# Module 01 – Exercise Class

# Data Structure

**Nguyen Quoc Thai**

# Objectives

## Python Basic

- ❖ Variable
- ❖ Operators
- ❖ Condition
- ❖ Function
- ❖ Built-in Function
- ❖ For/While Loop
- ❖ If-else

## Data Structure

- ❖ List
- ❖ Dictionary
- ❖ Tuple

## Algorithm

- ❖ Getting Max Over Kernel
- ❖ Character Counting
- ❖ Word Counting
- ❖ Levenshtein Distance

# Outline

SECTION 1

## Getting Max Over Kernel

SECTION 2

## Character Counting

SECTION 3

## Word Counting (File)

SECTION 4

## Levenshtein Distance

# Getting Max Over Kernel

! **Description**

**Problem 01:** Cho một list các số nguyên ***num_list*** và một sliding window (các bạn có thể tạm hiểu sliding window như là một list có kích thước nhỏ hơn ***num_list***) có kích thước size ***k*** di chuyển từ trái sang phải. Mỗi lần dịch chuyển 1 vị trí sang phải có thể nhìn thấy được ***k*** số trong ***num_list*** và tìm số lớn nhất trong ***k*** số này sau mỗi lần trượt. ***k*** phải lớn hơn hoặc bằng 1. Các bạn hãy viết chương trình Python giải quyết vấn đề trên.

**Example:**
- **Input:**
  num_list = [3, 4, 5, 1, −44, 5, 10, 12, 33, 1]
  k = 3
- **Output:** [5, 5, 5, 5, 10, 12, 33, 33]

# Getting Max Over Kernel

**!** **Solution**

Use *slicing* and *max()* function

k = 2

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

Max([3, 1]) = 3

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

Max([1, -2]) = 1

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

Max([-2, -1]) = -1

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

Max([-1, 5]) = 5

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

Max([5, 4]) = 5

# Getting Max Over Kernel

! **Get Kernels**

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

k = 2

| 3 | 1 |
|---|---|

| 1 | -2 |
|---|----|

| -2 | -1 |
|----|----|

| -1 | 5 |
|----|---|

| 5 | 4 |
|---|---|

```python
1 num_list = [3 , 1, -2, -1, 5, 4]
2 k = 2
3 sub_list = []
4
5 for element in num_list:
6     sub_list.append(element)
7
8     if len(sub_list) == k:
9         print(sub_list)
10        del sub_list[0]
```

```
[3, 1]
[1, -2]
[-2, -1]
[-1, 5]
[5, 4]
```

7

# Getting Max Over Kernel

## ! Get Kernels – Solution #1

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

| 3 | 1 |
|---|---|

| 1 | -2 |
|---|----|

| -2 | -1 |
|----|----|

| -1 | 5 |
|----|---|

| 5 | 4 |
|---|---|

k = 2

```
1 num_list = [3 , 1, -2, -1, 5, 4]
2 k = 2
3 result = []
4 sub_list = []
5
6 for element in num_list:
7     sub_list.append(element)
8
9     if len(sub_list) == k:
10        result.append(max(sub_list))
11        del sub_list[0]
12
13 print(result)
```

[3, 1, -1, 5, 5]

8

# Getting Max Over Kernel

**!** **Slicing – Solution #2**

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

k = 2

list[start:end]

| list[0:2] | list[1:3] | list[2:4] | list[3:5] | list[4:6] |

| 3 | 1 |
|---|---|

| 1 | -2 |
|---|----|

| -2 | -1 |
|----|----|

| -1 | 5 |
|----|---|

| 5 | 4 |
|---|---|

Start Index

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

0 => len(list) - k

End Index

| 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|

k => len(list)

9

**AI VIET NAM**
@aivietnam.edu.vn

! **Slicing – Solution #2**

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

k = 2

Start Index

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

0 => len(list) - k

End Index

| 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|

k => len(list)

```python
1 num_list = [3 , 1 , -2, -1, 5, 4]
2 k = 2
3 start_indexs = list(range(0, len(num_list)-k+1))
4 end_indexs = list(range(k, len(num_list)+1))
5 print(start_indexs)
6 print(end_indexs)
```

```
[0, 1, 2, 3, 4]
[2, 3, 4, 5, 6]
```

10

# Getting Max Over Kernel

**Slicing – Solution #2**

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|---|---|---|---|

k = 2

**Start Index**

| 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|

**End Index**

| 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|

ZIP

| 3 | 1 |
|---|---|

Max([3, 1]) = 3

| 1 | -2 |
|---|---|

Max([1, -2]) = 1

| -2 | -1 |
|---|---|

Max([-2, -1]) = -1

| -1 | 5 |
|---|---|

Max([-1, 5]) = 5

| 5 | 4 |
|---|---|

Max([5, 4]) = 5

# Getting Max Over Kernel

! **Slicing**

| 3 | 1 | -2 | -1 | 5 | 4 |
|---|---|----|----|---|---|

k = 2

```python
1 num_list = [3 , 1 , -2, -1, 5, 4]
2 k = 2
3 start_indexs = list(range(0, len(num_list)-k+1))
4 end_indexs = list(range(k, len(num_list)+1))
5
6 result = []
7 for start_index, end_index in zip(start_indexs, end_indexs):
8     sub_list = num_list[start_index: end_index]
9     result.append(max(sub_list))
10
11 print(result)
```

[3, 1, -1, 5, 5]

12

# Outline

# Character Counting

**!** **Description**

**Problem 01:** Viết thuật toán trả về một dictionary đếm số lượng chữ xuất hiện trong một từ, với key là chữ cái và value là số lần xuất hiện.

