


CS216: Program and Data Representation
University of Virginia Computer Science
Spring 2006 David Evans

Lecture 4:
Dynamic
Programming,
Trees



<http://www.cs.virginia.edu/cs216>

Schedule This Week


- Sections today and Tuesday
 - Thornton D rooms
- Wednesday, Feb 1:
 - Go to Ron Rivest's talk
 - 10:30am-noon, Newcomb Hall South Meeting Room
- Thursday, Feb 2:
 - Portman Wills, "Saving the World with a Computer Science Degree"
 - 3:30pm, MEC 205 (CS290/CS390 +)

UVa CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 2

RSA

(Rivest, Shamir, Adelman 1978)

- Public-Key Encryption
 - Allows parties who have not established a shared secret to communicate securely
 - Your web browser used it when you see
- Most important algorithm invented in past 100 years
- Link from notes



UVa CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 3

Wednesday's Class

- Ron Rivest, "Security of Voting Systems"
- Newcomb Hall South Meeting Room, 10:30am-noon
- This replaces CS216 class (so you don't have to walk out at 11am)

UVa CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 4

Sequence Alignment

- Brute force algorithm in PS1


```
bestAlignment (U, V) =
  base case if |U| == 0 or |V| == 0
  otherwise f(bestAlignment (U[1:], V),
              bestAlignment (U, V[1:]),
              bestAlignment (U[1:], V[1:]))
```
- Compare to Fibonacci:


```
fibonacci (n) =
  base case if n == 0 or n == 1
  g(fibonacci (n-1),
    fibonacci (n-2))
```

Running time $\in \Theta(\phi^n)$
 $\phi = 1.618\dots$

UVa CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 5

Sequence Alignment

- Input size = $n = |U| + |V|$

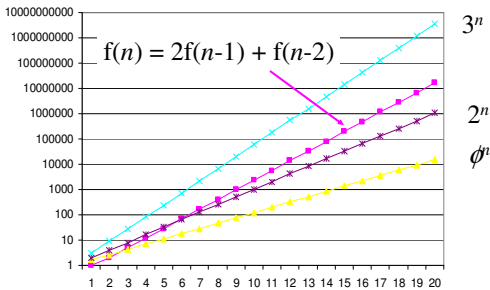
```
bestAlignment (U, V) =
  base case if |U| == 0 or |V| == 0
  otherwise
    f(bestAlignment (U[1:], V),           size = n-1
      bestAlignment (U, V[1:]),         size = n-1
      bestAlignment (U[1:], V[1:])     size = n-2)
```

$a(n) = a(n-1) + a(n-1) + a(n-2)$
 $> a(n-1) + a(n-2) \in \Theta(\phi^n)$

Running time of bestAlignment $\in \Omega(\phi^n)$

UVa CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 6

Growth of Best Alignment



BLAST

Basic Local Alignment Search Tool

<http://www.ncbi.nlm.nih.gov/BLAST/>

```
GGTTGCTGGCGTTTTCATAGGCTCCGCCCTGACGAGCATCAAAAATGACGCGGTGGCGAAACCCGAC
AGGACTATAAGATACCAGGGCTTCCCGCTGGAAGCTCCCTCGTTCGACCTGCCGTTACCGGATACCTGTC
CGGCTTTCCTCTCGGAAGCGTGGCTGCTCAGCGCTGACTTACTCAGTTGGGTAGGTGGTTCGCTCGAAGC
TGGCTGTGTCCCTTCCGAGCGCTCCGCTTATCCGTAACACTATGCTCTTGAGTCAACCCGGTAAAGTAG
GACAGTCCCGGACGCTCTGGTCAATTTTCGCGAGGACCGCTTCCCTGGAGTCCGCGTCCGCTTCCGCT
ATTCCGAATCTGCACGCCCTGCTCAAGCCTTCTGCTCAAAAGCTTTCGGCGAAGCGGCTATTATCGCGC
GCATGGCGCCGACGCGCTGGGCTGGGCTCCGACGCGAGGCTGGATGGCTTCCCATATGATTTCTTCGC
TTCGCGGCCCGCGTTCAGGCCATGCTGTCAGGAGGATGATGACGACATCAGGACAGCTTCAACGGCTCT
TACCAGCTAACTTCGATCACTGGACCGCTGATGCTCACGGGATTTATGCCACATGGACGCTTCTGGCGTT
TTTCCATAGGCTCCGCCCTGACGAGCATCAAAAAGTCAAGGTTGGCGAAACCCGACAGCACTATAAAGT
ACGAGGCTTTCCTCCCTGGAAGCGCTCTCTGCTCCGCTTACCGGATACCTGCGCTTCTCCCTT
CGGCTTCTCAATGCTCAGCTGTAGGTATCTCAGTTGGGTAGGCTGCTCCCTCAAGCTGACGAAACCCGCG
TTACGCCGACCGCTCCGCTTATCCGTAACACTGCTTGGTCCAAACGACTTAAAGGTTGGATGGATTG
TAGGCGCCCGCTATACCTTGTCTGCTCCCGCGGTGATGGAGCCGCGCACCTGCACTGAATGGAGCGCGC
GGCACTCGTAAACGCGCAAGAAATGGAGCAATCAATCTTCCGAGAACTGTAATGGCGAAACCAACCTTGG
CCATCGCTCCGCTTCCAGCAGCGCACGGCGCATCTCGGCGAGCTTGGTCTCT
```

