

How Do We Compare Biological Sequences?

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- **Penalizing Insertions and Deletions in Sequence Alignment**
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment



Activate
Go to Settings

0:06

0:11 / 5:27



Naive Scoring of Indels

We previously assigned a fixed penalty σ to each indel

However, this fixed penalty may be too severe for a series of 100 consecutive indels

A series of k indels often represents a single evolutionary event (**gap**) rather than k events:

two gaps (assign lower score)	GATCCAG	GATCCAG	a single gap (assign higher score)
	GA-C-AG	GA--CAG	



Activate Microphone
Go to Settings

More Adequate Gap Penalties

Affine gap penalty for a gap of length k : $\sigma + \varepsilon \cdot (k-1)$

σ - the **gap opening penalty**

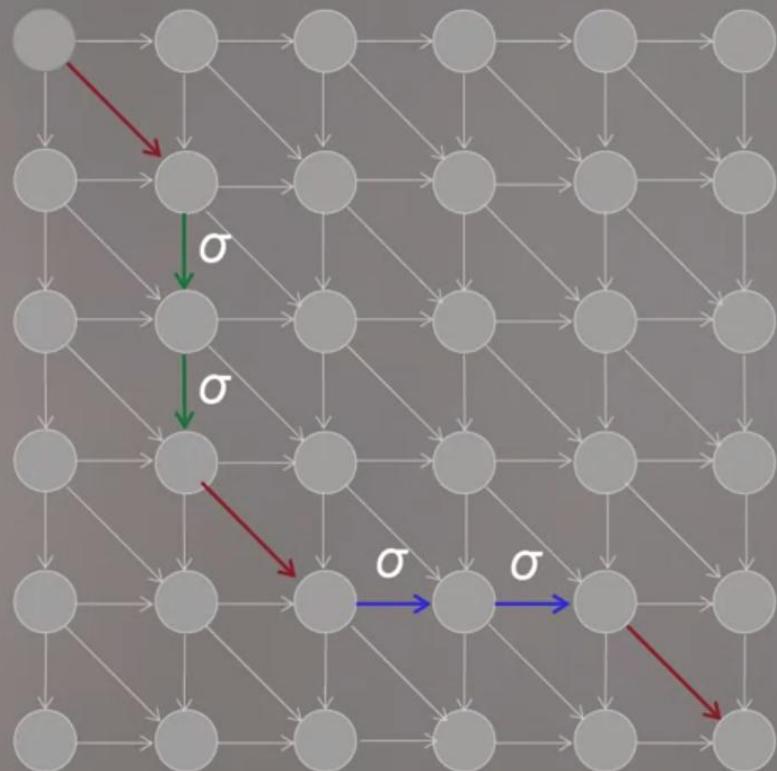
ε - the **gap extension penalty**

$\sigma > \varepsilon$, since starting a gap should be penalized more than extending it.



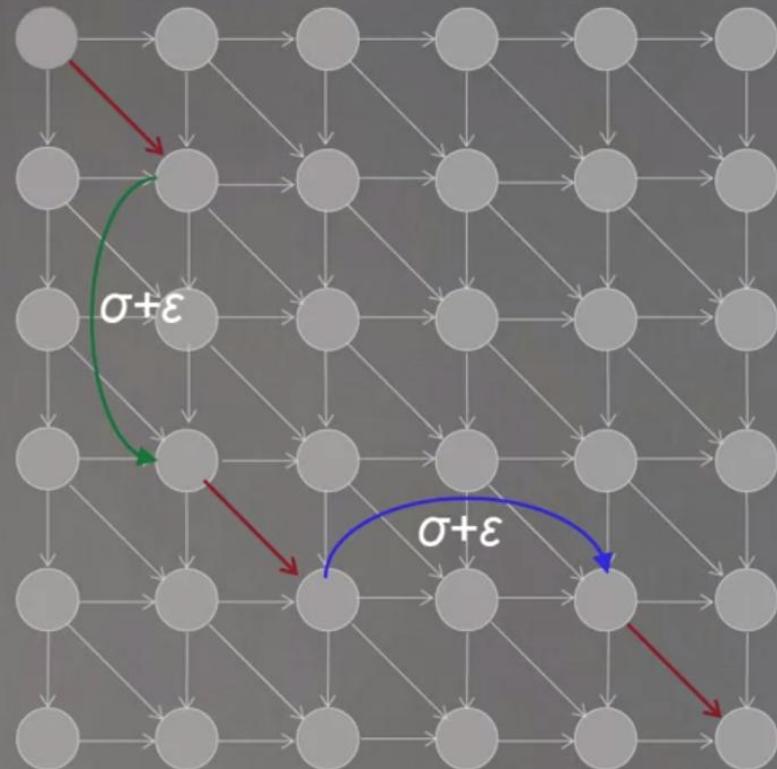
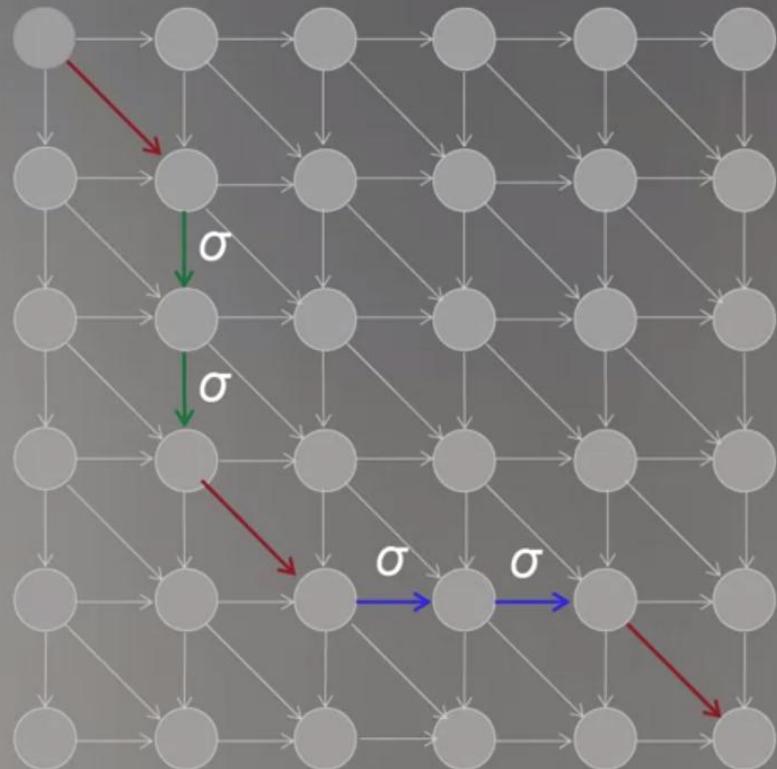
Activate Live Chat
Go to Settings

Modeling Affine Gap Penalties by Long Edges



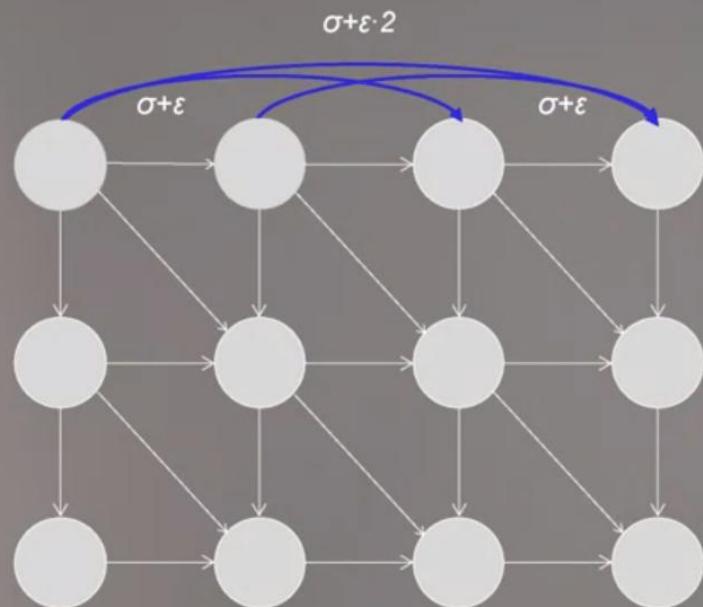
Activate microphone
Go to Settings

Modeling Affine Gap Penalties by Long Edges



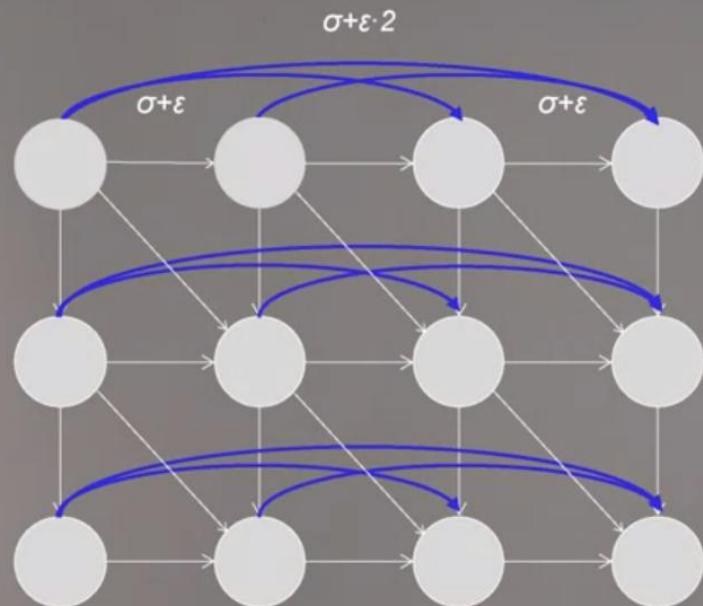
Activate Windows
Go to Settings to activate Windows.

