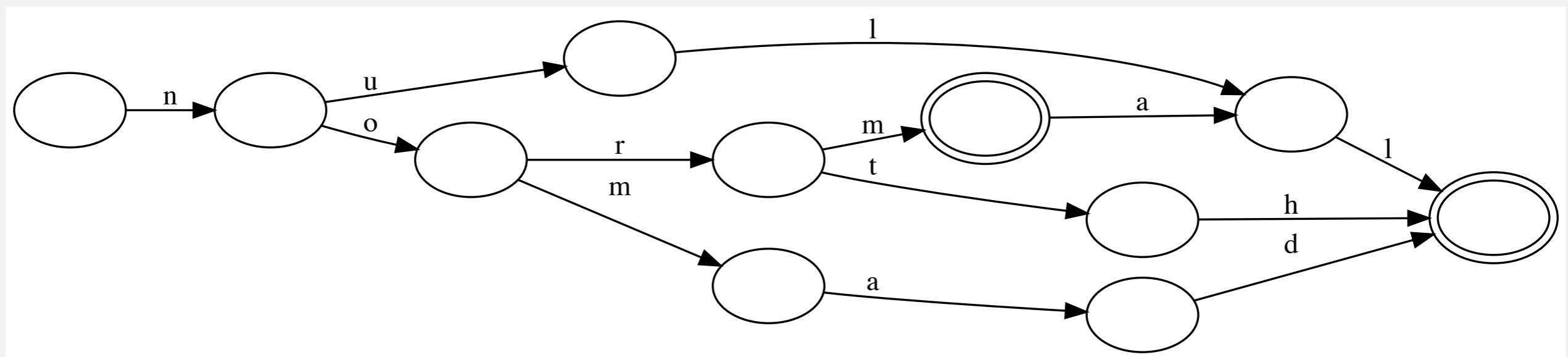
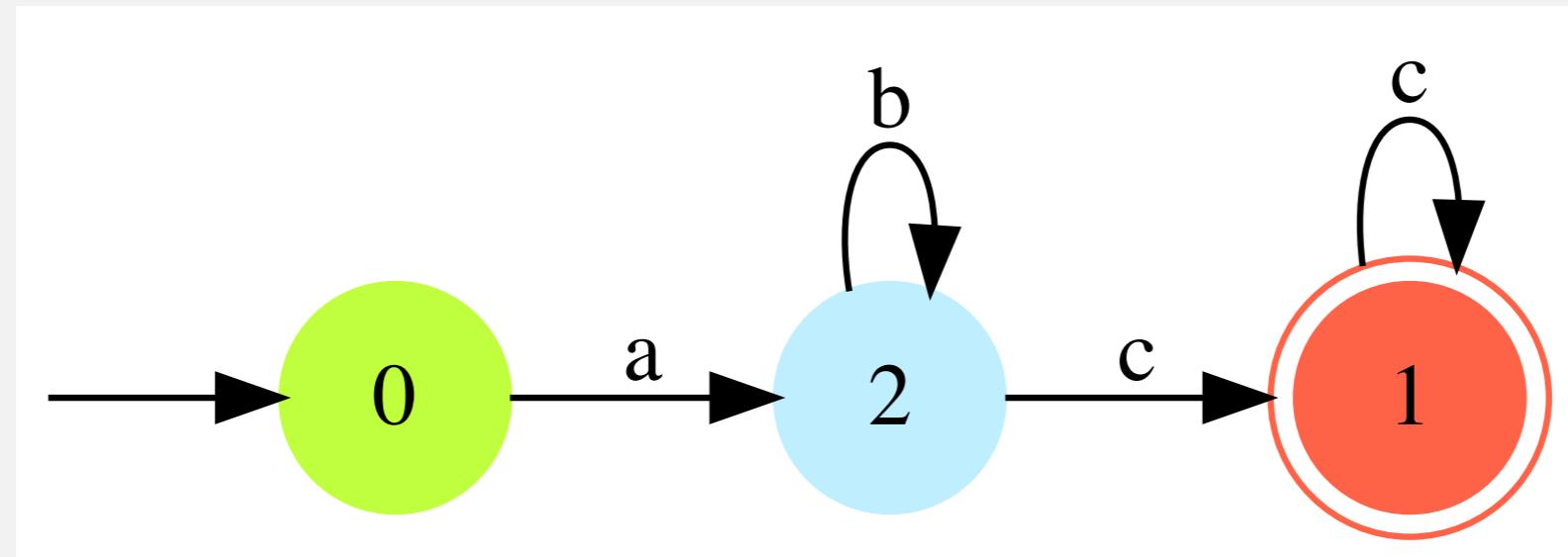
malorca

All Images Maps News Shopp

About 114.000.000 results (0,64 seconds)

Showing results for **mallorca**
Search instead for **malorca**

Finite State Automata



Parsing