       **Input:** một từ

       **Output:** dictionary đếm số lần các chữ xuất hiện

       **Note:** Giả sử các từ nhập vào đều có các chữ cái thuộc [a-z] hoặc [A-Z]

**Example:**
- **Input:**

       **word** = 'baby'
- **Output:**

       `{'b': 2, 'a': 1, 'y': 1}`

# Character Counting

**Solution**

name = 'baby'

index     0    1    2    3

name = | b | a | b | y |

```
1 word = 'baby'
2
3 for character in word:
4     print(character)
```
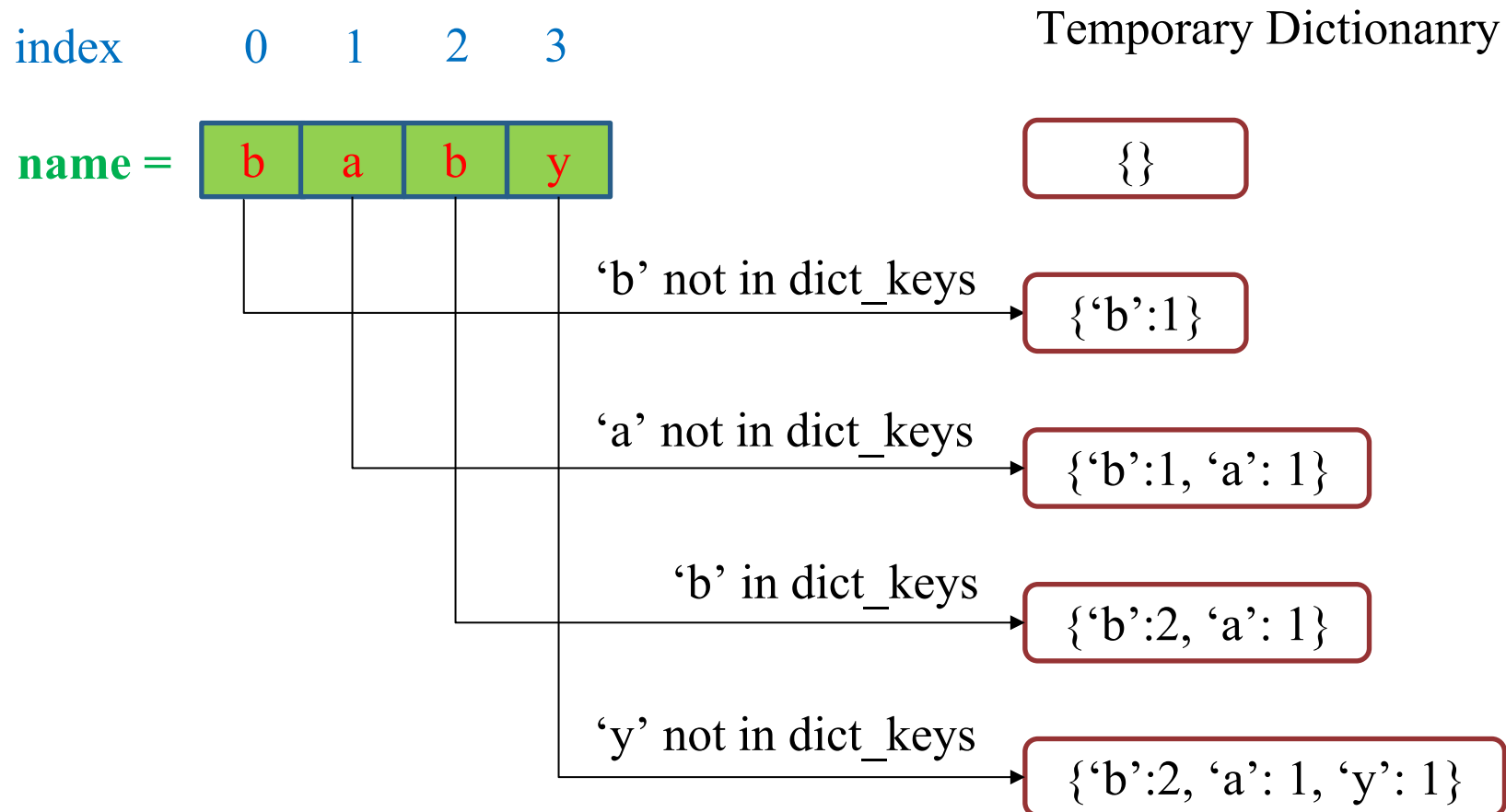
b
a
b
y

# Character Counting

**Solution**

name = 'baby'

index     0    1    2    3         Temporary Dictionanry

name = | b | a | b | y |

{}

'b' not in dict_keys        {'b':1}

'a' not in dict_keys        {'b':1, 'a': 1}

'b' in dict_keys        {'b':2, 'a': 1}

'y' not in dict_keys        {'b':2, 'a': 1, 'y': 1}

16

# Character Counting

**!** **Solution**

```python
1  character_statistic = {}
2
3  word = 'baby'
4
5  for character in word:
6      if character in character_statistic:
7          character_statistic[character] += 1
8      else:
9          character_statistic[character] = 1
10
11 print(character_statistic)
```

```
{'b': 2, 'a': 1, 'y': 1}
```

# Character Counting

**! Extension**

‘Baby’      {'B': 1, 'a': 1, 'b': 1, 'y': 1}

≠

‘baby’      {'b': 2, 'a': 1, 'y': 1}

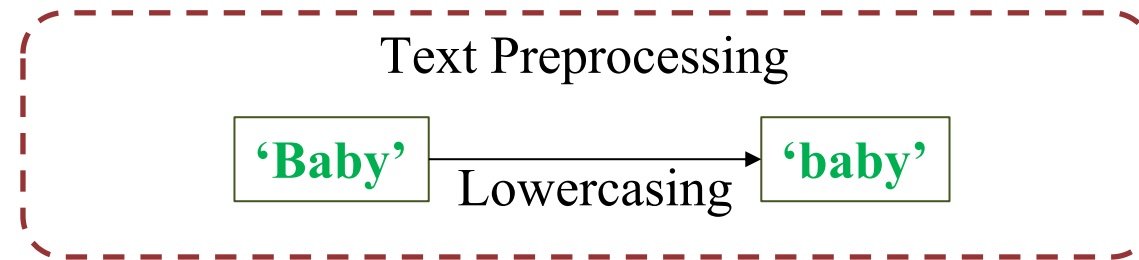‘Baby’ and ‘baby’: the same meaning in text