Building Manhattan with Affine Gap Penalties



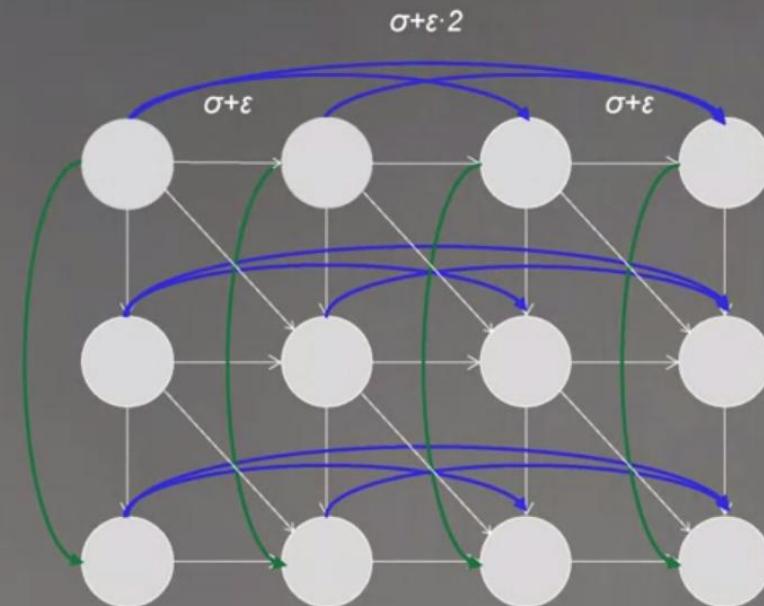
Activate []
Go to Settings []

Building Manhattan with Affine Gap Penalties



Activate
Go to Settings

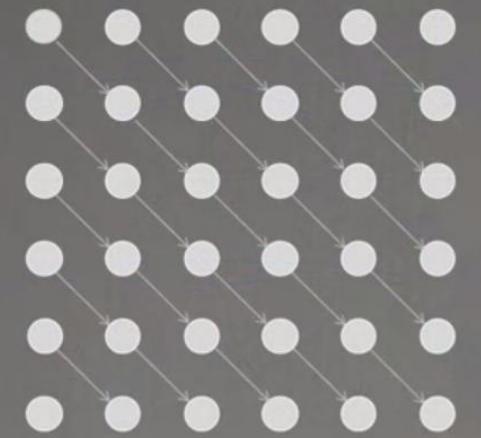
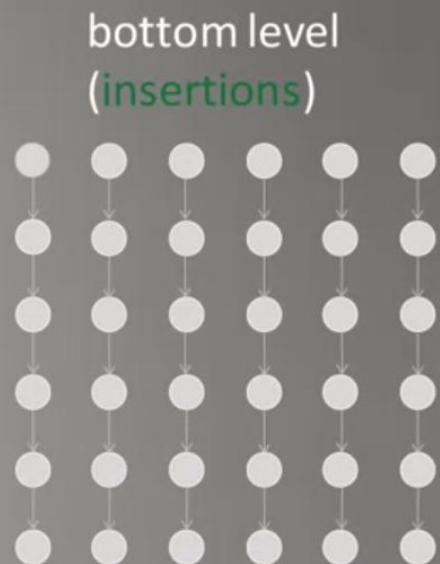
Building Manhattan with Affine Gap Penalties



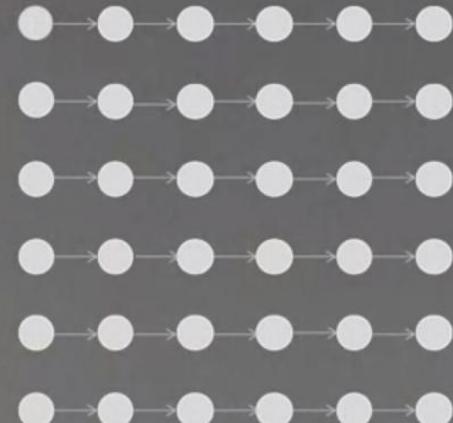
Activate Windows
Go to Settings to activate Windows.

Instead:

Building Manhattan on 3 levels

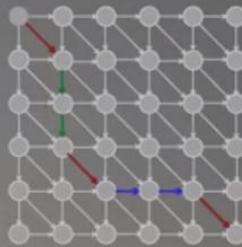


middle level
(**matches/mismatches**)

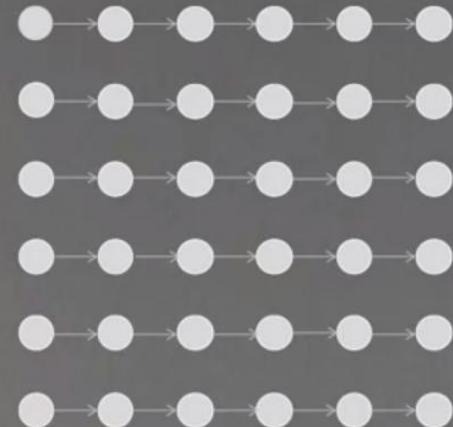
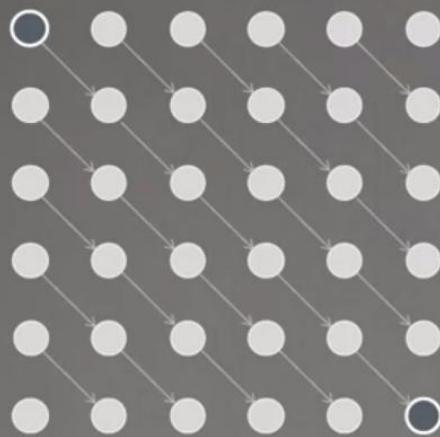
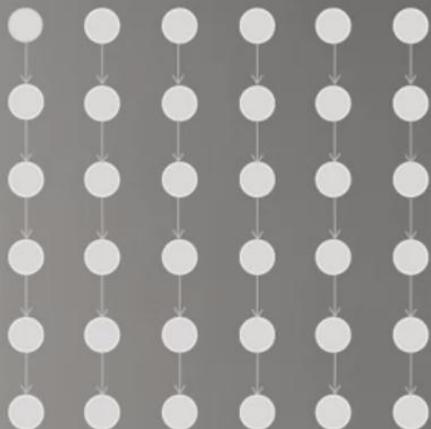


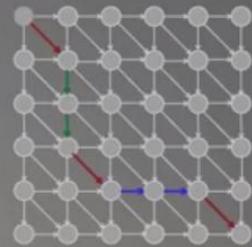
upper level
(**deletions**)

Activate Windows
Go to Settings to activate Windows.

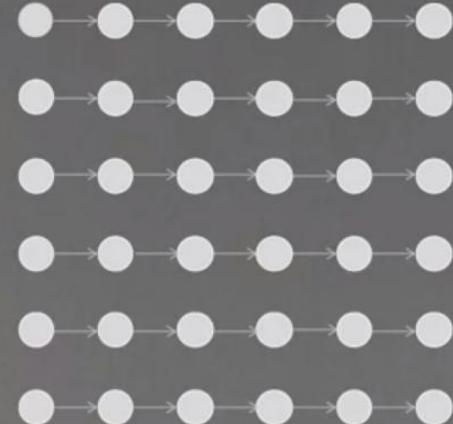
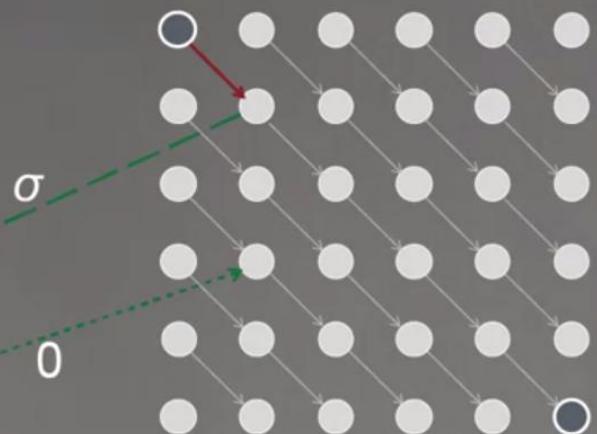
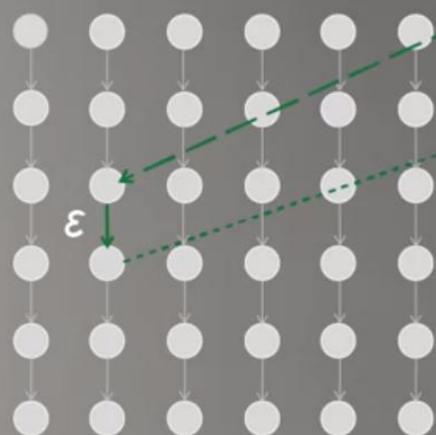


- How can we emulate this path in the 3-level Manhattan?

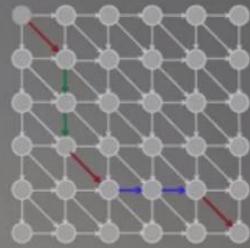




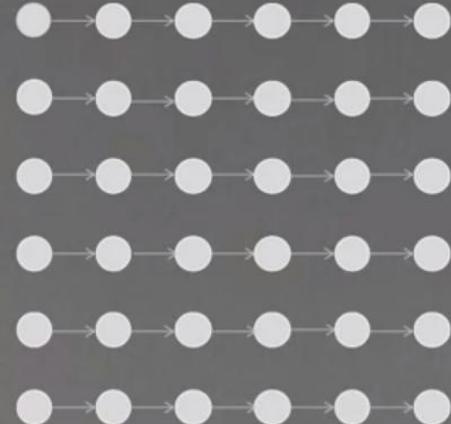
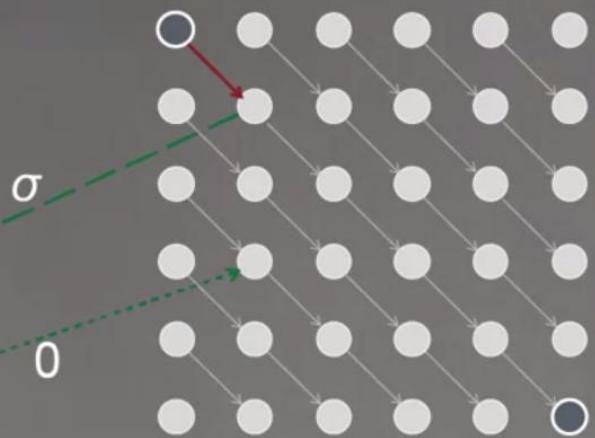
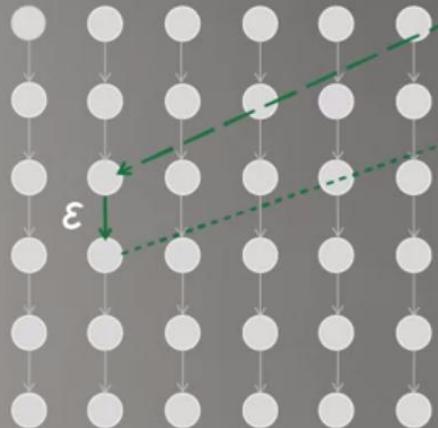
How can we emulate this path in the 3-level Manhattan?



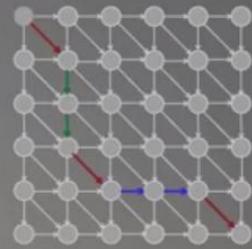
Activate Windows
Go to Settings to activate Windows.



