# Sequence alignment

Correspondence between bases of two DNA sequences, or between amino acids of two protein sequences

Alignment :  2 x k matrix ( k ≥ m, n )

$V$ = ACCTGGTAAA     n = 10

$W$ = ACTGCGTATA     m = 10

8  matches
1  mismatches
1  deletions
1  insertions

| V | A | C | C | T | G | — | G | T | A | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|
| W | A | C | — | T | G | C | G | T | A | T | A |

# "Goodness" of alignments

Given two sequences, there are many possible alignments

```
ATTTTCCC

ATTTACGC
```
distance=2

```
ATTT-TCCC

ATTTA-CGC
```
distance=3

```
ATTTTCCC————————

————————ATTTACGC
```
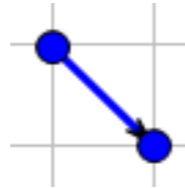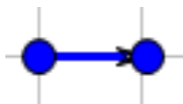distance=16

**Edit distance**: the total number of substitutions, insertions and deletions needed to transform one sequence to another

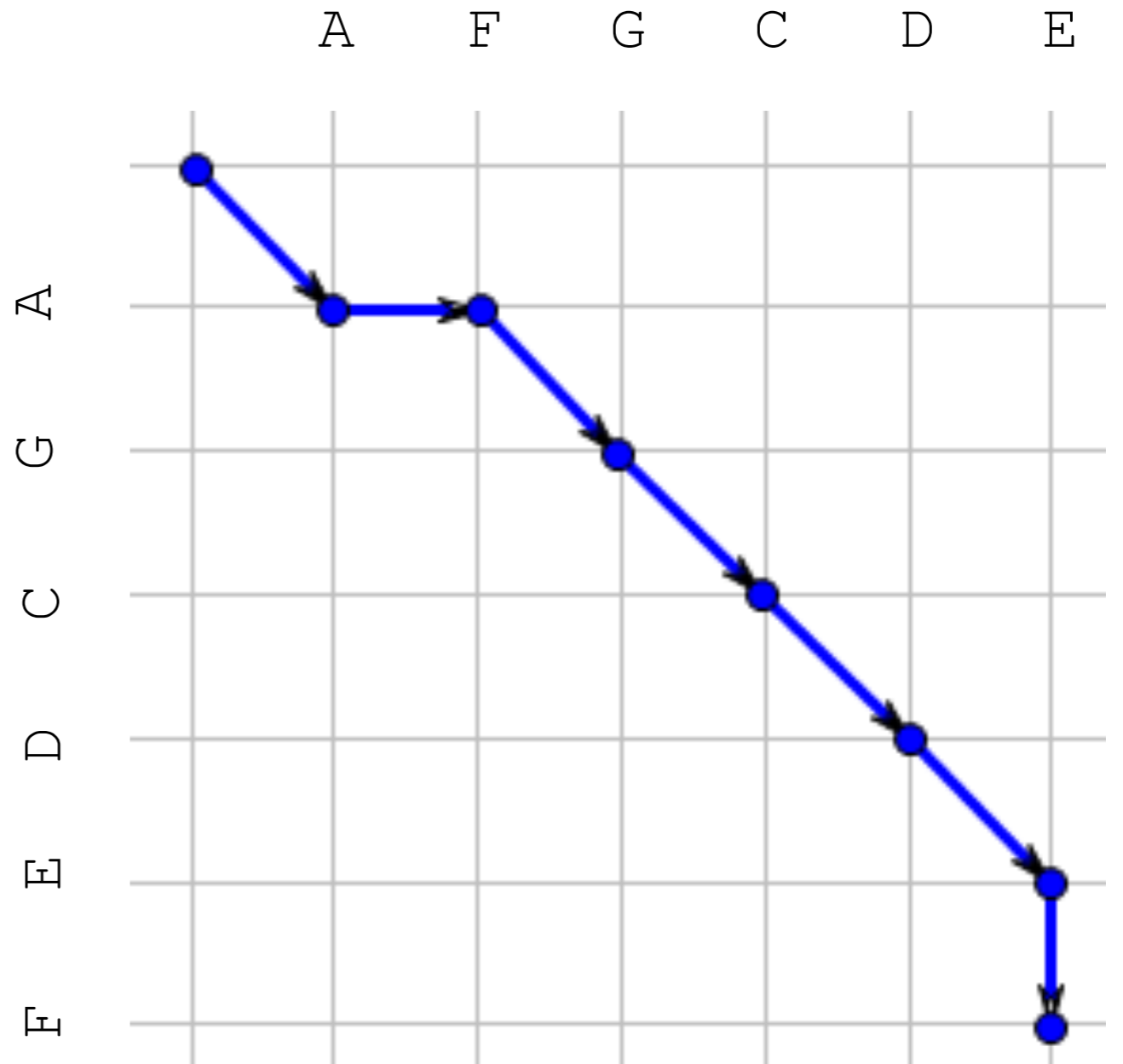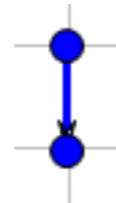# Enumeration of all possible alignments
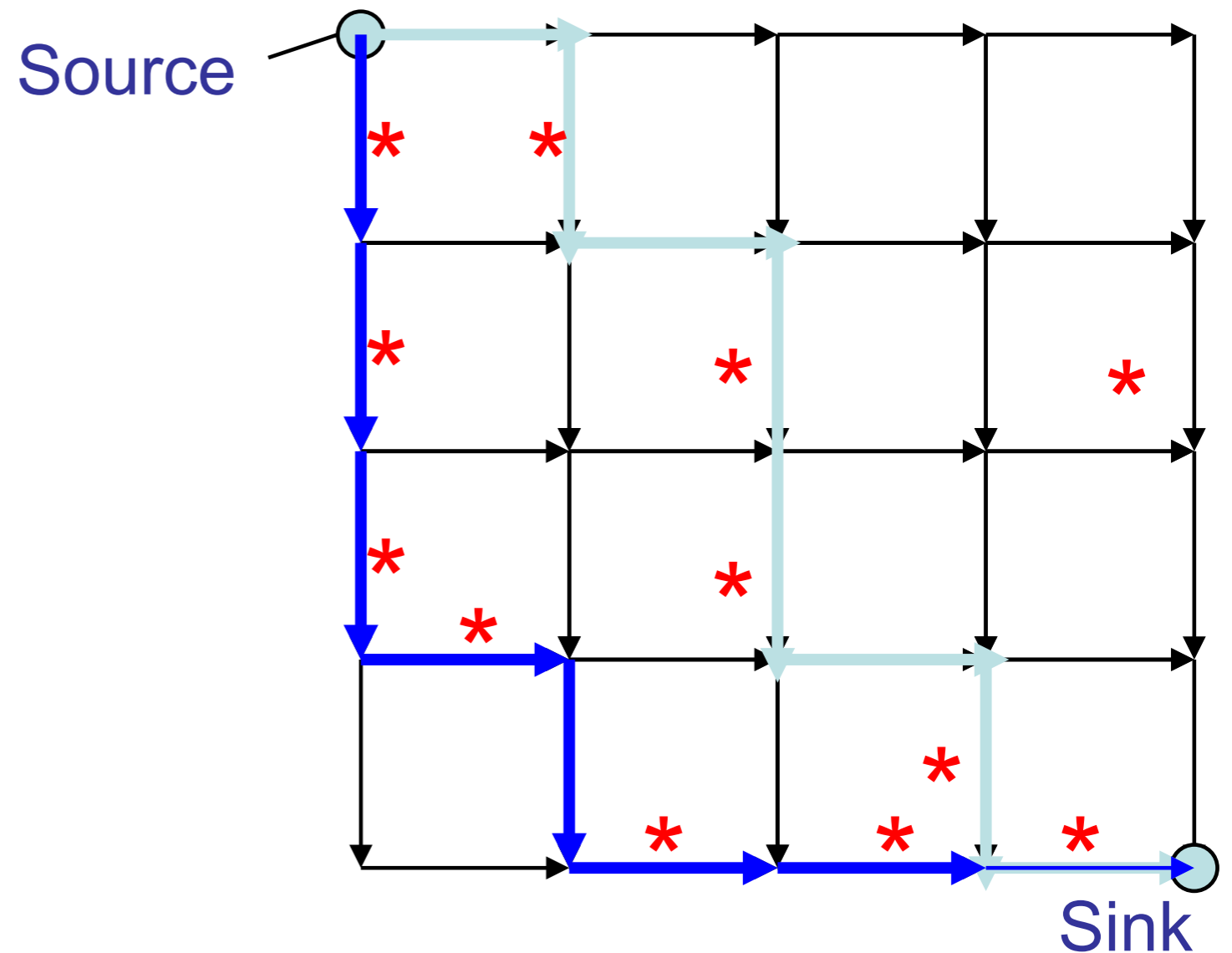
Match