Text Preprocessing

‘Baby’ ⟶ ‘baby’

18

**! Extension**

Text Preprocessing

| **'Baby'** | → Lowercasing → | **'baby'** |

```
1 word = 'Baby'
2 print(word)
3 word = word.lower()
4 print(word)
```

```
Baby
baby
```

# Character Counting

! **Extension**

```python
1 character_statistic = {}
2
3 word = 'Baby'
4 word = word.lower()
5
6 for character in word:
7     if character in character_statistic:
8         character_statistic[character] += 1
9     else:
10         character_statistic[character] = 1
11
12 print(character_statistic)
```

{'b': 2, 'a': 1, 'y': 1}

20

# Outline

**SECTION 1**

## Getting Max Over Kernel

**SECTION 2**

## Character Counting

**SECTION 3**

## Word Counting (File)

**SECTION 4**

## Levenshtein Distance

# Word Counting

**Description**

**Problem:** Viết thuật toán đọc các câu trong một file txt, đếm số lượng các từ xuất hiện và trả về một dictionary với key là từ và value là số lần từ đó xuất hiện.

  **Input:** Đường dẫn đến file txt

  **Output:** Dictionary đếm số lần các từ xuất hiện

  **Note:**

  Giả sử các từ trong file txt đều có các chữ cái thuộc [a-z] hoặc [A-Z]

  Không cần các thao tác xử lý string phức tạp nhưng cần xử lý các từ đều là viết thường

  File: https://drive.google.com/uc?id=1IBScGdW2xlNsc9v5zSAya548kNgiOrko

# Word Counting

**! Example**

```
2  !gdown   https://drive.google.com/uc?id=1IBScGdW2xlNsc9v5zSAya548kNgiOrko
3  file_path = '/content/P1_data.txt'
4  word_count(file_path)
5  >>{'a': 7,
6   'again': 1,
7   'and': 1,
8   'are': 1,
9   'at': 1,
10  'be': 1,
11  'become': 2,
12  ...}
```

## ! Read File

```
1 !gdown https://drive.google.com/uc?id=1IBScGdW2xlNsc9v5zSAya548kNgiOrko
```

```
Downloading...
From: https://drive.google.com/uc?id=1IBScGdW2xlNsc9v5zSAya548kNgiOrko
To: /content/P1_data.txt
100% 747/747 [00:00<00:00, 1.81MB/s]
```

```python
1 with open('/content/P1_data.txt', 'r') as f:
2     sentences = f.readlines()
3 type(sentences)
```

```
list
```

```python
1 sentences[:2]
```

```
['He who conquers himself is the mightiest warrior\n',
 'Try not to become a man of success but rather become a man of value\n']
```

24

**AI VIET NAM**
@aivietnam.edu.vn

## ! Read File

```python
1 with open('/content/P1_data.txt', 'r') as f:
2     document = f.read()
3 type(document)
```

str

```python
1 document
```
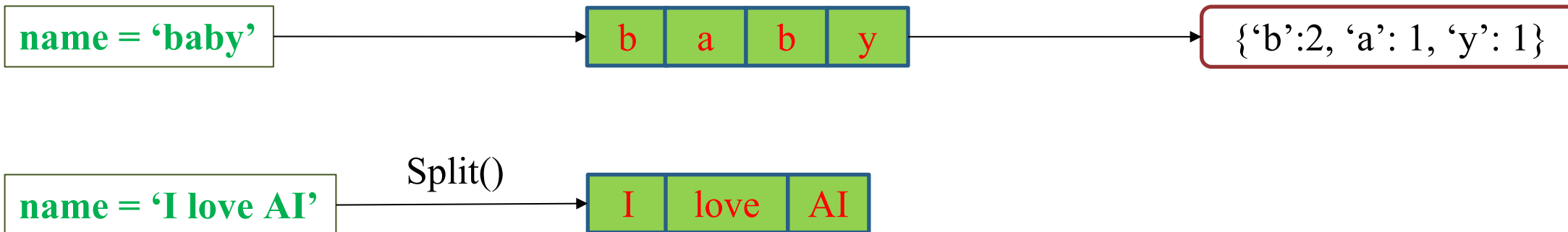
'He who conquers himself is the mightiest warrior\nTry not to become a man of success but rather become a man of value\nOne man with courage makes a majority\nOne secret of success in life is for a man to be ready for his opportunity when it comes\nThe successful man will profit from his mistakes and try again in a different way \nA successful man is one who can lay a firm foundation with the bricks others have thrown at him\nSuccess usually comes to those who are too busy looking for it\nWe cannot solve problems with the kind of thinking we employed when we came up with them\nJust one small positive thought in the morning can change your whole day\nYou can get everything in life you want if you will just help enough other people get what they want'
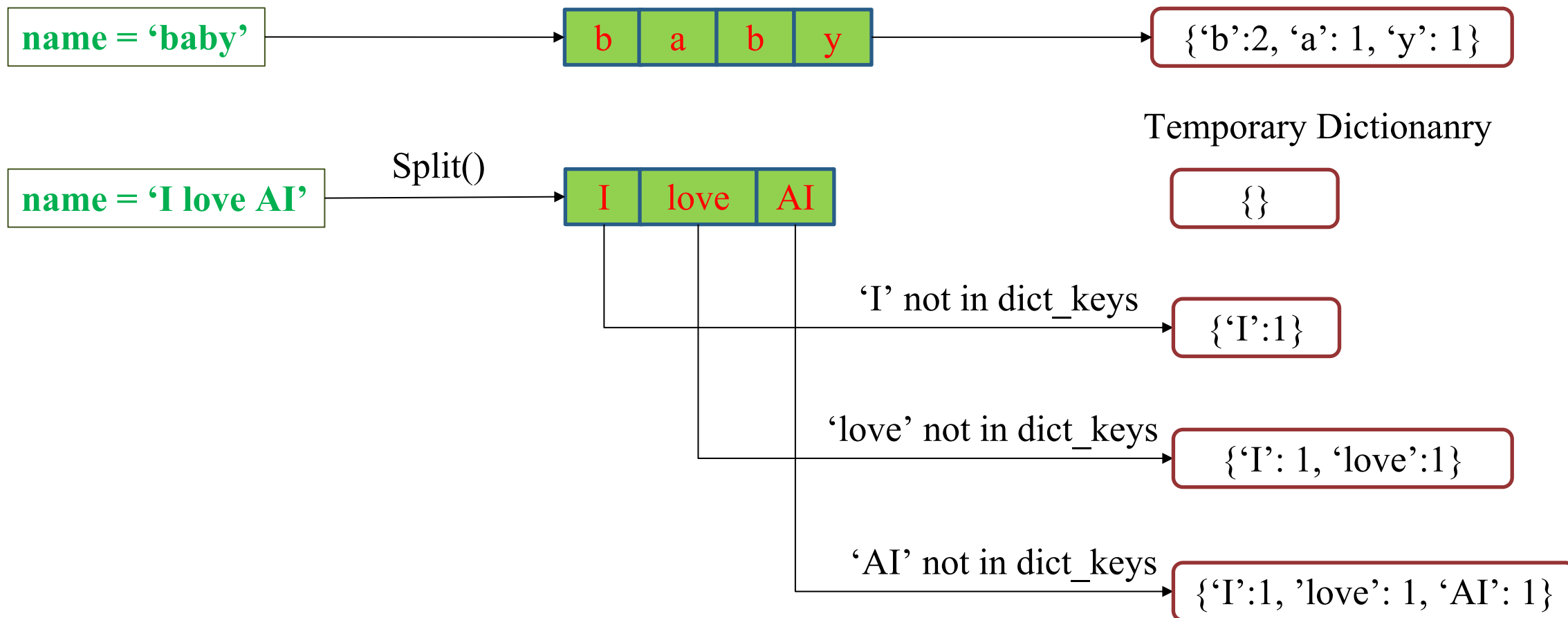
# Word Counting