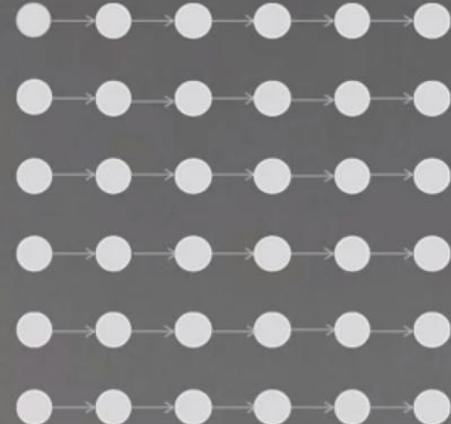
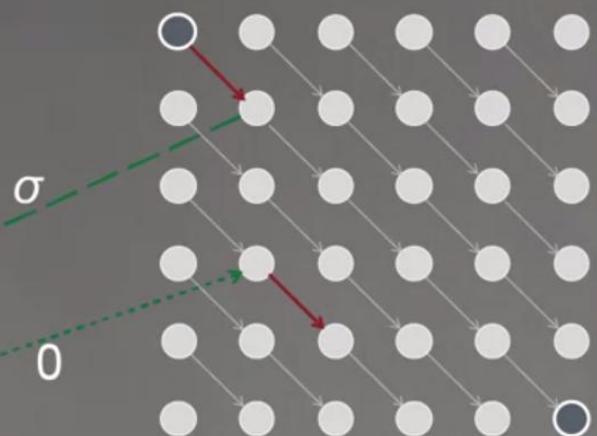
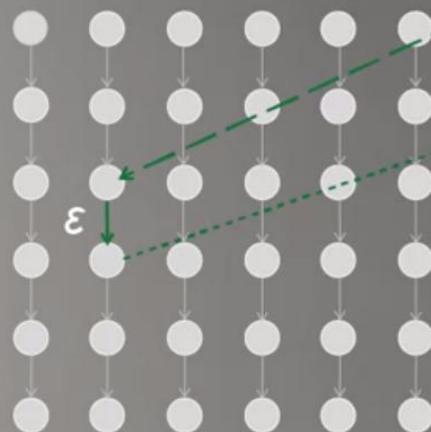
How can we emulate this path in the 3-level Manhattan?



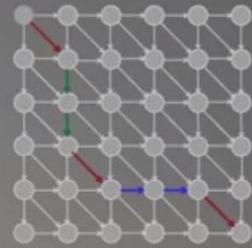
Activate Windows
Go to Settings to activate Windows.



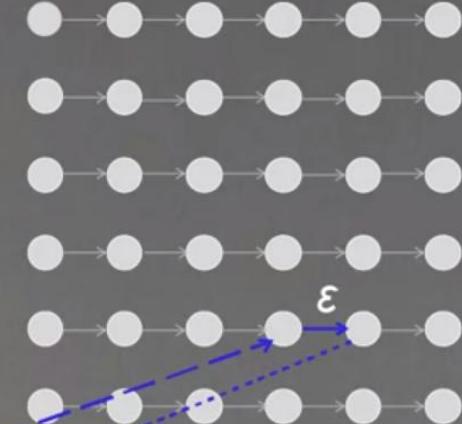
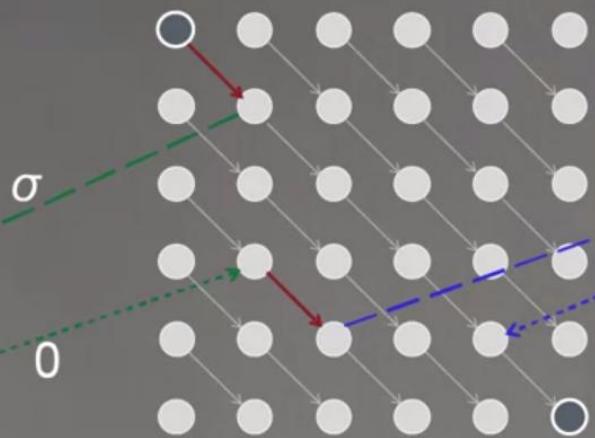
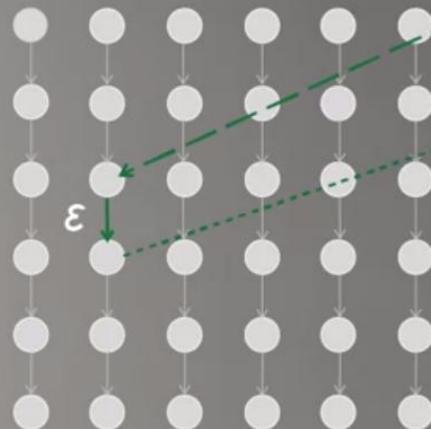
How can we emulate this path in the 3-level Manhattan?



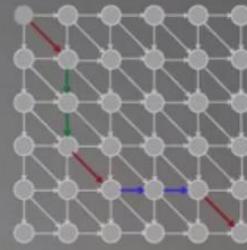
Activate Windows
Go to Settings to activate Windows.



How can we emulate this path in the 3-level Manhattan?

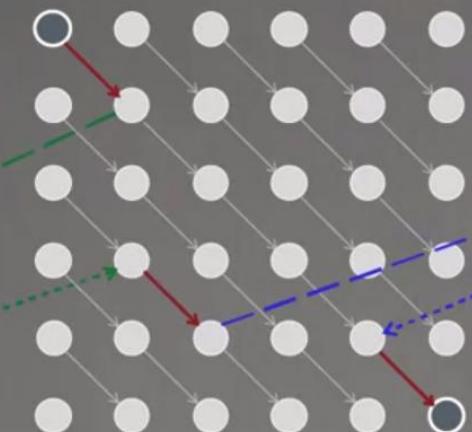
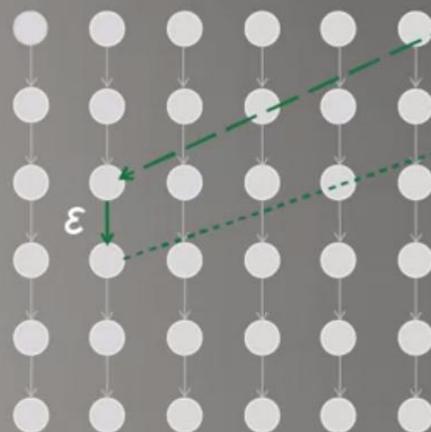


Activate Windows
Go to Settings to activate Windows.

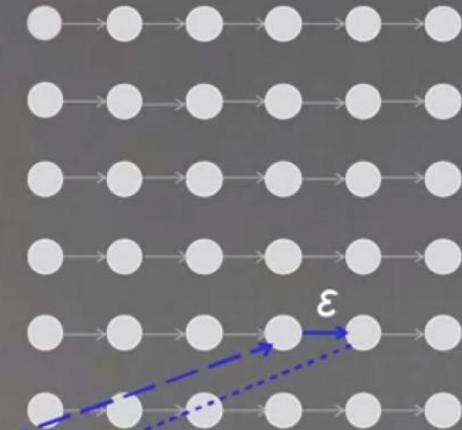


How can we emulate this path in the 3-level Manhattan?

$$lower_{i,j} = \max \left\{ lower_{i-1,j} - \varepsilon, middle_{i-1,j} - \sigma \right\}$$



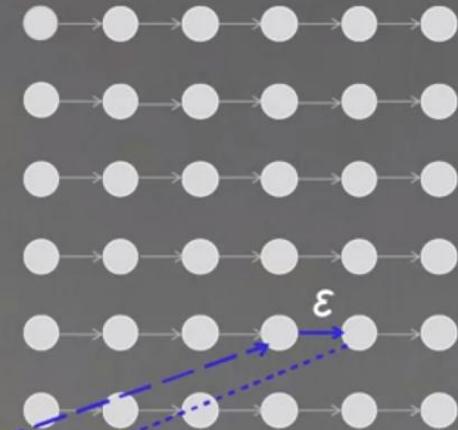
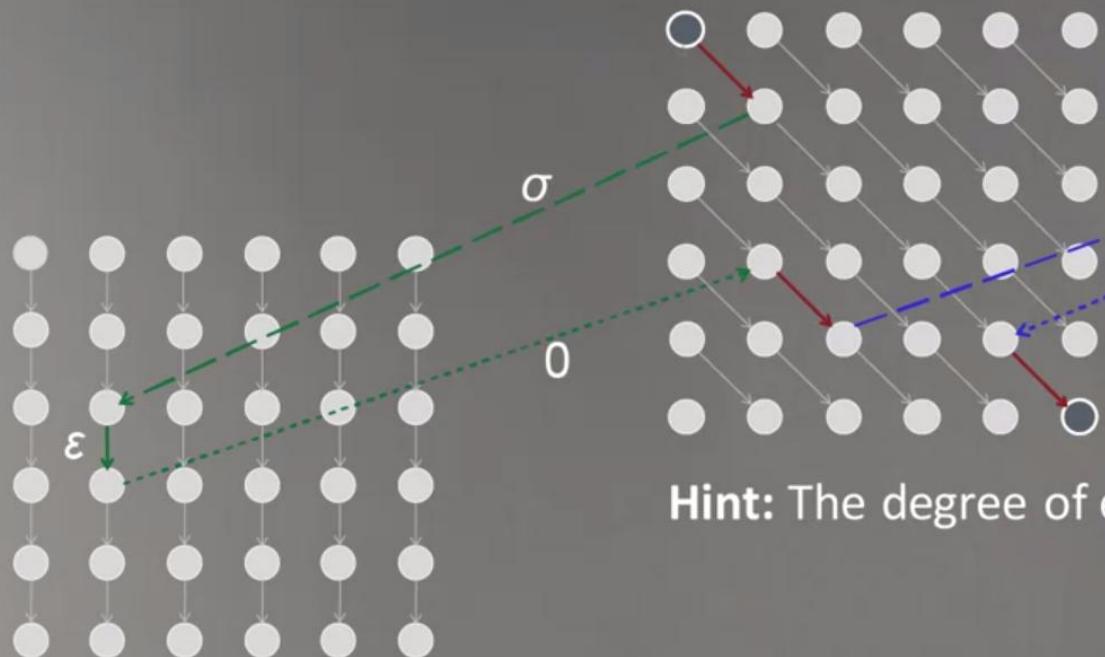
$$middle_{i,j} = \max \left\{ lower_{i,j}, middle_{i-1,j-1} + score(v_i, w_j), upper_{i,j} \right\}$$



$$upper_{i,j} = \max \left\{ upper_{i,j-1} - \varepsilon, middle_{i,j-1} - \sigma \right\}$$

Activate Windows
Go to Settings to activate Windows.

How many edges does the 3-level Manhattan have?



Hint: The degree of each node is small

Activate Windows
Go to Settings to activate Windows.

HOW DO WE COMPARE BIOLOGICAL SEQUENCES?

DYNAMIC PROGRAMMING AND DIVIDE AND CONQUER ALGORITHMS

part #9

PHILLIP COMPEAU AND PAVEL PEVZNER.
BIOINFORMATICS ALGORITHMS: AN ACTIVE LEARNING APPROACH
©2013 BY COMPEAU AND PEVZNER. ALL RIGHTS RESERVED

[MUSIC]

Activate Windows
Go to Settings to activate Windows.