Insertion_X

Insertion_Y

A    F    G    C    D    E

A

G

C
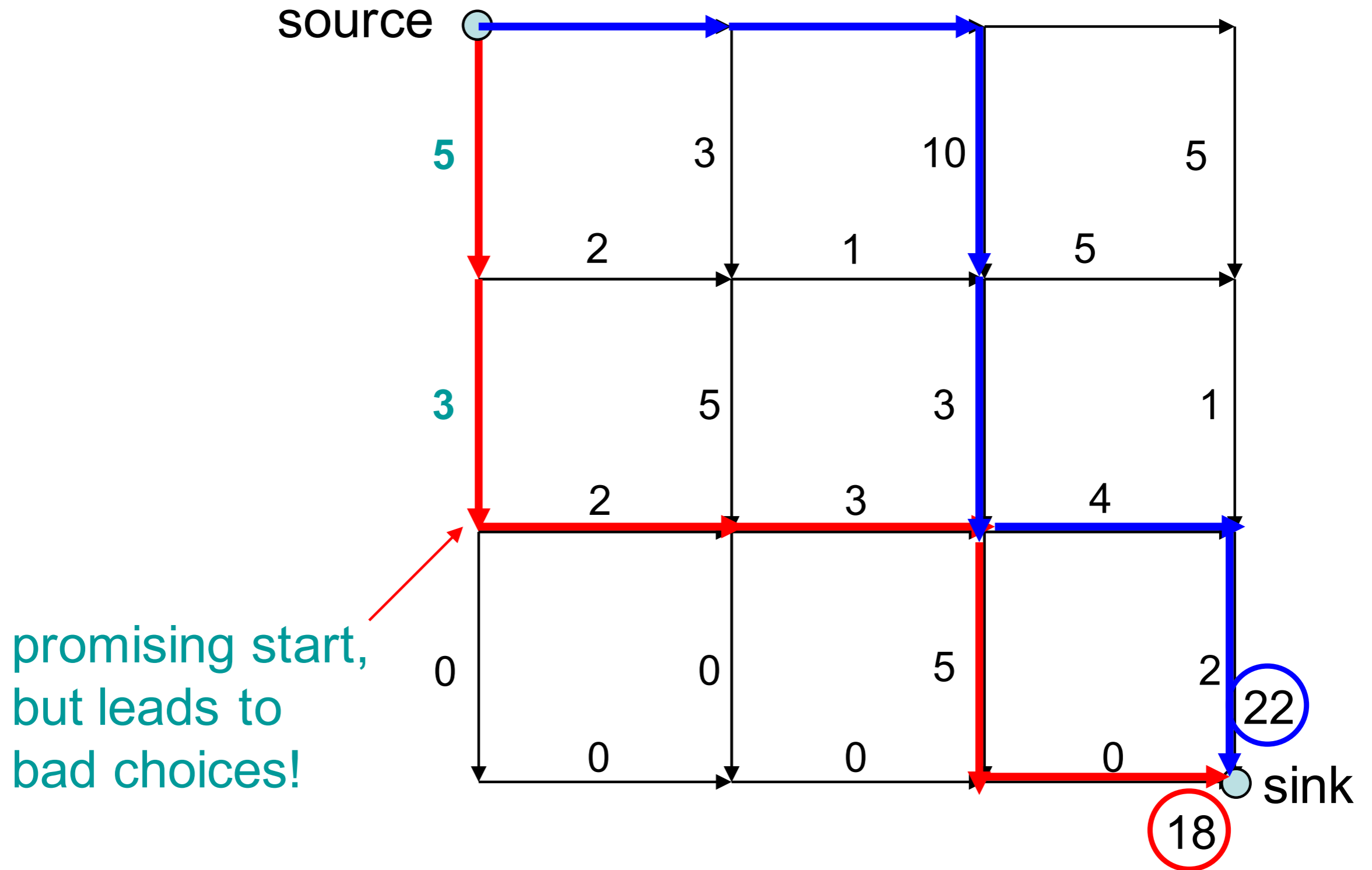
D

E

F

**A-GCDEF**

**AFGCDE-**

Very expensive

# Manhattan tourist problem

Imagine seeking a path (from source to sink) to travel (only eastward and southward) with the most number of attractions (*) in the Manhattan grid

# Would a greedy algorithm work?



source

5

3

3

10

5

2

1

5

5

3

1

2

3

4

promising start,
but leads to
bad choices!