Dinosaur DNA from Jurassic Park (p. 103)

NCBI formatting BLAST

Nucleotide Protein Translations Retrieve results for an RIB

Your request has been successfully submitted and put into the Blast Queue.

Query = (1200 letters)

The request ID is 1138550623-16164-102585974167.BLASTQ4

Format or **Reset**

The results are estimated to be ready in 10 seconds but may be done sooner.

Please press "FORMAT" when you wish to check your results. You may change the formatting options for your "FORMAT" again. You may also request results of a different search by entering any other valid request ID to use

NCBI results of BLAST

BLASTN 2.2.13 [Nov-27-2005]

Reference:
Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jihong Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1990). "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

RID: 1138550623-16164-102585974167.BLASTQ4

Database: All GenBank+EMBL+DDBJ+PDB sequences (1000000000, 975, 025, environmental samples on August 15, 1 or 2 2005 sequences), 3,735,524 sequences, 16,505,094,612 total letters

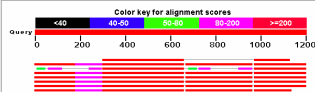
|U| = 1200
|V| > 16.5 Billion

If you have any problems or questions with the results of this search please refer to the [BLAST Page](#).
[Taxonomy reports](#)

Query: Seqopt=1200

Distribution of 495 Blast Hits on the Query Sequence

Mouse over to see the details, click to show alignments

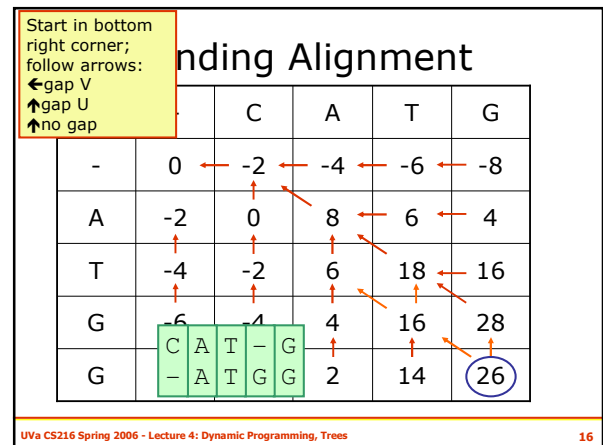
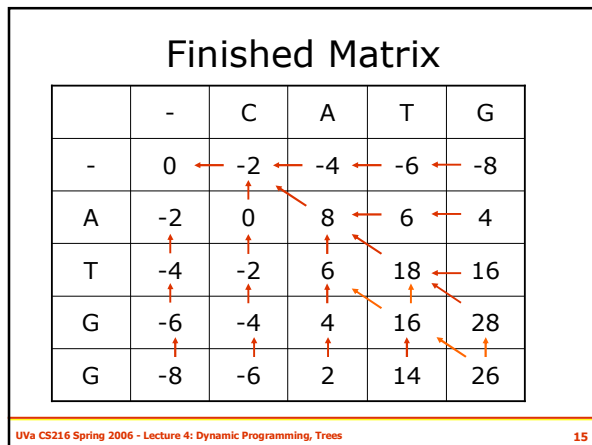
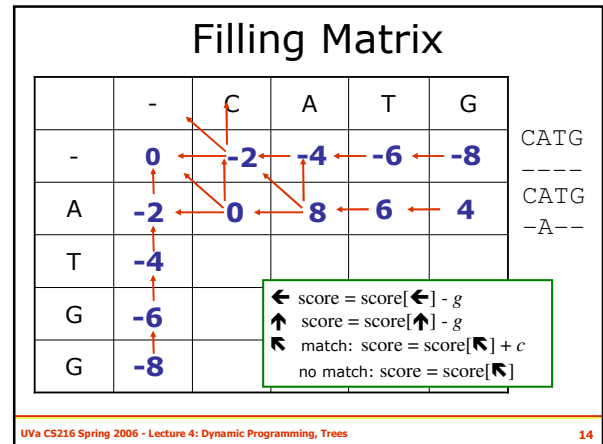
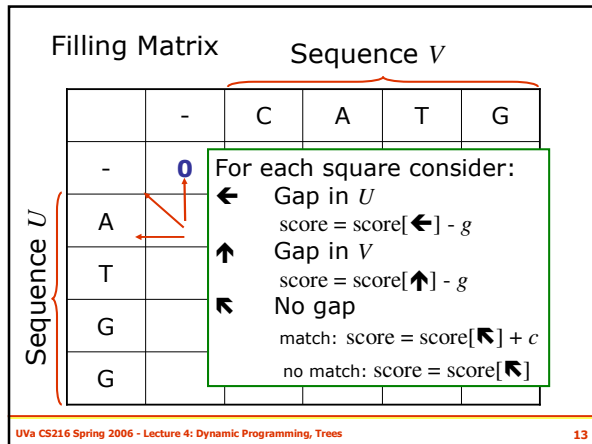


Needleman-Wunsch Algorithm

- Avoid effort duplication by storing intermediate results
- Computing a matrix showing the best possible score up to each pair of positions
- Find the alignment by following a path in that matrix

N-W Initialization Sequence V

		-	C	A	T	G
Sequence U	-	0				
	A					
	T					
	G					
	G					



- ### N-W Correctness
- Informal argument:
 - Fills each cell by picking the best possible choice
 - Finds the best possible path using the filled in matrix