How Do We Compare Biological Sequences?

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- **Space-Efficient Sequence Alignment**
- Multiple Sequence Alignment



implement sequence alignment in a more efficient way.

Activate Windows
Go to Settings to activate Windows

Can We Now Align NRP Synthetases from Two Different Bacteria?

The bottleneck: NRP synthetases are long ($\approx 20,000$ amino acids)

Alignment **runtime**: proportional to #edges (quadratic)

Alignment **memory**: proportional to #nodes (quadratic)

Memory is often a bottleneck when comparing long sequences



proportional to the number of edges, which is quadratic.

Activate Windows
Go to Settings to activate Windows

Can We Now Align NRP Synthetases from Two Different Bacteria?

The bottleneck: NRP synthetases are long ($\approx 20,000$ amino acids)

Alignment **runtime**: proportional to #edges (quadratic)

Alignment **memory**: proportional to #nodes (quadratic)

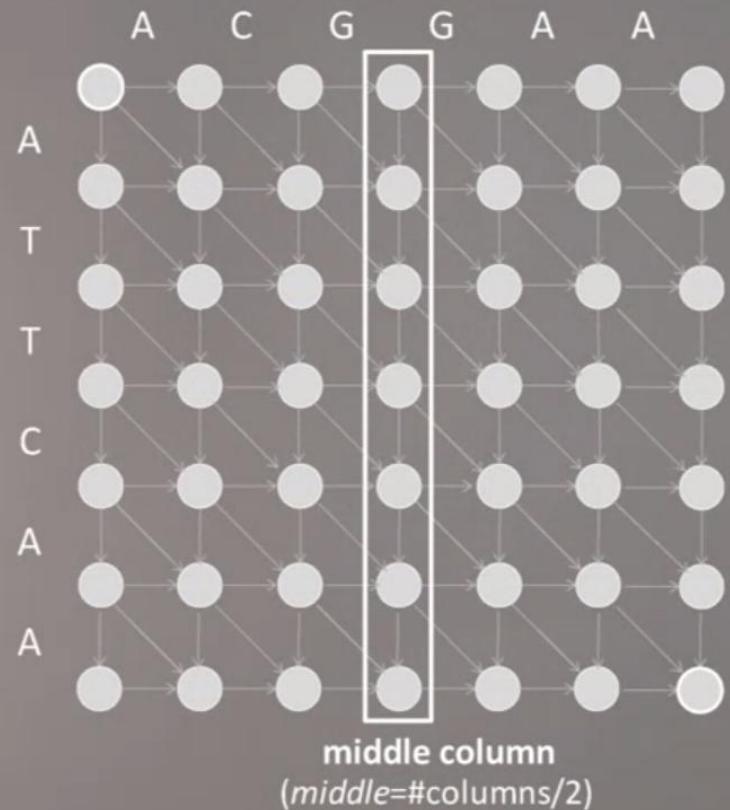
Memory is often a bottleneck
when comparing long sequences



consumption of alignment algorithm is also
proportional to

Activate Windows
Go to Settings to activate Windows

Middle Column of the Alignment



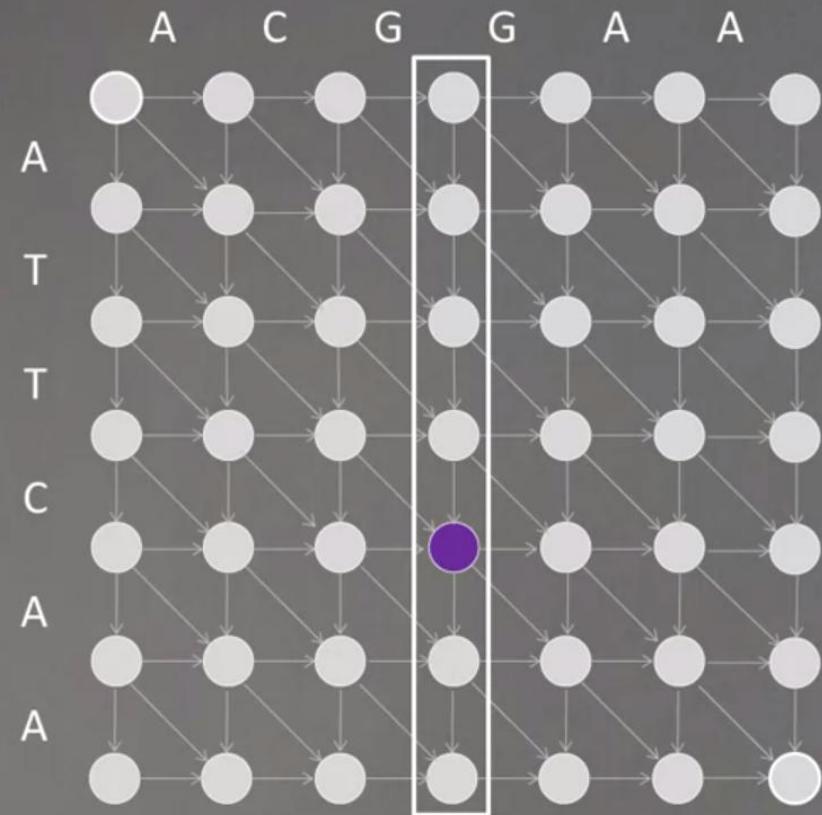
Let's define middle column of our
Manhattan grid as a column in the middle,

Activate Windows
Go to Settings to activate Windows

Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source*, *sink*)

find *MiddleNode*



Let's design algorithm to find sequence alignment.

Activate Windows
Go to Settings to activate Windows.

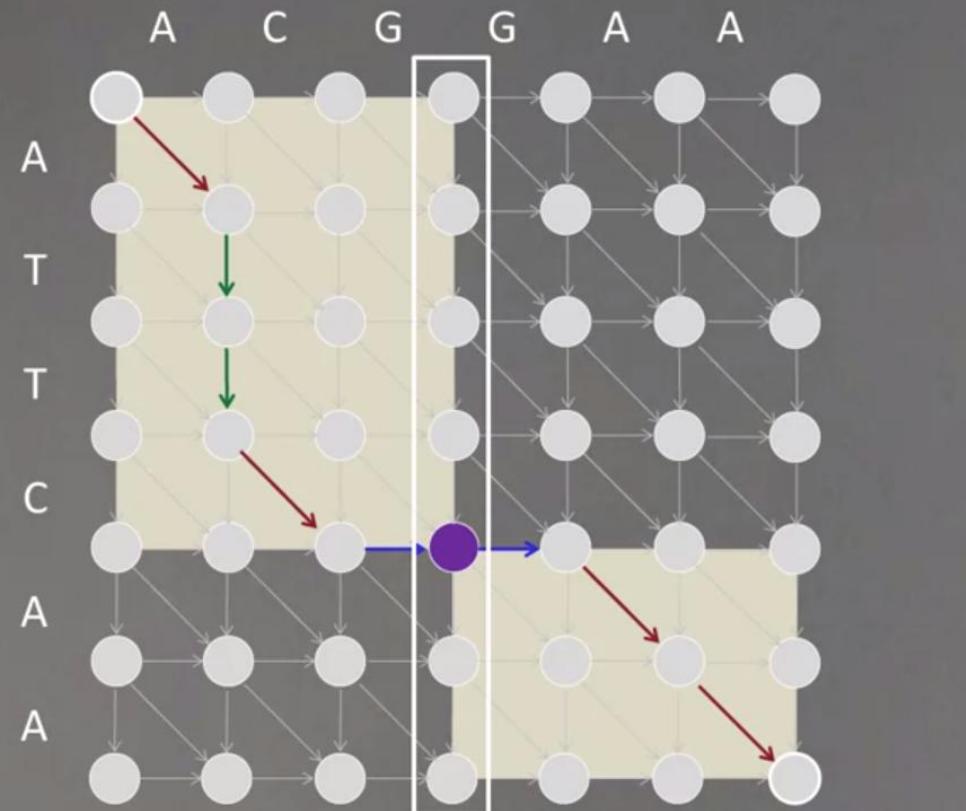
Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source*, *sink*)

 find *MiddleNode*

AlignmentPath(*source*, *MiddleNode*)

AlignmentPath(*MiddleNode*, *sink*)



since we know all parts of a path, we can
do it recursively.

Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source*, *sink*)

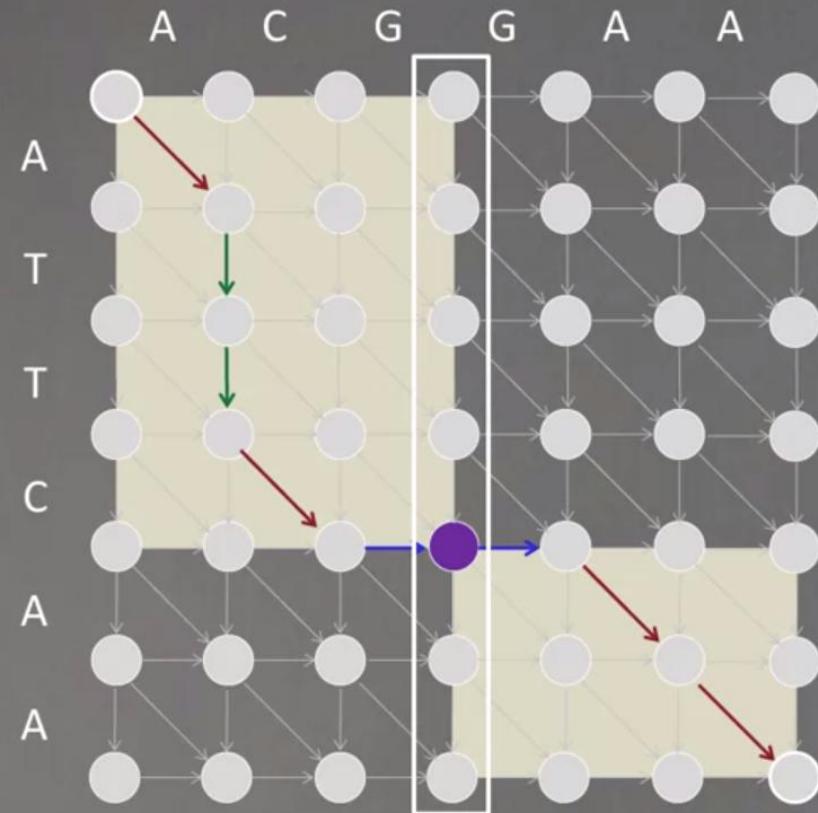
 find *MiddleNode*

AlignmentPath(*source*, *MiddleNode*)

AlignmentPath(*MiddleNode*, *sink*)

The only problem left is how to find this middle node in **linear space**!