0

0

5

2

0

0

0

22

18

sink

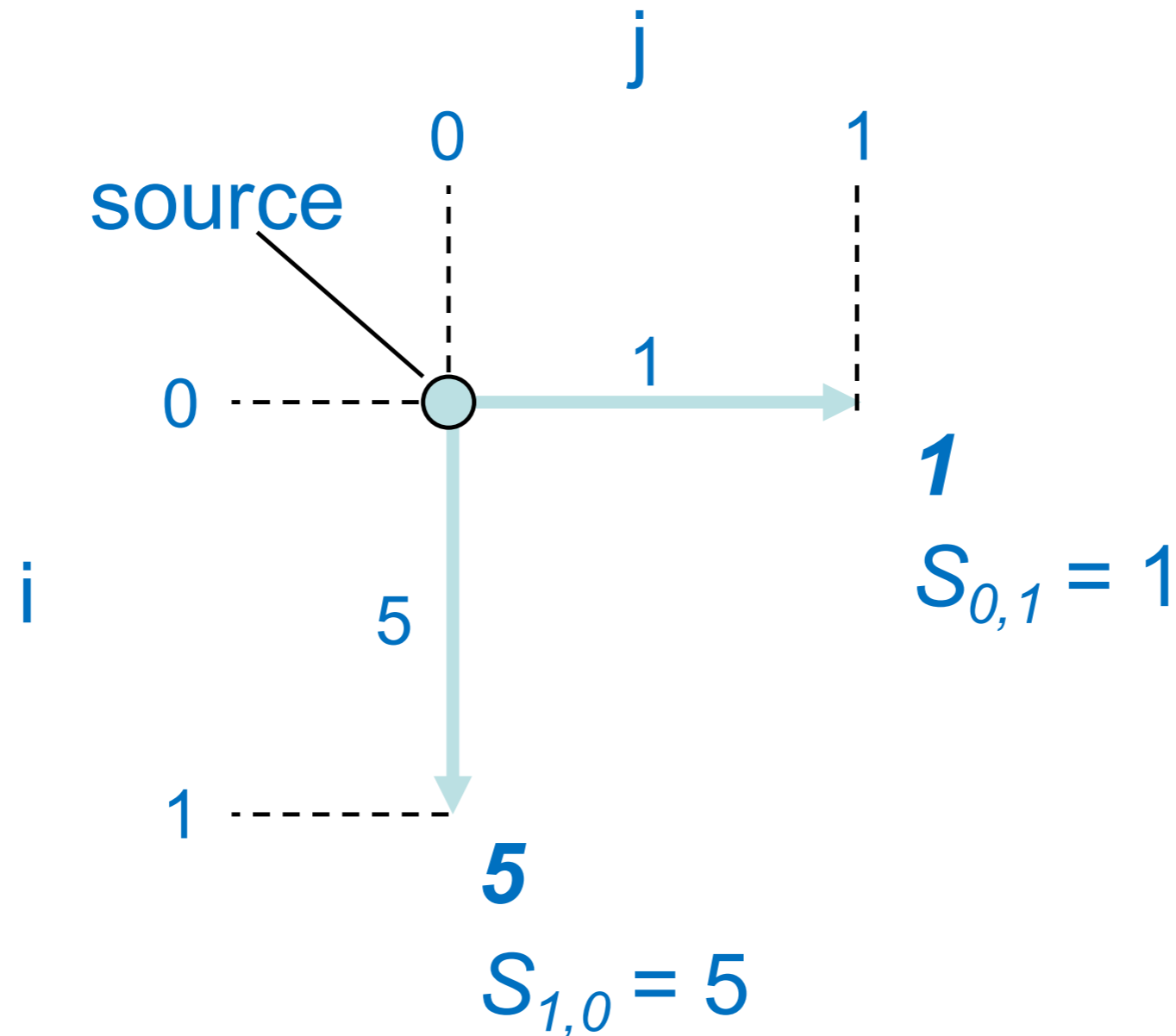# How about a recursive solution?

Function **MT**(*n,m*)

1. *x = MT(n-1,m)+*
   *weight of the edge from (n-1,m) to (n,m)*
2. *y = MT(n,m-1)+*
   *weight of the edge from (n,m-1) to (n,m)*
3. **return** *max{x,y}*