 - Guaranteed to find the best possible alignment since all possibilities are considered
- UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 17

- ### N-W Analysis
- What is the space usage?
Need to store the matrix:
= $(|U| + 1) * (|V| + 1)$
- UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 18

N-W Running Time

- Time to fill matrix
 - Each square $\in O(1)$
 - Assumes:
 - Lookups are $O(1)$
 - Time scales with number of cells: $\Theta(|U| * |V|)$
- Time to find alignment
 - One decision for each entry in answer $\Theta(|U| + |V|)$
- Total running time $\in \Theta(|U| * |V|)$

\leftarrow score = score[\leftarrow] - g
 \uparrow score = score[\uparrow] - g
 \blacktriangleright match: score = score[\blacktriangleright] + c
 no match: score = score[\blacktriangleright]

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 19

Good enough?

$|U| = 1200$
 $|V| > 16.5 \text{ Billion}$

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 20

Heuristic Alignment

- BLAST needs to be faster
- No asymptotically faster algorithm is guaranteed to find best alignment
 - Only way to do better is reject some alignments without considering them
- BLAST uses heuristics:
 - Looks for short sequences (~ 3 proteins = 9 nucleotides) that match well without gaps
 - Extend those sequences into longer sequences

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 21

PS2

- Part 1 (1-3): List representations
- Part 2 (4-5): Dynamic programming implementation of sequence alignment
- Part 3 (6-10): Brute force implementation of phylogeny
 - Like alignment, brute force doesn't scale

PS2 is longer and harder than PS1.
You have 10 days to complete it -
Get started early!

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 22

Data Structures

- If we have a good list implementation, do we need any other data structures?
- For computing: **no**
 - We can compute everything with just lists (actually even less). The underlying machine memory can be thought of as a list.
- For *thinking*: **yes**
 - Lists are a very limited way of thinking about problems.