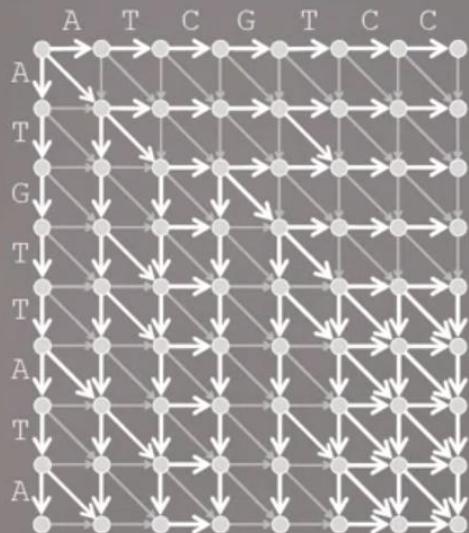
problem is how to find this middle node, because we don't know yet.



Activate Windows
Go to Settings to activate Windows.

Computing Alignment Score in Linear Space

Finding the **longest path** in the alignment graph
requires storing all backtracking pointers – $O(nm)$ memory



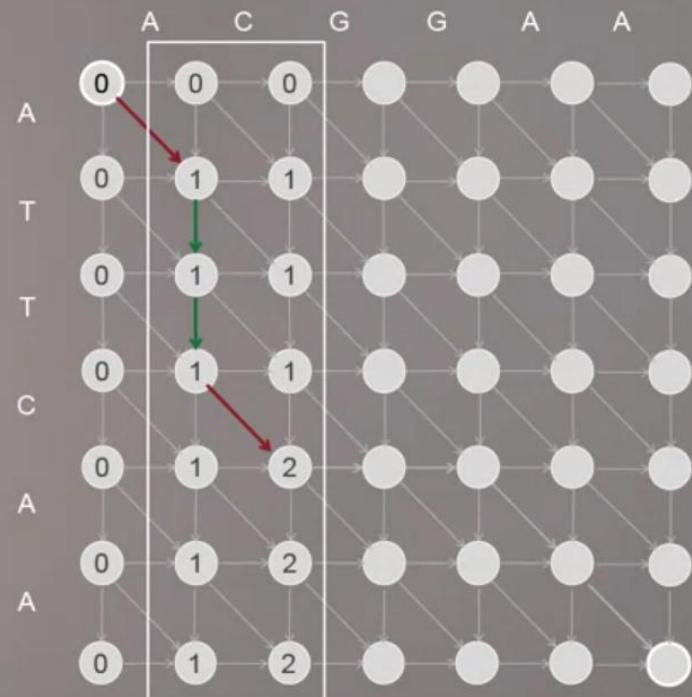
Computing the **length of the longest path** does not
require storing any backtracking pointers – $O(n)$ memory



Can we skip

Activate Windows
Go to Settings to activate Windows

Computing Alignment Score in Linear Space

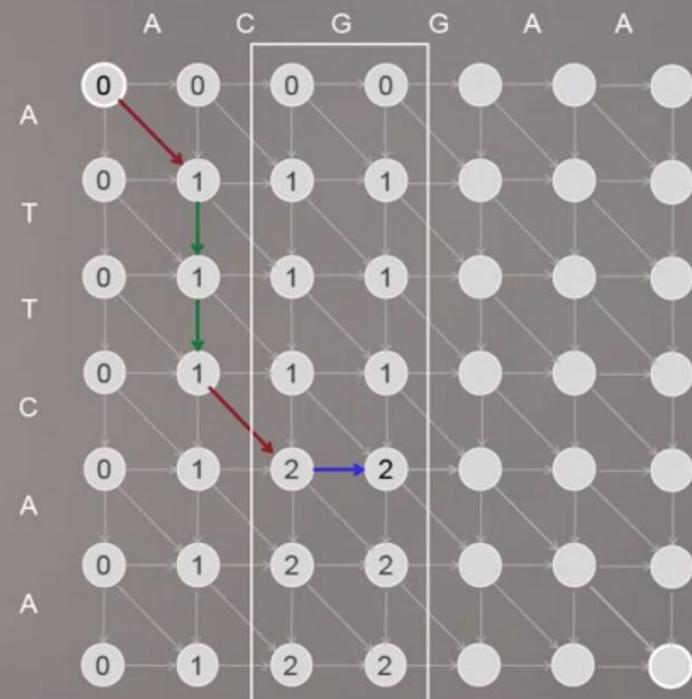


of the first column, and we can simply forget them.



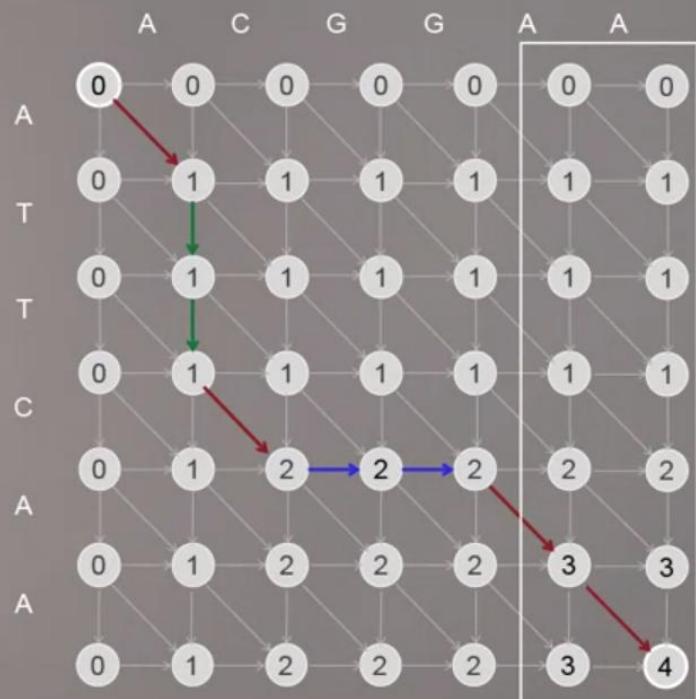
Activate Windows
Go to Settings to activate Windows

Computing Alignment Score in Linear Space



Activate Windows
Go to Settings to activate Windows

Computing Alignment Score in Linear Space



and so on and so on and so on.
And at each moment of



Computing Alignment Score in Linear Space

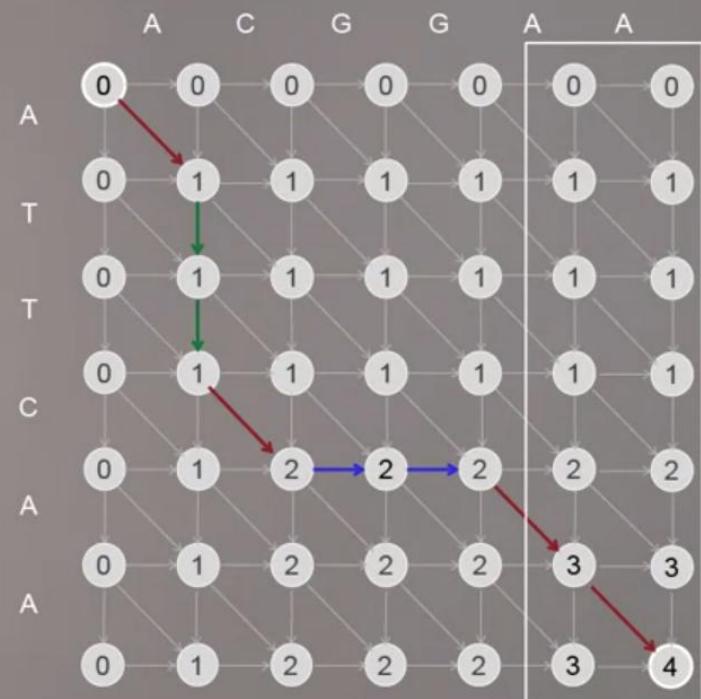


And we don't need to store all scores of all nodes in our Manhattan grid.

Activate Windows
Go to Settings to activate Windows

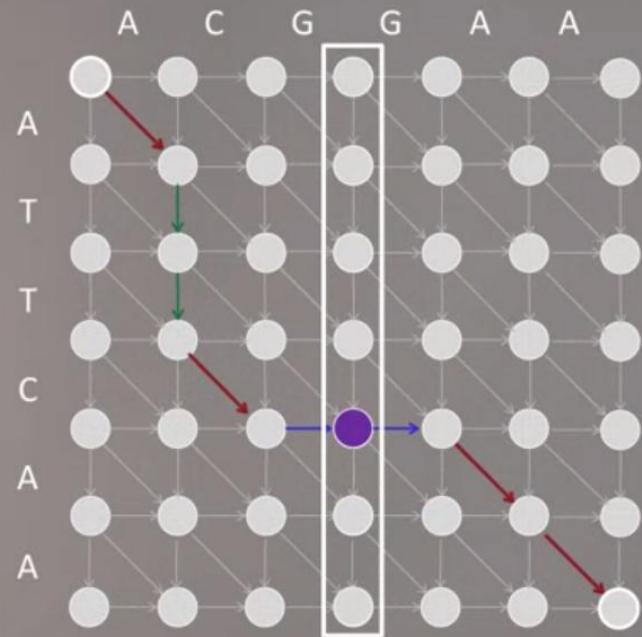
Press Esc to exit full screen

Computing Alignment Score in Linear Space



this algorithms, we need to store only two columns.