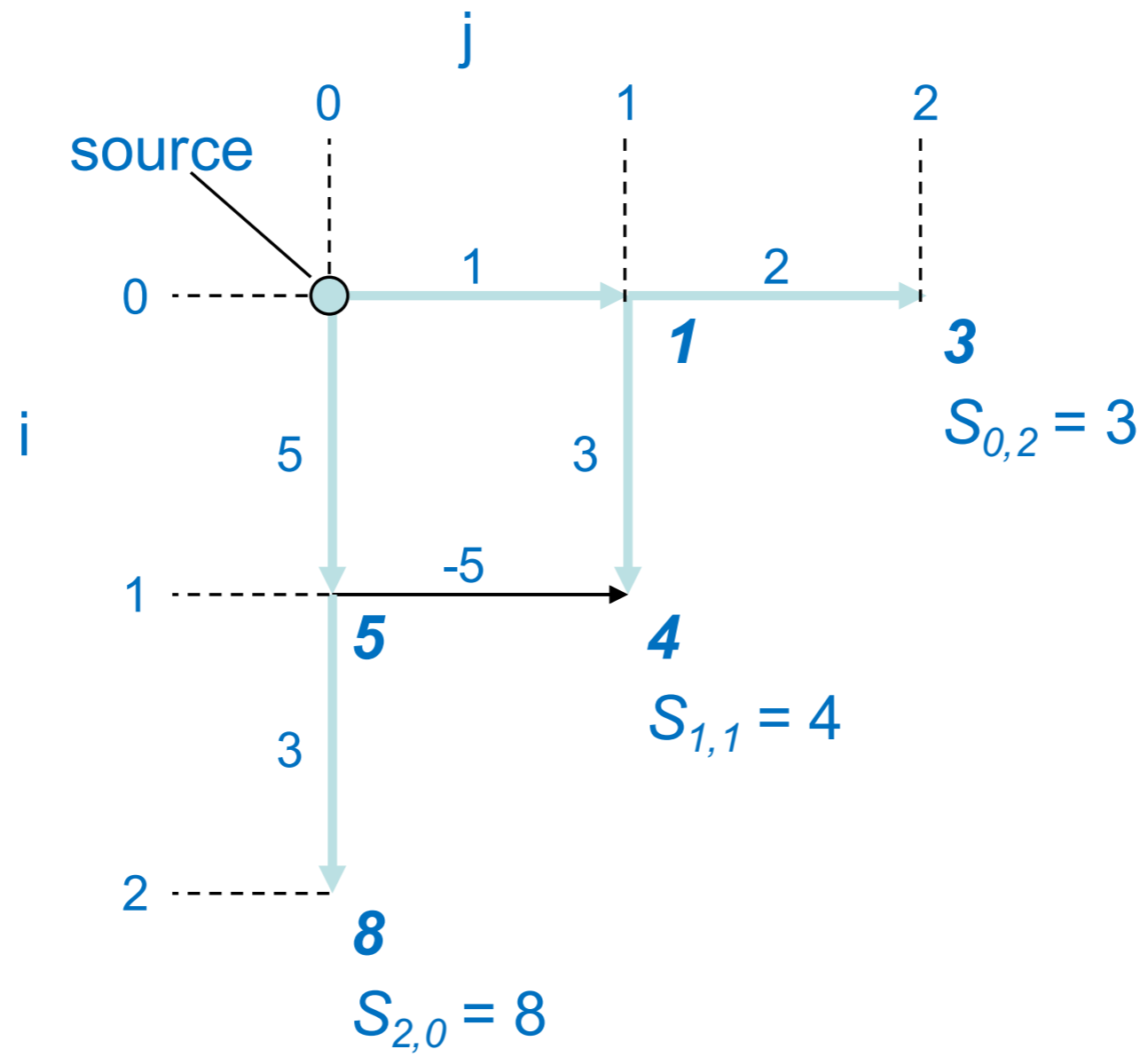**MT**(x, y) returns the "most weighted" path
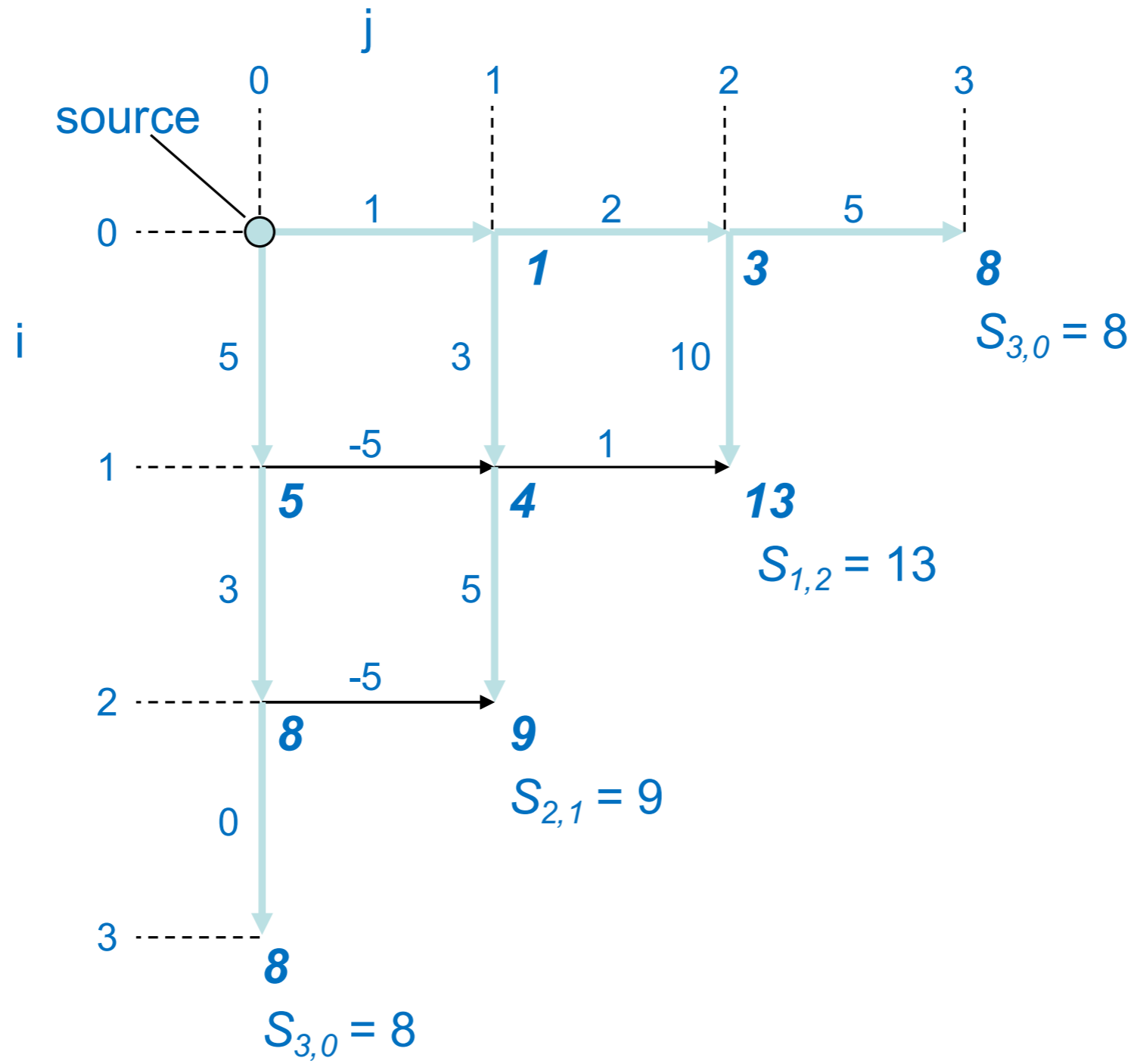from point (x, y) to the "sink".