**!** **Counting**

name = 'baby' → | b | a | b | y | → {'b':2, 'a': 1, 'y': 1}

name = 'I love AI' — Split() → | I | love | AI |

# Word Counting

## ! Counting

name = 'baby' → | b | a | b | y | → {'b':2, 'a': 1, 'y': 1}

Temporary Dictionanry

name = 'I love AI' —Split()→ | I | love | AI | → {}

'I' not in dict_keys → {'I':1}

'love' not in dict_keys → {'I': 1, 'love':1}

'AI' not in dict_keys → {'I':1, 'love': 1, 'AI': 1}

# Word Counting

**!** **Counting**

```python
1 sentence = 'I love AI'
2 words = sentence.split()
3
4 counter = {}
5 for word in words:
6     if word in counter:
7         counter[word] += 1
8     else:
9         counter[word] = 1
10
11 print(counter)
```

```
{'I': 1, 'love': 1, 'AI': 1}
```

# Word Counting

## ! Word Counting

```python
1 words = document.split()
2
3 counter = {}
4 for word in words:
5     if word in counter:
6         counter[word] += 1
7     else:
8         counter[word] = 1
9
10 print(counter)
```

```
{'He': 1, 'who': 3, 'conquers': 1, 'himself': 1, 'is': 3, 'the': 4, 'mightiest': 1
```

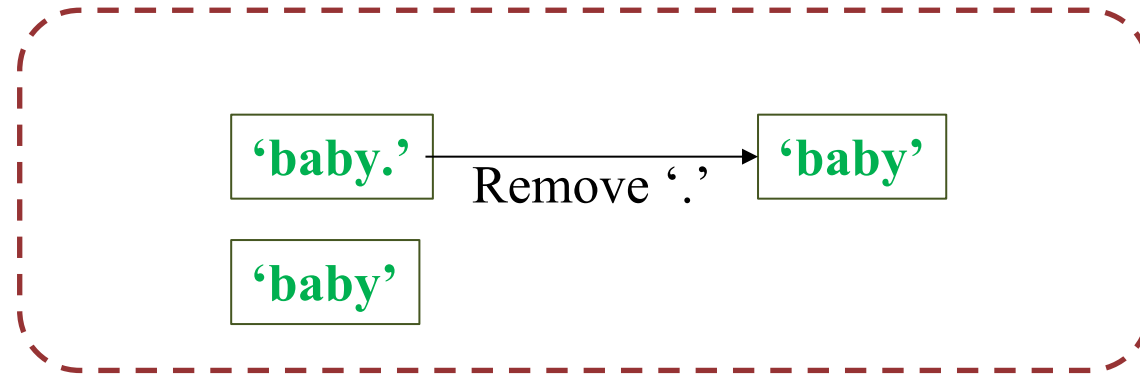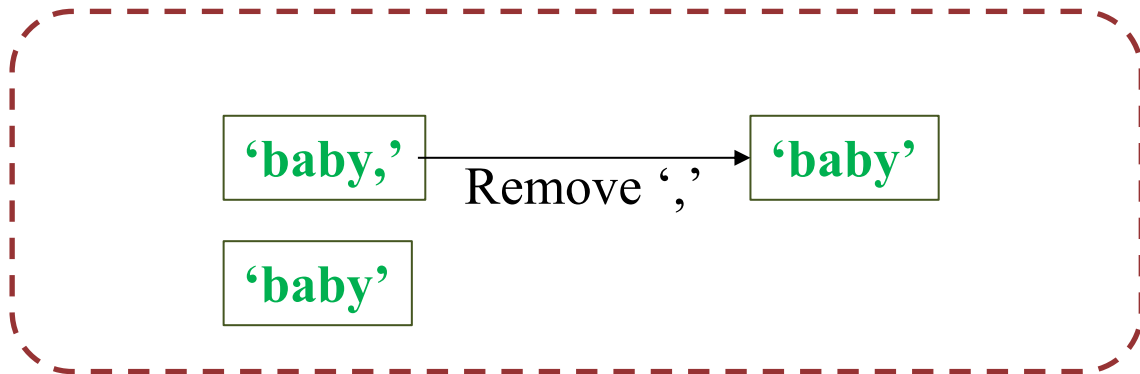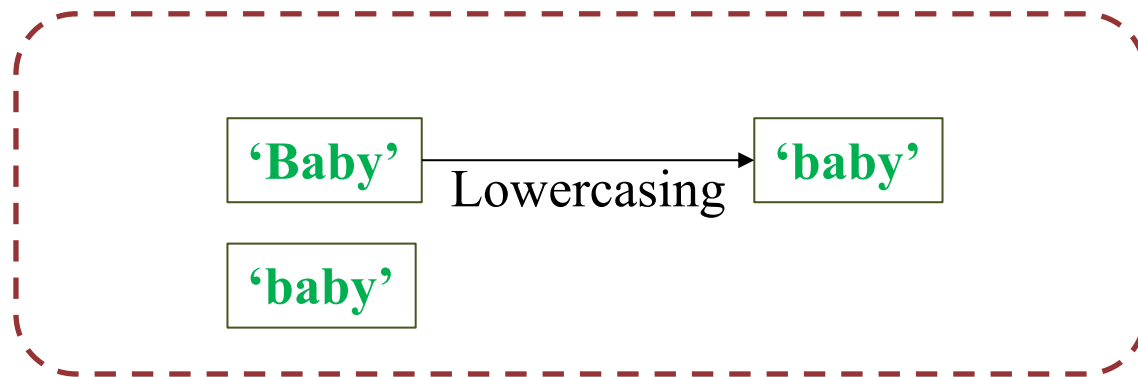! **Word Counting**

```python
1 counter = {}
2
3 for sentence in sentences:
4     words = sentence.split()
5     for word in words:
6         if word in counter:
7             counter[word] += 1
8         else:
9             counter[word] = 1
10
11 print(counter)
```