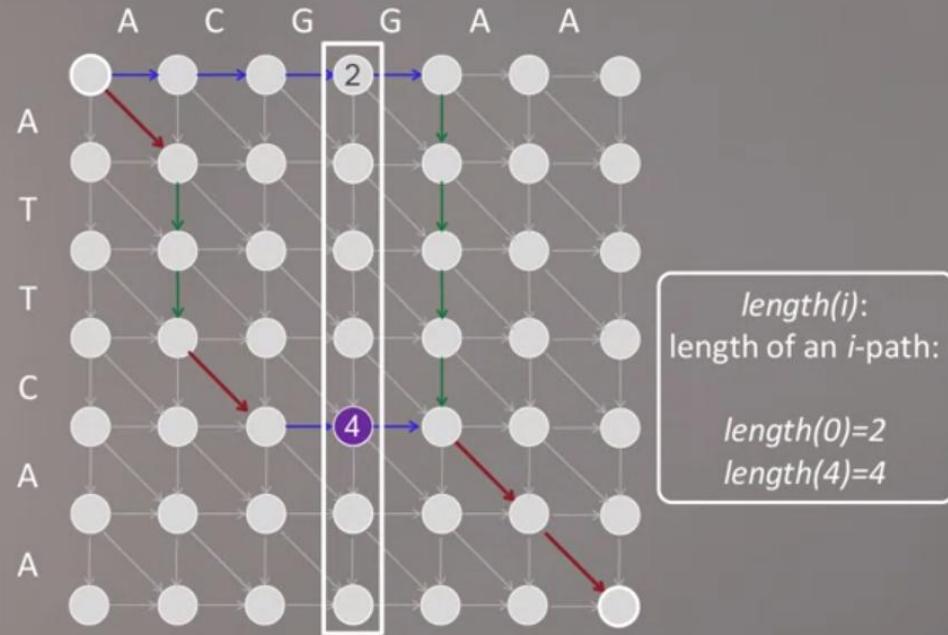
Can We Find the Middle Node without Constructing the Longest Path?



our middle column and middle node.
So, middle column was a column in the



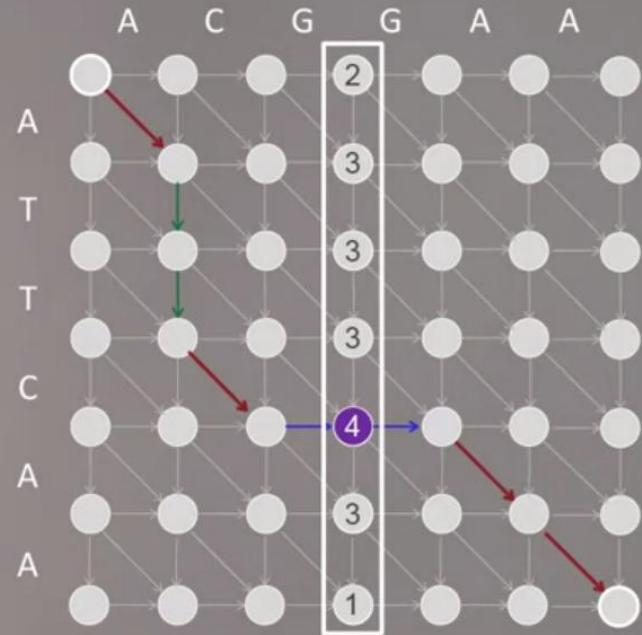
Can We Find The Lengths of All i -paths?



is two here, and lengths of four paths is
four.

Activate Windows
Go to Settings to activate Windows

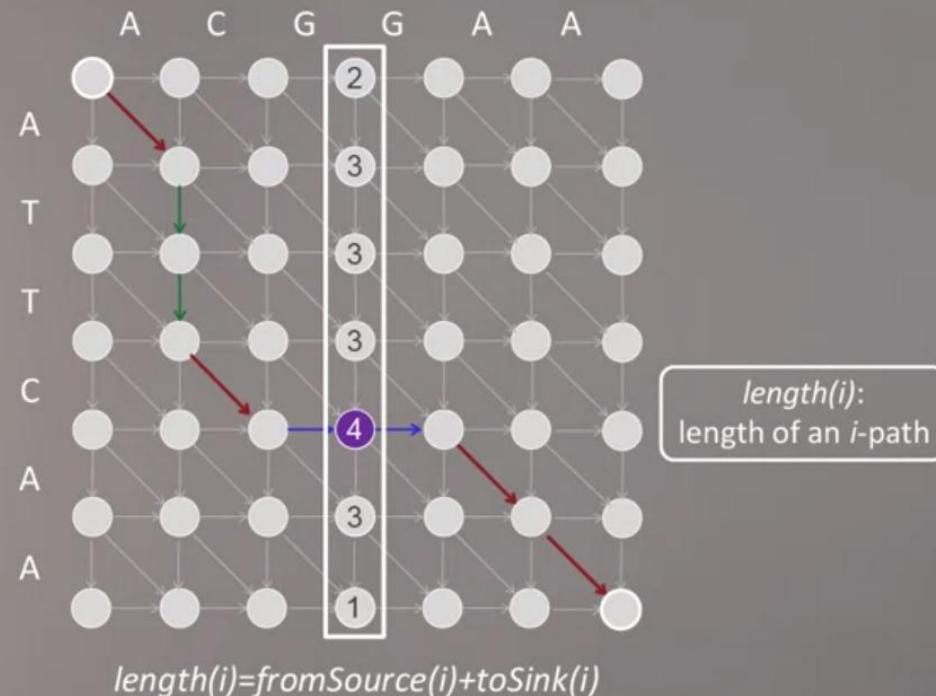
Can We Find the Lengths of All i -paths?



For example here, in this graph, it's two,
three, three, four, three, and one.

Activate Windows
Go to Settings to activate Windows

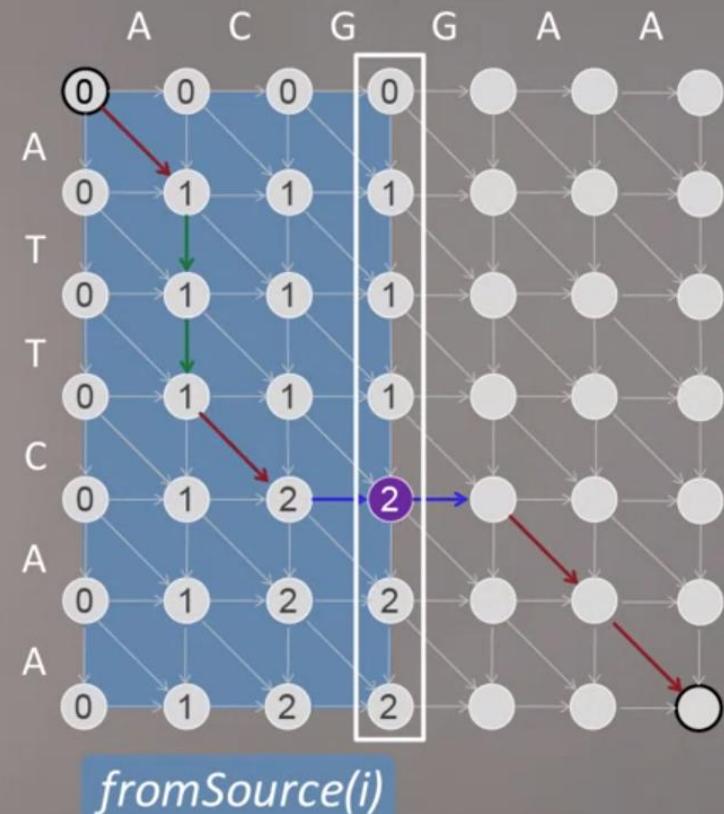
Can We Find The Lengths of i -paths?



this node, and from this node, node to the sink.

Activate Windows
Go to Settings to activate windows

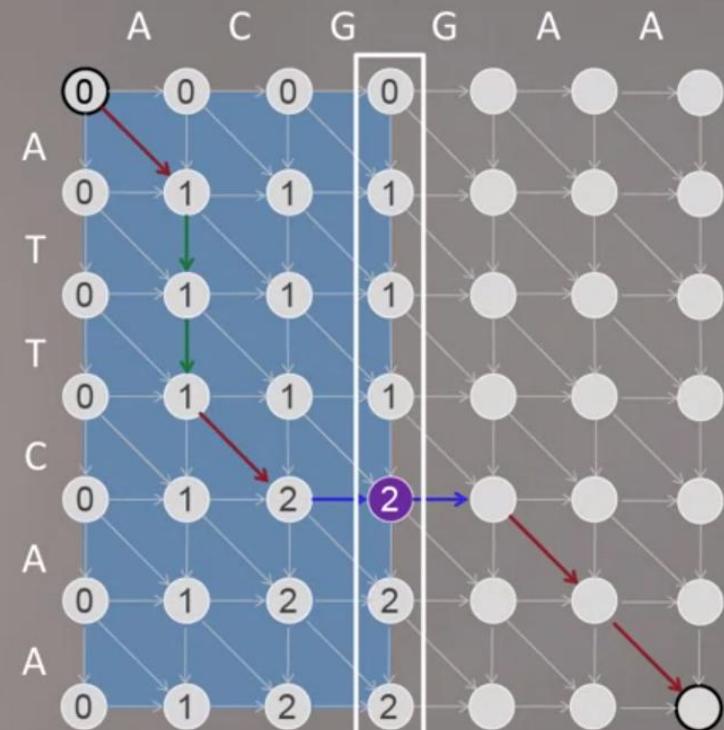
Computing *FromSource*



from source *i* just by going through the
first part of



Computing *FromSource*



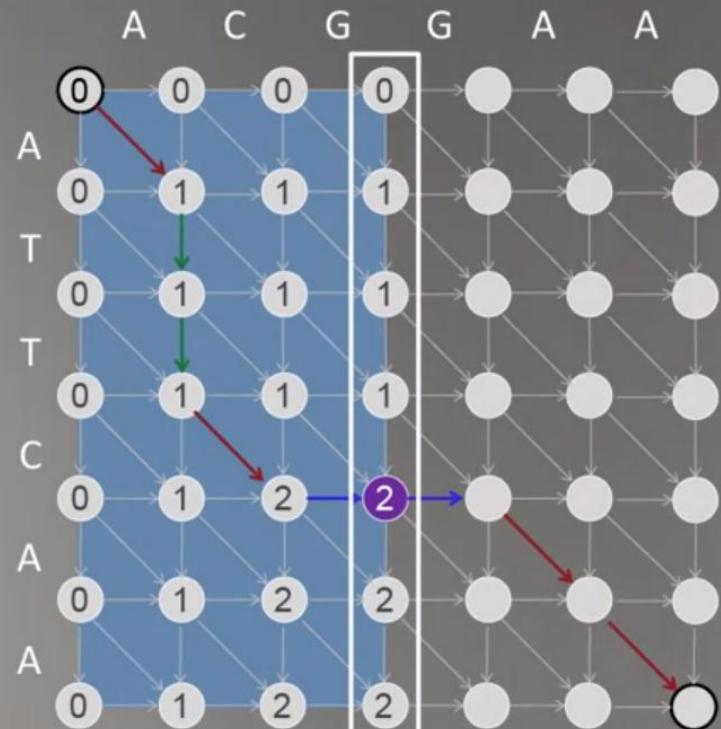
fromSource(i)

from source i just by going through the
first part of



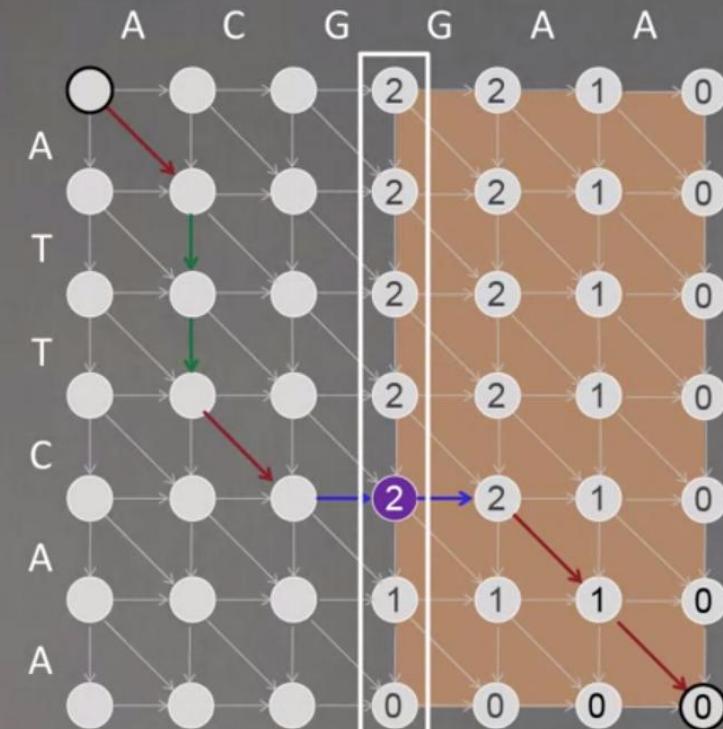
Activate Windows
Go to Settings to activate Windows 8