# Why this is not efficient?

- MT(n,m) needs MT(n, m-1) and MT(n-1, m)
- Both of these need MT(n-1, m-1)
- So MT(n-1, m-1) will be computed at least twice
- Dynamic programming: the same idea as this recursive algorithm, but keep all intermediate results in a table and reuse
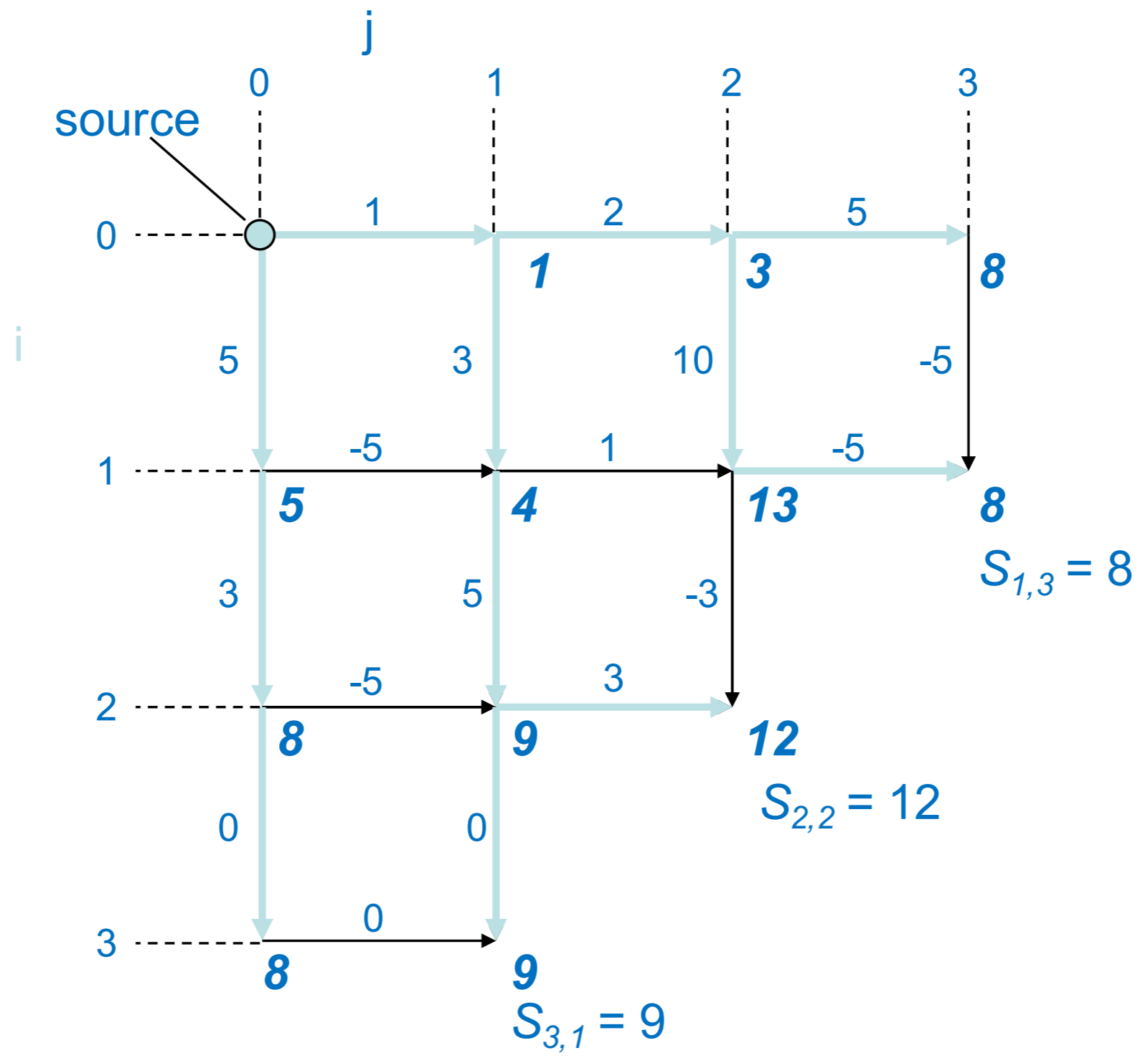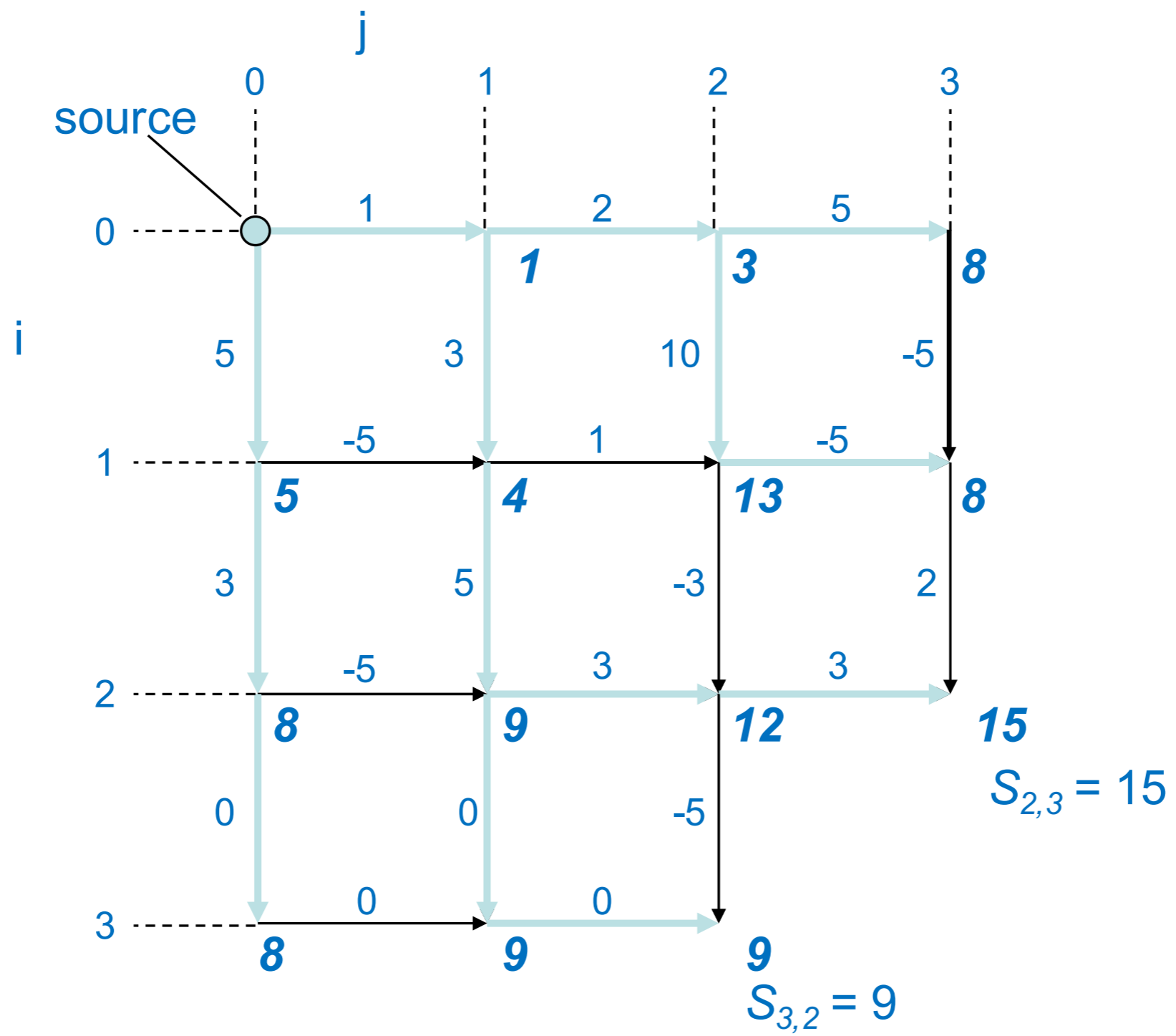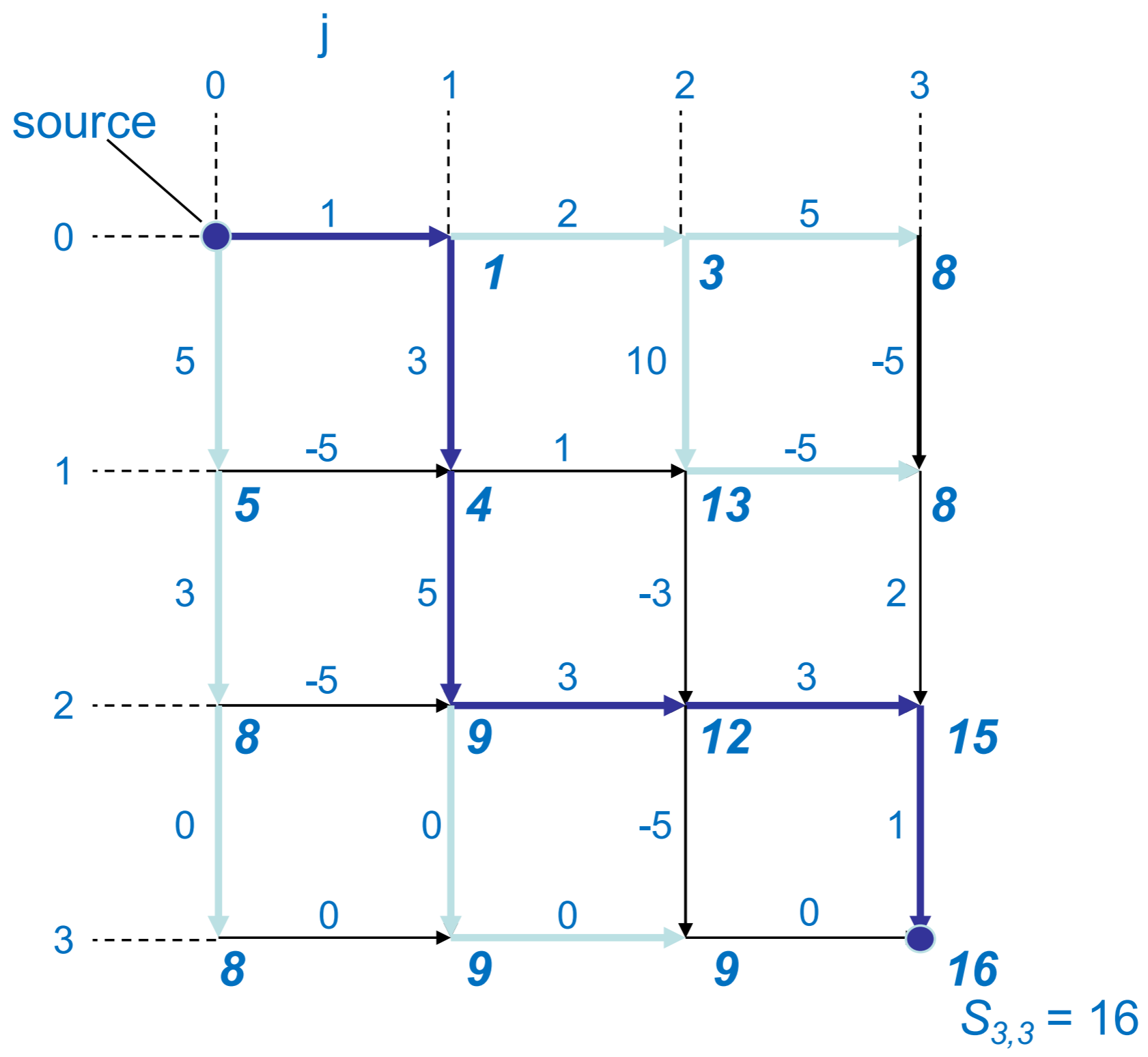
# How to avoid redundant calculations