{'He': 1, 'who': 3, 'conquers': 1, 'himself': 1, 'is': 3, 'the': 4, 'mightiest': 1

**AI VIET NAM**
@aivietnam.edu.vn

**!** **Text Preprocessing**

'Baby' → Lowercasing → 'baby'

'baby'

'baby,' → Remove ',' → 'baby'

'baby'

'baby.' → Remove '.' → 'baby'

'baby'

**AI VIET NAM**
@aivietnam.edu.vn

**!** **Text Preprocessing**

```python
1 def preprocess_text(sentence):
2     sentence = sentence.lower()
3     sentence = sentence.replace('.', '').replace(',', '')
4     words = sentence.split()
5     return words
6
7 sentence = 'I love AI. AI is not easy'
8 preprocess_text(sentence)
```

```
['i', 'love', 'ai', 'ai', 'is', 'not', 'easy']
```

! **Text Preprocessing**

```python
 1 words = preprocess_text(document)
 2
 3 counter = {}
 4 for word in words:
 5     if word in counter:
 6         counter[word] += 1
 7     else:
 8         counter[word] = 1
 9
10 print(counter)
```

```
{'he': 1, 'who': 3, 'conquers': 1, 'himself': 1, 'is': 3, 'the': 5, 'mightiest': 1
```

## Text Preprocessing

```python
1  counter = {}
2
3  for sentence in sentences:
4      words = preprocess_text(sentence)
5      for word in words:
6          if word in counter:
7              counter[word] += 1
8          else:
9              counter[word] = 1
10
11 print(counter)
```

{'he': 1, 'who': 3, 'conquers': 1, 'himself': 1, 'is': 3, 'the': 5, 'mightiest': 1
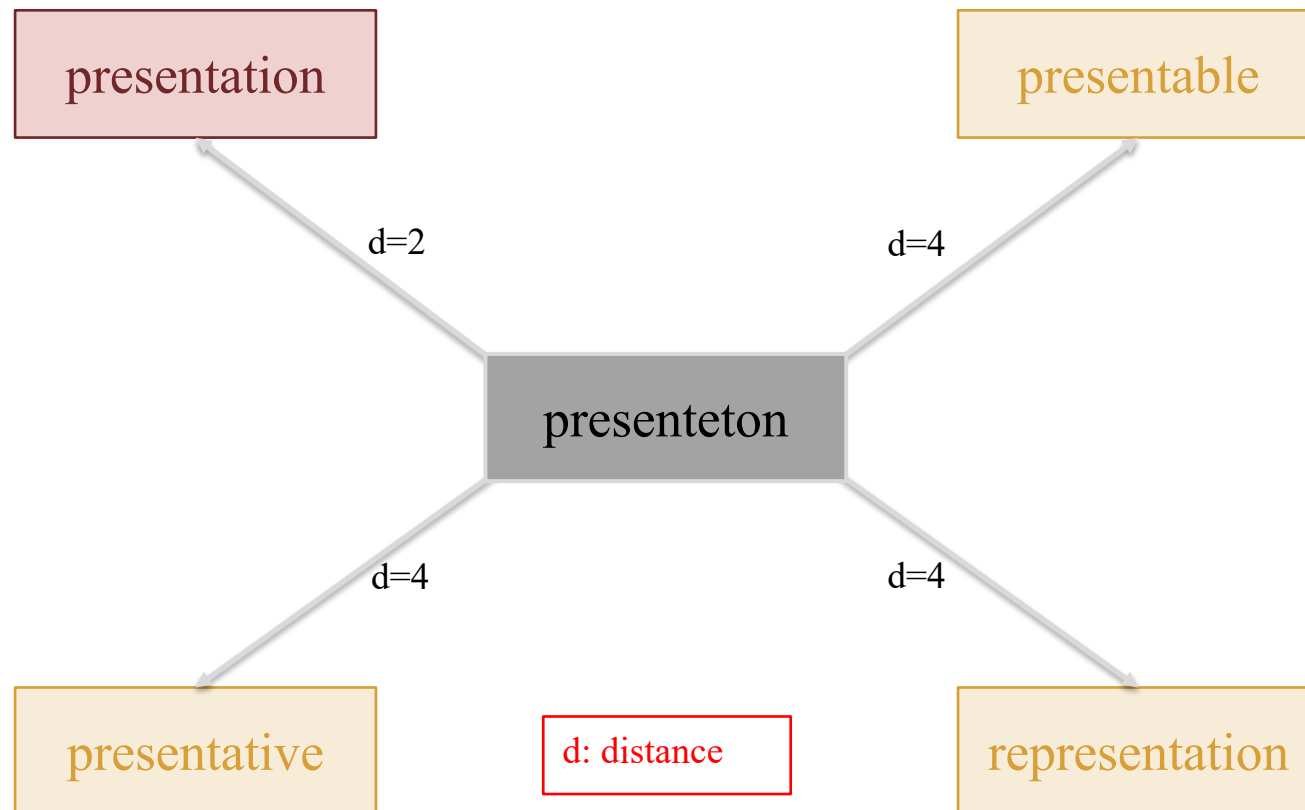
# Outline