Computing *FromSource* and *toSink*



fromSource(i)

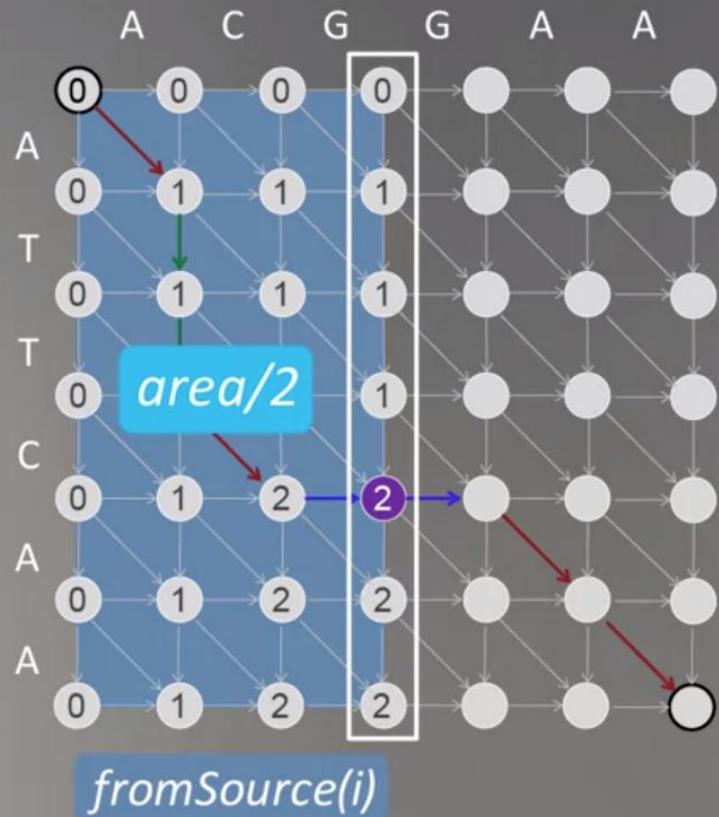
failures from sink by going in the second part



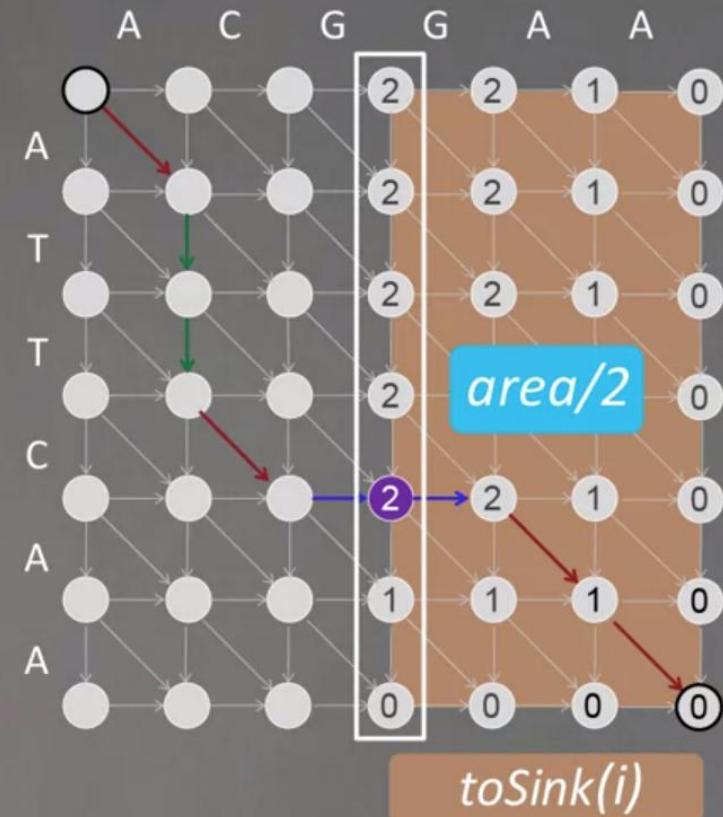
toSink(i)

Activate Windows
Go to Settings to activate Windows.

How Much Time Did It Take to Find the Middle Node ?



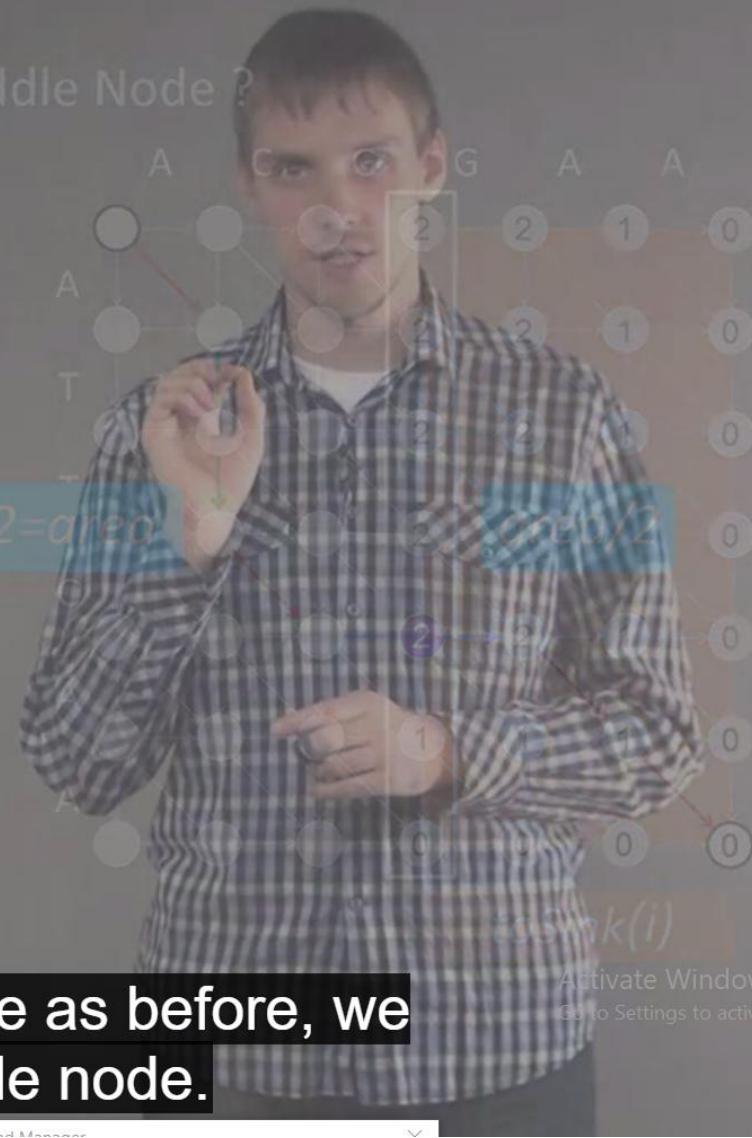
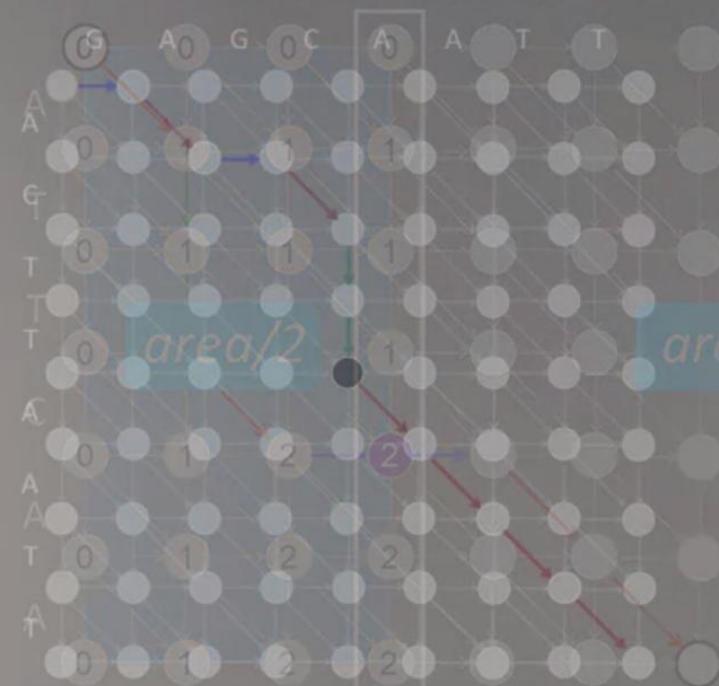
from source, we need to visit area divided
by two of our whole Manhattan grid.



Activate Windows
Go to Settings to activate Windows.

Laughable Progress: $O(n^2)$ to Find the Middle Node?

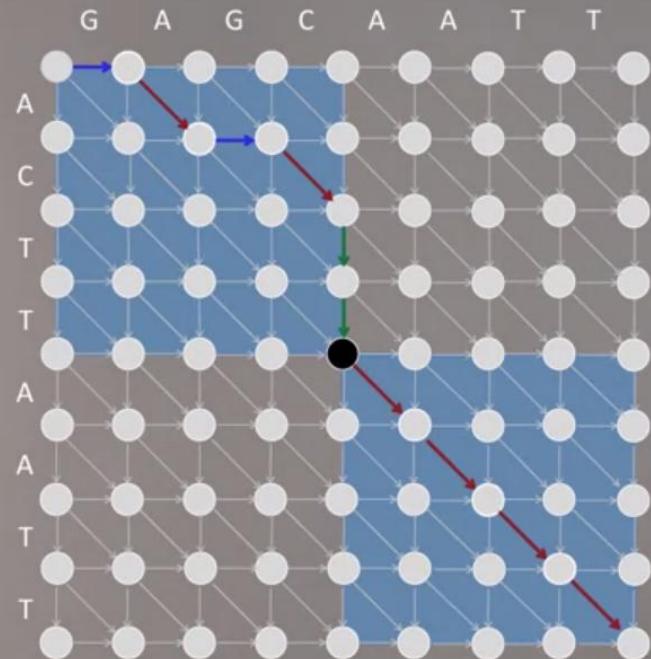
Time to Find ONE Nodes!