- Calculate optimal path score for each vertex in the graph
- Each vertex's score is the maximum of the prior vertices score plus the weight of the respective edge in between
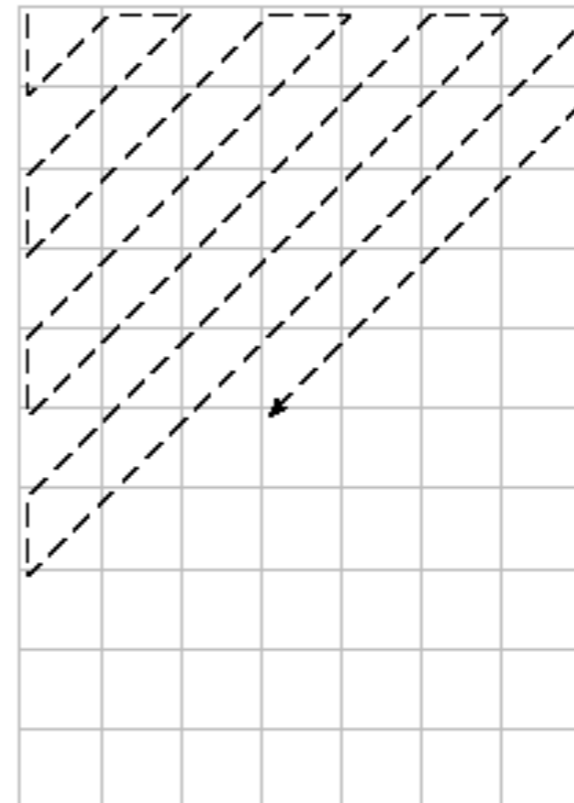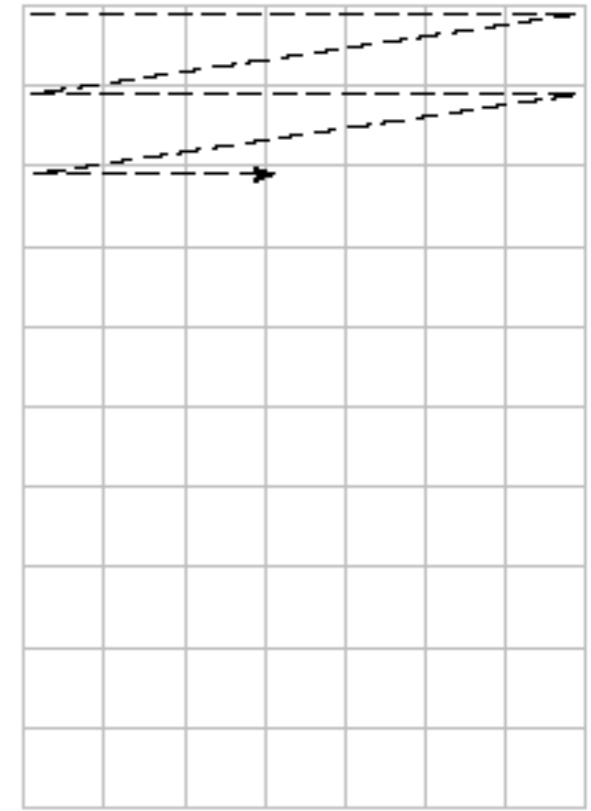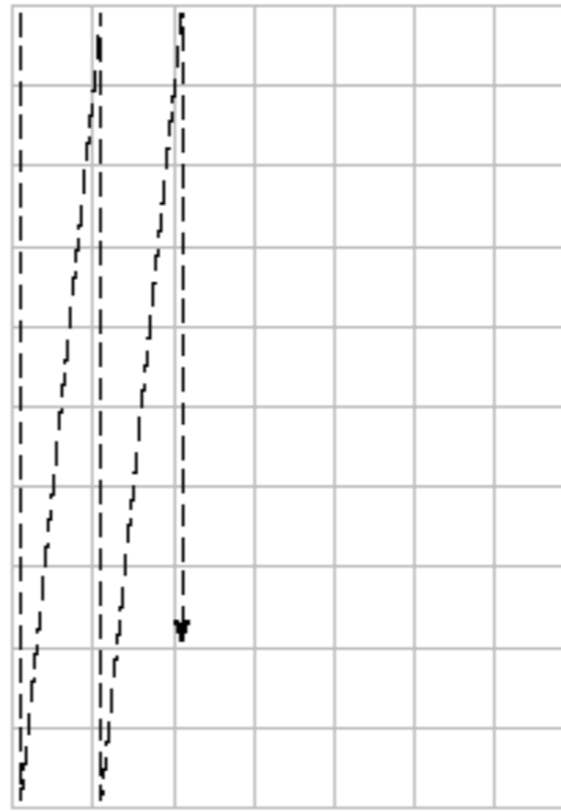
# To ensure the correctness