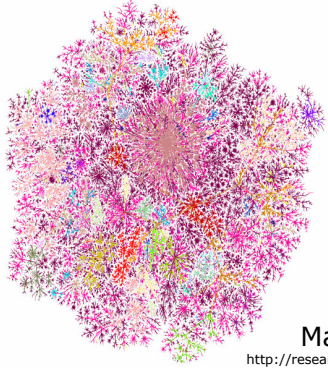
UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 23

List Limitations

Node Node Node

In a list, every element has direct relationships with only two things:
predecessor and successor

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 24



Complex Relationships

Bill Cheswick's
Map of the Internet
<http://research.lumeta.com/ches/map/gallery/>

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 25

List → Tree

- List: each element has relationships with up to **2** other elements:

Predecessor

Element

Successor
- Binary Tree: each element has relationships with up to **3** other elements:

Parent

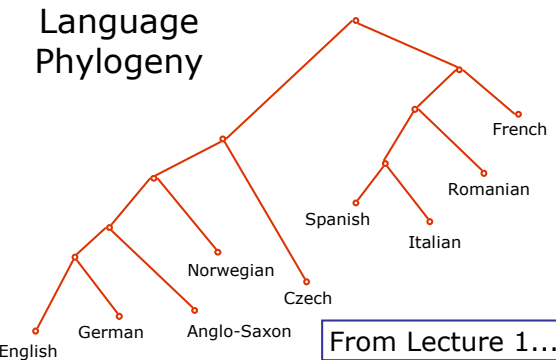
Element

Left Child

Right Child

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 26

Language Phylogeny

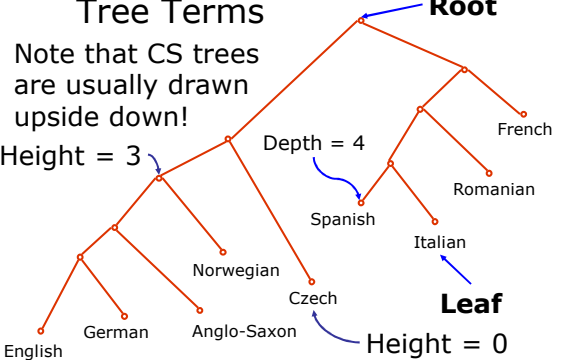


From Lecture 1...

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 27

Tree Terms

Note that CS trees are usually drawn upside down!



UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 28

Tree Terms

- Root:** a node with no parent
– There can only be one root
- Leaf:** a node with no children
- Height of a Node:** length of the longest path from that node to a leaf
- Depth of a Node:** length of the path from the Root to that node
- Height of a Tree:** maximum depth of a node in that tree = height of the root

UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 29

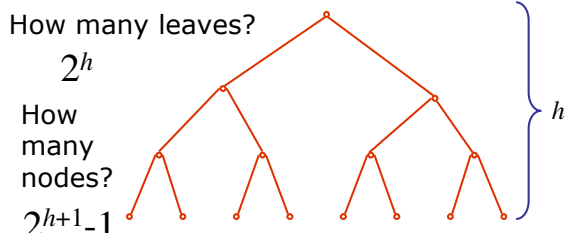
Perfect Binary Tree

How many leaves?
 2^h

How many nodes?
 $2^{h+1} - 1$

All leaves have the same depth

Nodes = $1 + 2 + 2^2 + \dots + 2^h$



UvA CS216 Spring 2006 - Lecture 4: Dynamic Programming, Trees 30

Tree Node Operations

- getLeftChild (Node)
- getRightChild (Node)
- getInfo (Node)

```
def isLeaf (self):  
    return self.getLeftChild () == None  
        and self.getRightChild () == None
```

Calculating Height

```
def height (self):  
    if self.isLeaf ():  
        return 0  
    else:  
        return 1  
            + max(self.getLeftChild().height(),  
                self.getRightChild().height())
```

Height of a Node:
length of the
longest path from
that node to a
leaf

Analyzing Height

- What is the asymptotic running time or our height procedure?

```
def height (self):  
    if self.isLeaf ():  
        return 0  
    else:  
        return 1  
            + max(self.getLeftChild().height(),  
                self.getRightChild().height())
```

Charge

- Today and tomorrow:
 - Sections in Thorton D classrooms
- Wednesday:
 - Instead of CS216, go to Ron Rivest's talk
 - 10:30am, Newcomb Hall South Meeting Room
- Get started on PS2!