Using the same run time as before, we found only middle node.

Laughable Progress: $O(nm)$

Time to Find **ONE** Nodes!



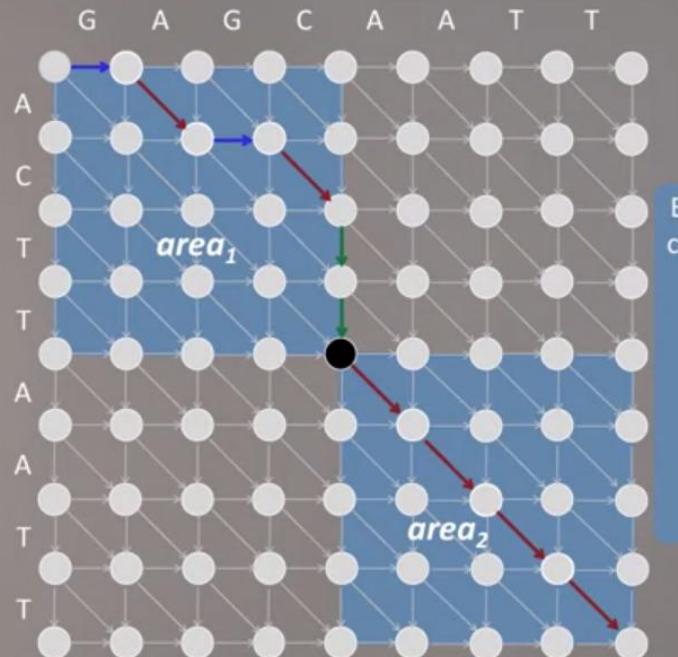
How much time would it take to conquer 2 subproblems?

here should be visited after you've found
middle node, because our



Activate Windows
Go to Settings to activate Windows

Laughable Progress: $O(nm)$
Time to Find **ONE** Nodes!



Each subproblem
can be conquered
in time
proportional to
its area:

$$area_1 + area_2 = area/2$$

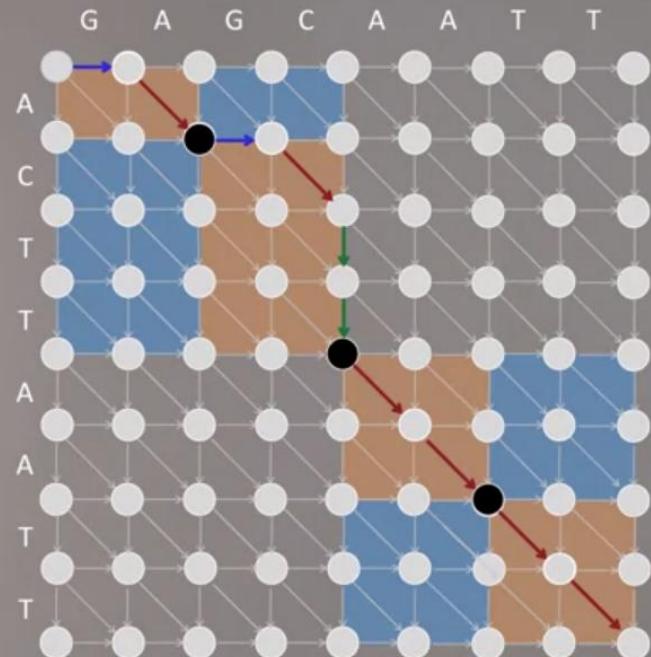
How much time would it take to conquer 2 subproblems?

grid so, we need only half of the time
which were required to find



Activate Windows
Go to Settings to activate Windows

Laughable Progress: $O(nm+nm/2)$
Time to Find **THREE** Nodes!



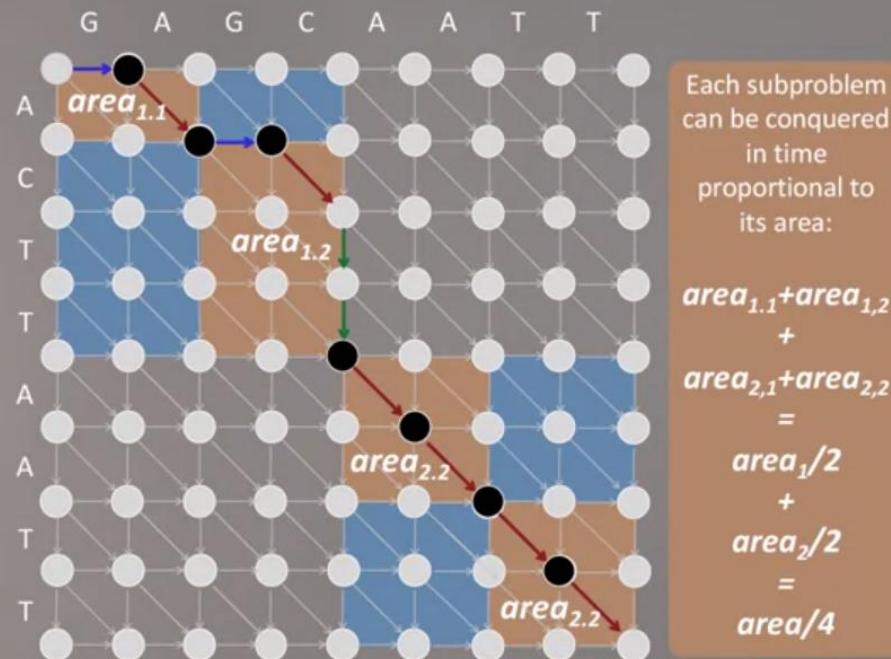
How much time would it take to conquer 4 subproblems?

step we will split our area into smaller rectangles



Activate Windows
Go to Settings to activate Windows

Laughable Progress: $O(nm+nm/2)$
Time to Find **THREE** Nodes!



How much time would it take to conquer 4 subproblems?

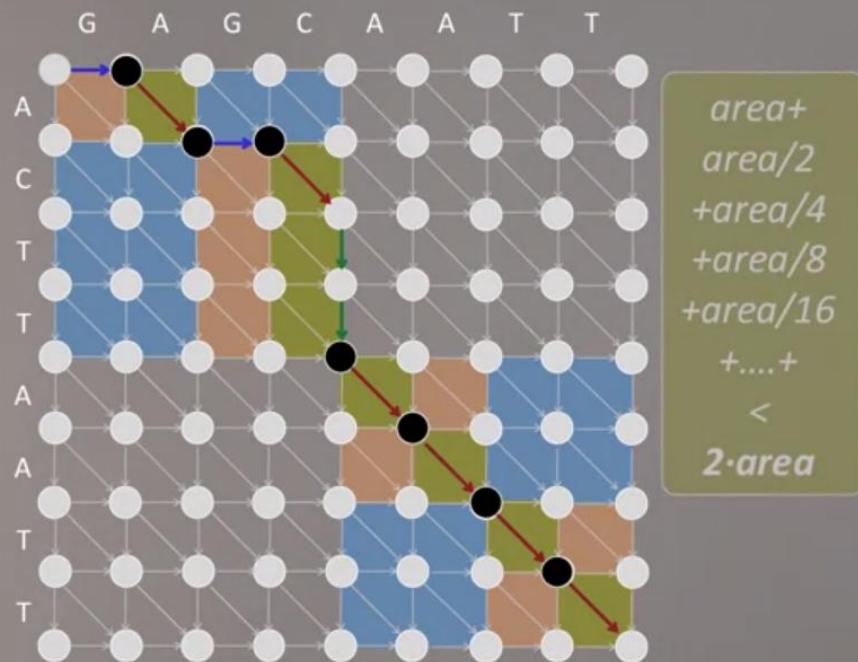
we need, to take a look on, is
proportional to the



Activate Windows
Go to Settings to activate Windows

$$O(nm + nm/2 + nm/4)$$

Time to Find **NEARLY ALL** Nodes!



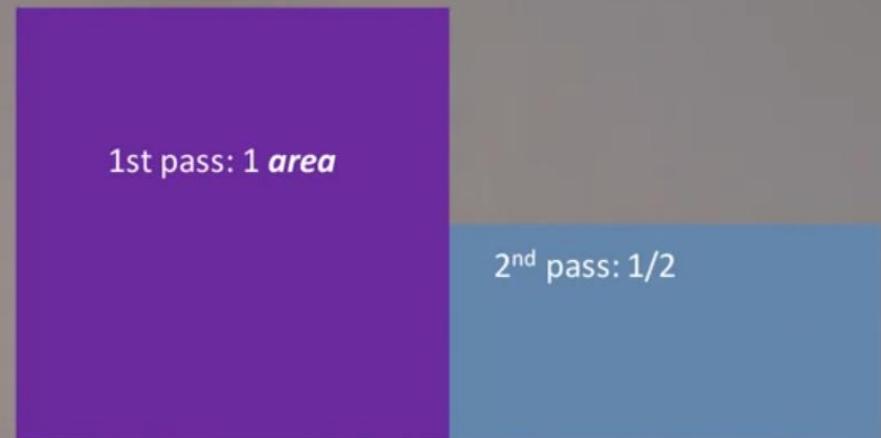
How much time would it take to conquer ALL subproblems?

And on the next step we see that area divided by four, and so on and so on.



Total Time:

$$area + area/2 + area/4 + area/8 + area/16 + \dots$$

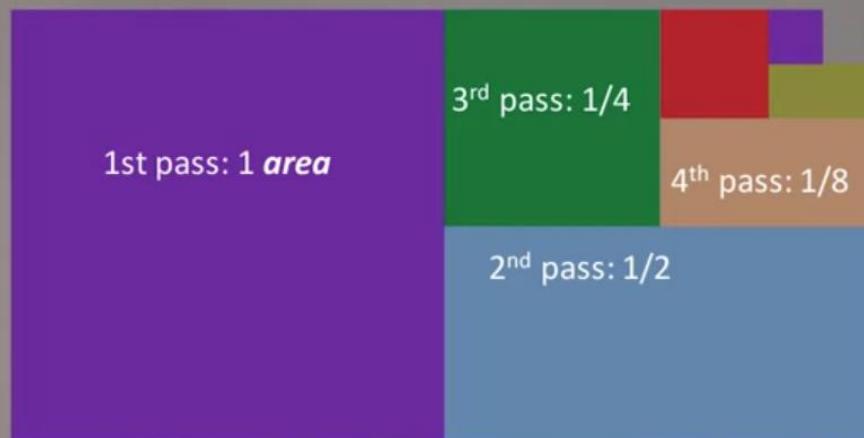


On the second pass we need to take a look
on that half of area.

Activate Windows
Go to Settings to activate Windows

Total Time:

$$area + area/2 + area/4 + area/8 + area/16 + \dots$$



It's fourth, and so on and so on.

Activate Windows
Go to Settings to activate Windows