- By the time the vertex x is analyzed, the values $s_y$ for all its predecessors y should be computed – otherwise we are in trouble.

- We need to traverse the vertices in some order

- For a grid, can traverse vertices row by row, column by column, or diagonal by diagonal

3 different strategies:
  a) Column by column
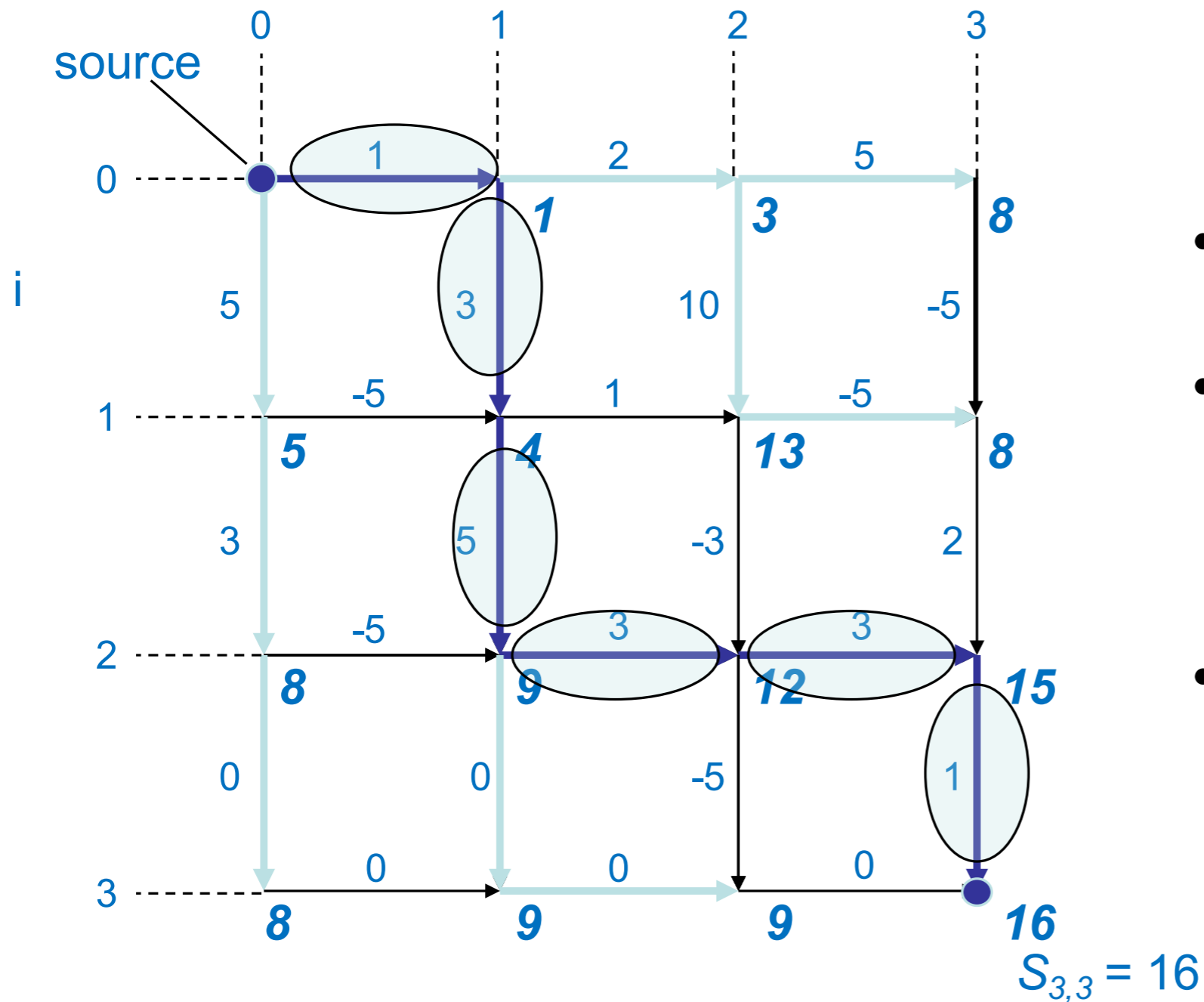  b) Row by row
  c) Along diagonals