# Levenshtein Distance

**AI VIET NAM**

! **Getting Started**

presenteton

✦ Có phải ý của bạn là: *presentation*

🎤  🔊

How to measure the similarity or gap between two strings ?

presentation

presentable

d=2

d=4

presenteton

d=4

d=4

presentative

d: distance

representation

# Levenshtein Distance

**AI VIET NAM**

! **The Minimum Edit Distance Algorithm**

| Source: | hola |
|---|---|

| Target: | hello |
|---|---|

Delete cost = 1

Insert cost = 1

Substitution cost = 1

insert →

|  |  | Target | | | | | |
|---|---|---|---|---|---|---|---|
|  |  | **#** | **h** | **e** | **l** | **l** | **o** |
| Source | **#** | 0 | 1 | 2 | ③ | 4 | 5 |
|  | **h** | ① |  |  |  |  |  |
|  | **o** | 2 |  |  | sub |  |  |
|  | **l** | ③ |  |  |  |  |  |
|  | **a** | 4 |  |  |  |  |  |

delete ↓

# → #hel
1. insert h
2. insert e
3. insert l

#h → #
1. delete h

#hol → #
1. delete h
2. delete o
3. delete l

38

# Levenshtein Distance

**Algorithm**

$$D[i,j] = min \begin{cases} D[i-1,j] + \text{delcost(source[i])} \\ D[i,j-1] + \text{inscost(target[j])} \\ D[i-1,j-1] + \text{subcost(source[i], target[j])} \end{cases}$$

insert

| j \ i | | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|
| | | # | h | e | l | l | o |
| # | | 0 | 1 | 2 | 3 | 4 | 5 |
| h | | 1 | | | | | |
| o | | 2 | | | | | |
| l | | 3 | | | | | |
| a | | 4 | | | | | |

**Target**

**Source**

delete

Delete cost = 1

Insert cost = 1

Substitution cost = 1

$$D[1,1] = min \begin{cases} D[0,1] + \text{delcost(source[1])} = 1 + 1 = 2 \\ D[1,0] + \text{inscost(target[1])} = 1 + 1 = 2 \\ D[0,0] + \text{subcost}(source[1], \text{target[1]}) = 0 + 0 = 0 \end{cases}$$

Note:
If source[i] == target[j]:
    subcost(source[i] ,target[j])=0

# Levenshtein Distance

! **Algorithm**

insert →

| | | | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|---|
| | | j | | | Target | | | |
| | i | | # | h | e | l | l | o |
| | | # | 0 | 1 | 2 | 3 | 4 | 5 |
| Source | | h | 1 | 0 | | | | |
| | | o | 2 | | | | | |
| | | l | 3 | | | | | |
| | | a | 4 | | | | | |

delete

Delete cost = 1

Insert cost = 1

Substitution cost = 1

$$D[2,1] = min \begin{cases} D[1,1] + \text{delcost}(source[2]) = 0 + 1 = 1 \\ D[2,0] + \text{inscost}(target[1]) = 2 + 1 = 3 \\ D[1,0] + \text{subcost}(source[2], target[1]) = 1 + 1 = 2 \end{cases}$$

Note:
If source[i] == target[j]:
    subcost(source[i] ,target[j])=0

$$D[1,2] = min \begin{cases} D[0,1] + \text{delcost}(source[1]) = 2 + 1 = 3 \\ D[1,1] + \text{inscost}(target[2]) = 0 + 1 = 1 \\ D[0,1] + \text{subcost}(source[1], target[2]) = 1 + 1 = 2 \end{cases}$$

# Levenshtein Distance

## ! Algorithm



insert →

| i \ j | | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|
| | | | | **Target** | | | |
| | | # | h | e | l | l | o |
| | # | 0 | 1 | 2 | 3 | 4 | 5 |
| **Source** | h | 1 | 0 | 1 | 2 | 3 | 4 |
| | o | 2 | 1 | 1 | 2 | 3 | 3 |
| | l | 3 | 2 | 2 | 1 | 2 | 3 |
| | a | 4 | 3 | 3 | 2 | 2 | 3 |

delete
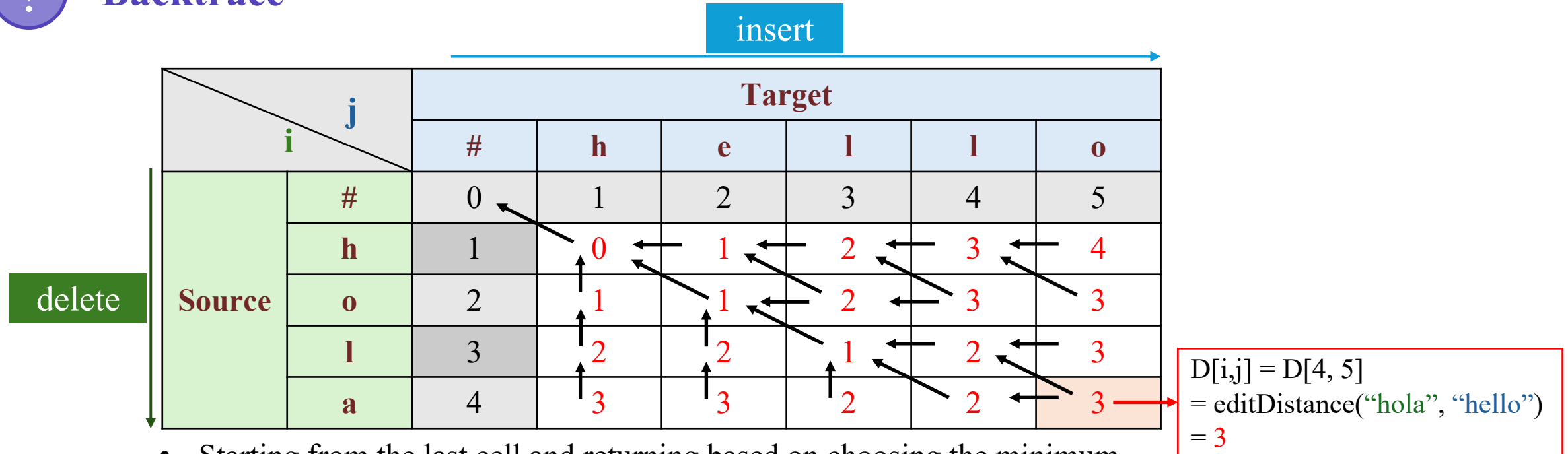
D[i,j] = D[4, 5]
= editDistance("hola","hello")
= 3

When going down each step, store back pointers in each cell to serve for the backtrace phase.

# Levenshtein Distance

**Backtrace**

insert →

| | j | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|
| i | | | | Target | | | |
| # | # | 0 | 1 | 2 | 3 | 4 | 5 |
| Source | h | 1 | 0 | 1 | 2 | 3 | 4 |
| | o | 2 | 1 | 1 | 2 | 3 | 3 |
| | l | 3 | 2 | 2 | 1 | 2 | 3 |
| | a | 4 | 3 | 3 | 2 | 2 | 3 |

delete ↓

D[i,j] = D[4, 5]
= editDistance("hola", "hello")
= 3

- Starting from the last cell and returning based on choosing the <u>minimum cell value.</u>
- Each cell may have <u>multiple</u> to return to because they have the same <u>minimum value</u>.

# Levenshtein Distance

**! Minimum edit distance path**

insert

| j \ i | | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|
| | | **Target** | | | | | |
| | # | 0 | 1 | 2 | 3 | 4 | 5 |
| **Source** | h | 1 | 0 | 1 | 2 | 3 | 4 |
| | o | 2 | 1 | 1 | 2 | 3 | 3 |
| | l | 3 | 2 | 2 | 1 | 2 | 3 |
| | a | 4 | 3 | 3 | 2 | 2 | 3 |

delete

D[i,j] = D[4, 5]
= editDistance("hola", "hello")
= 3

- This is one of the minimum edit distance paths.
- Modify steps (going in the <u>reverse</u> direction with the backtrace path):
  - sub(h, h) => hola; cost = 0
  - sub(o, e) => hela; cost = 1
  - sub(l, l) => hela; cost = 0
  - sub(a, l) => hell; cost = 1
  - insert(o) => hello; cost = 1

43

# Levenshtein Distance

**Backtrace**



insert

| j / i | | # | h | e | l | l | o |
|---|---|---|---|---|---|---|---|
| | | | | **Target** | | | |
| | # | 0 | 1 | 2 | 3 | 4 | 5 |
| **Source** | h | 1 | 0 | 1 | 2 | 3 | 4 |
| | o | 2 | 1 | 1 | 2 | 3 | 3 |
| | l | 3 | 2 | 2 | 1 | 2 | 3 |
| | a | 4 | 3 | 3 | 2 | 2 | 3 |

delete

D[i,j] = D[4, 5]
= editDistance("hola", "hello")
= 3

- This is another minimum edit distance path.
- Modify steps (going in the reverse direction with the backtrace path):
  - sub(h, h) => hola; cost = 0
  - sub(o, e) => hela; cost = 1
  - sub(l, l) => hela; cost = 0
  - insert(l) => hella; cost = 1
  - sub(a, o) => hello; cost = 1
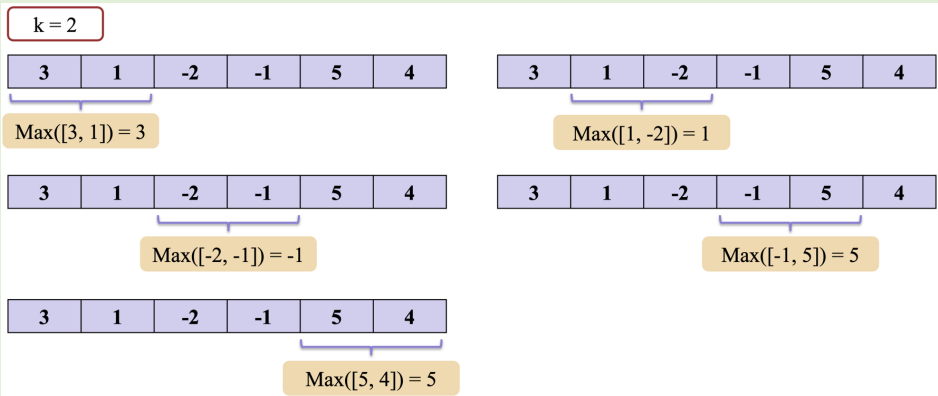
# Levenshtein Distance