Pseudocode?


Runtime?

# Recursive algorithm -> Dynamic programming ?

Function **MT**(*n,m*)

1. *x = MT(n-1,m)+*
   weight of the edge from *(n-1,m) to (n,m)*
2. *y = MT(n,m-1)+*
   weight of the edge from *(n,m-1) to (n,m)*
3. **return** *max{x,y}*

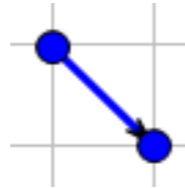**MT**(x, y) returns the "most weighted" path
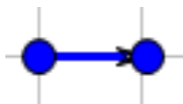from point (x, y) to the "sink".

# How to find the optimal path



- Start from Sink.

- Find which of the two edges gave the "max". Take it.

- Repeat.

$S_{3,3} = 16$

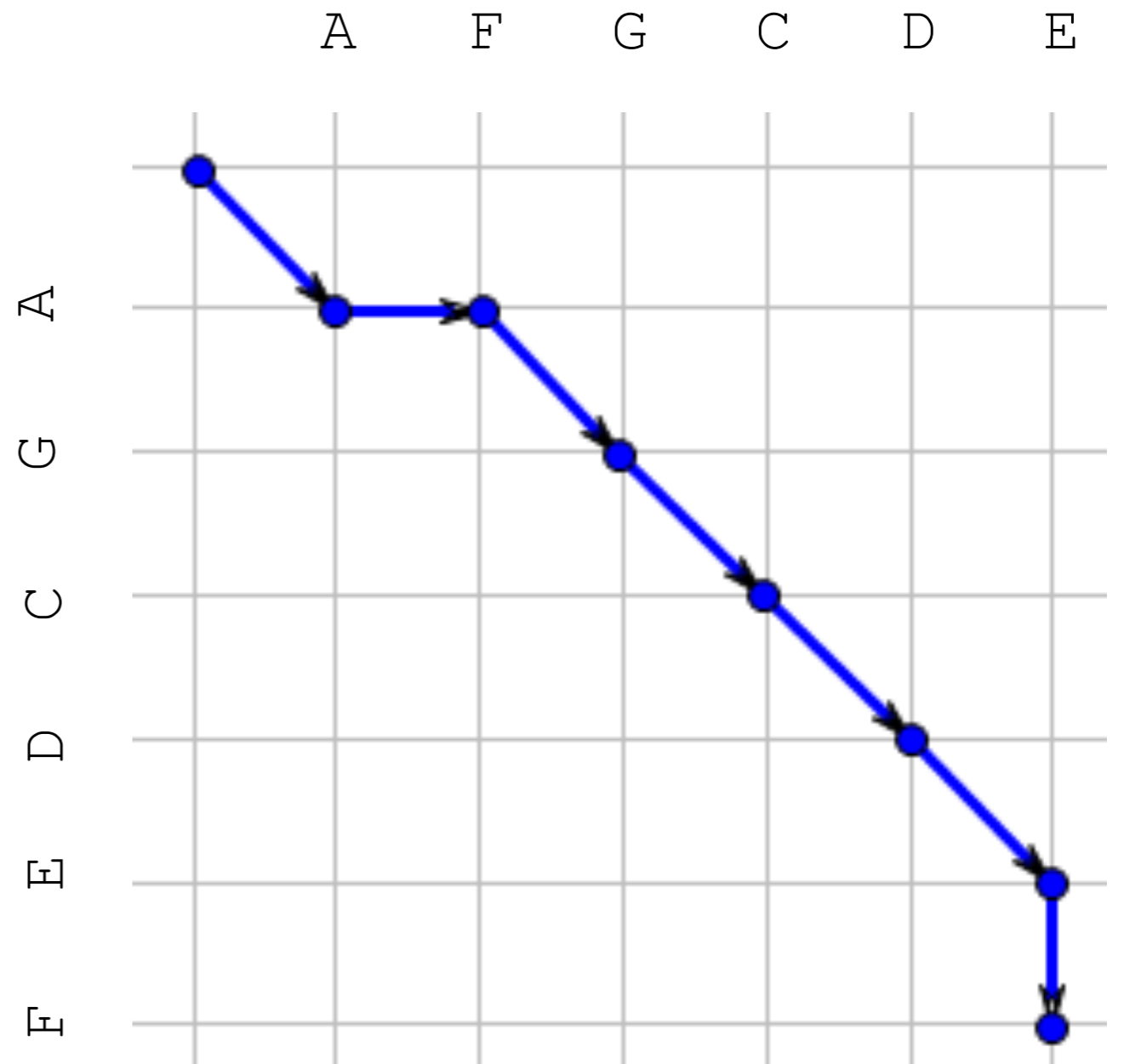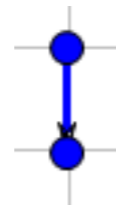# The edit distance problem



Match

Insertion_X

Insertion_Y

A F G C D E

A
G
C
D
E
F

**A-GCDEF**
**AFGCDE-**

# Recipe

1. Identify subproblems

2. Write down recursions

3. Make it dynamic-programming!