**!** **Code**

```python
1 def levenshtein_distance(token1, token2):
2     distances = [[0]*(len(token2)+1) for i in range(len(token1)+1)]
3
4     for t1 in range(len(token1) + 1):
5         distances[t1][0] = t1
6
7     for t2 in range(len(token2) + 1):
8         distances[0][t2] = t2
9
10    a = 0
11    b = 0
12    c = 0
```

```python
13
14    for t1 in range(1, len(token1) + 1):
15        for t2 in range(1, len(token2) + 1):
16            if (token1[t1-1] == token2[t2-1]):
17                distances[t1][t2] = distances[t1 - 1][t2 - 1]
18            else:
19                a = distances[t1][t2 - 1]
20                b = distances[t1 - 1][t2]
21                c = distances[t1 - 1][t2 - 1]
22
23                if (a <= b and a <= c):
24                    distances[t1][t2] = a + 1
25                elif (b <= a and b <= c):
26                    distances[t1][t2] = b + 1
27                else:
28                    distances[t1][t2] = c + 1
29
30    return distances[len(token1)][len(token2)]
31
32 assert levenshtein_distance("hi", "hello") == 4
33 print(levenshtein_distance("hola", "hello"))
```
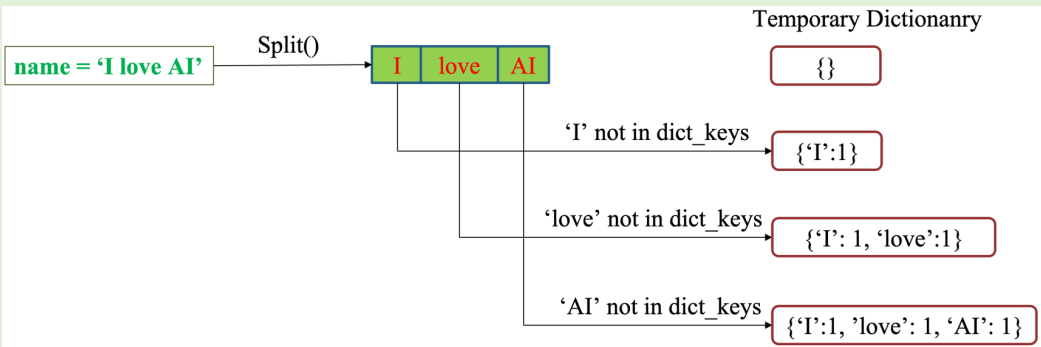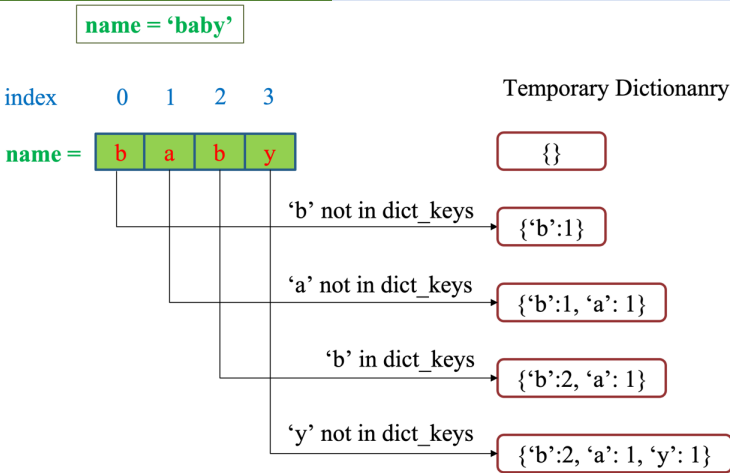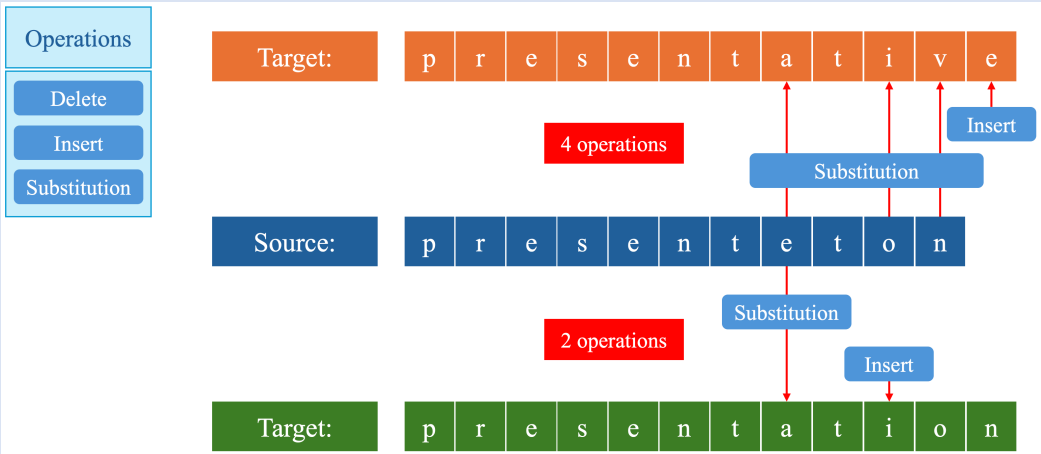
3

45

# Summary

## Getting Max Over Kernel

k = 2

| 3 | 1 | -2 | -1 | 5 | 4 |

Max([3, 1]) = 3

| 3 | 1 | -2 | -1 | 5 | 4 |

Max([1, -2]) = 1

| 3 | 1 | -2 | -1 | 5 | 4 |

Max([-2, -1]) = -1

| 3 | 1 | -2 | -1 | 5 | 4 |

Max([-1, 5]) = 5

| 3 | 1 | -2 | -1 | 5 | 4 |

Max([5, 4]) = 5

## Word Counting

name = 'I love AI'  Split() → | I | love | AI |

Temporary Dictionanry

{}

'I' not in dict_keys → {'I':1}

'love' not in dict_keys → {'I': 1, 'love':1}

'AI' not in dict_keys → {'I':1, 'love': 1, 'AI': 1}

## Character Counting

name = 'baby'

| index | 0 | 1 | 2 | 3 |

Temporary Dictionanry

name = | b | a | b | y |

{}

'b' not in dict_keys → {'b':1}

'a' not in dict_keys → {'b':1, 'a': 1}

'b' in dict_keys → {'b':2, 'a': 1}

'y' not in dict_keys → {'b':2, 'a': 1, 'y': 1}

## Levenshtein

Operations
- Delete
- Insert
- Substitution

Target: | p | r | e | s | e | n | t | a | t | i | v | e |

Insert

4 operations

Substitution

Source: | p | r | e | s | e | n | t | e | t | o | n |

Substitution

2 operations

Insert

Target: | p | r | e | s | e | n | t | a | t | i | o | n |

46

# Thanks!

## Any questions?